# SAMPLE EFFICIENT ROBUST OFFLINE SELF-PLAY FOR MODEL-BASED REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

## ABSTRACT

Multi-agent reinforcement learning (MARL), as a thriving field, explores how multiple agents independently make decisions in a shared dynamic environment. Due to environmental uncertainties, policies in MARL must remain robust to tackle the sim-to-real gap. Although robust RL has been extensively explored in single-agent settings, it has seldom received attention in self-play, where strategic interactions heighten uncertainties. We focus on robust two-player zero-sum Markov games (TZMGs) in offline RL, specifically on tabular robust TZMGs (RTZMGs) with a given uncertainty set. To address sample scarcity, we introduce a model-based algorithm (RTZ-VI-LCB) for RTZMGs, which integrates robust value iteration considering uncertainty level and applies a data-driven penalty to the robust value estimates. We establish the finite-sample complexity of RTZ-VI-LCB by accounting for distribution shifts in the historical dataset. Our algorithm is capable of learning under partial coverage and environmental uncertainty. An information-theoretic lower bound is developed to show that learning RTZMGs is at least as difficult as standard TZMGs when the uncertainty level is sufficiently small. This confirms the tightness of our algorithm's sample complexity, which is optimal regarding both state and action spaces. To the best of our knowledge, our algorithm is the first to attain this optimality and establishes a new benchmark for offline RTZMGs. We also extend our algorithm to multi-agent general-sum Markov games, achieving a breakthrough in breaking the curse of multiagency.

032

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

#### 1 INTRODUCTION

033 Multi-agent reinforcement learning (MARL), which aims to develop algorithms for multiple agents 034 to learn and make decisions in dynamic environments, has gained significant attention in areas such as game playing (Silver et al., 2017), autonomous driving (Bhalla et al., 2020), and Path 035 Planning (Cao et al., 2020). Under the constraints on time or resources, a key challenge in applying MARL to real-world scenarios is the restricted ability to interact or explore the environment. Offline 037 MARL, also named as batch MARL, addresses this issue by utilizing historical data collected from past interactions, often generated by unknown behavior policies. Researchers hope that this data can offer valuable insights into the optimal policy without the need for further exploration (Lambert 040 et al., 2022). Beyond seeking to maximize the expected total rewards, a critical challenge lies 041 in addressing environmental uncertainties stemming from model mismatches, system noise, and 042 the gap between simulation and real-world situations. Standard MARL algorithms that train in 043 ideal conditions are highly sensitive and prone to catastrophic failure when faced with even small 044 adversarial perturbations in the deployment environment (Zhang et al., 2020; Yeh et al., 2021; Zeng et al., 2022). However, historical data is often gathered under the assumption of model stability, which is unrealistic due to the time-varying and non-stationary nature of real-world systems. Thus, 046 the robust guarantee is critical in offline settings, leading to the formulation of offline robust MARL. 047

As a specific setting of MARL, two-player zero-sum Markov games (TZMGs) are a fascinating area of research, thus leading the field of robust TZMGs (RTZMGs) following from robust MARL.
The inherent solution concepts for RTZMGs encompass equilibria not just between the two players but also between their adversaries, who select the worst-case environments from a predefined uncertainty set for each player. This structure inherently offers greater robustness and stability when facing unmodeled disruptions. Despite recent efforts (Kardeş et al., 2011; Blanchet et al., 2024; Zhang et al., 2020; Ma et al., 2023), there is still a lack of fundamental understanding in learning for

RTZMGs. For a tabular RTZMG with horizon length H, states S, actions  $\{A, B\}$ , and uncertainty sizes  $\{\sigma^+, \sigma^-\}$  for the two players, the best sample complexity for offline setting so far is achieved by  $P^2M^2PO$  (Blanchet et al., 2024) with a near-optimal sample complexity on H, S,  $\{A, B\}$ , where however the influence of uncertainty levels is overlooked. Notably, historical data often only offers partial and limited coverage of the state-action space, leading to poor estimates of model parameters and, in turn, unreliable policy learning outcomes. We summarize previous works and present them along with our results in Table 1. Consequently, current solutions lack an algorithm with optimal sample complexity under partial coverage. Thus, we explore the unresolved question as follows:

> Can we achieve effective sample complexity with robustness to learn Nash policy under partial and limited coverage in TZMGs simultaneously?

Table 1: A comparison between RTZ-VI-LCB and  $P^2M^2PO$  (Blanchet et al., 2024) on finding an  $\varepsilon$ -optimal robust Nash policy in finite-horizon offline RTZMGs with  $f(\sigma^+, \sigma^-, H) = \min\left\{\frac{(H\sigma^+ - 1 + (1-\sigma^+)^H)}{(\sigma^+)^2}, \frac{(H\sigma^- - 1 + (1-\sigma^-)^H)}{(\sigma^-)^2}, H\right\}$ , where the uncertainty set is quantified by total variation (TV) distance. The sample complexities omit all logarithmic factors.

Algorithm	Sample complexity	Uncertainty level		
$P^2M^2PO$	$\frac{C_{\rm r}H^5S^2AB}{\varepsilon^2}$	not consider		
RTZ-VI-LCB (Ours)	$\frac{C_{\mathbf{r}}^{\star}H^4S(A+B)}{\varepsilon^2}f(\sigma^+,\sigma^-,H)$	full range		
Lower bound	$\frac{C_{\rm r}^{\star}SH^4(A+B)}{\varepsilon^2}$	$\min\left\{\sigma^+, \sigma^-\right\} \lesssim \frac{1}{H}$		
Lower bound	$\frac{C_{\rm r}^{\star} SH^3(A+B)}{\varepsilon^2 \min\{\sigma^+, \sigma^-\}}$	$\min\left\{\sigma^+,\sigma^-\right\}\gtrsim \frac{1}{H}$		

#### 1.1 CONTRIBUTION

062

063

064 065

067

079 080

082

084

085

090

092

093

094

097

098

We aim to understand and achieve effective sample complexity under partial convergence in RTZMGs. Our contributions are outlined as follows.

- We introduce a concept to evaluate the quality of historical data, which is the robust unilateral clipped concentrability coefficient  $C_r^{\star} \in \left[\frac{1}{S(A+B)}, \infty\right)$ . This coefficient captures the distribution shift between the behavior policy  $(\mu^n, \nu^n)$  and the single optimal robust policies  $(\mu, \nu^{\star})$  and  $(\mu^{\star}, \nu)$  under model perturbations, without requiring full coverage of the state-action space by the behavior policy. In contrast,  $P^2M^2PO$  (Blanchet et al., 2024) measures distribution mismatch using the maximum density ratio  $C_r$ , which is less tight than our robust unilateral clipped concentrability coefficient  $C_r^{\star}$ .
- We design a new model-based algorithm for offline RTZMGs, an optimistic variant of robust value iteration (VI) for RTZMGs named RTZ-VI-LCB. Specifically, RTZ-VI-LCB incorporates a plug-in estimator of the nominal transition kernel (Iyengar, 2005) and introduces a data-informed penalty to the robust value estimates. armed with TV distance, we show that this algorithm achieves an  $\varepsilon$ -optimal robust Nash equilibrium (NE) policy up to some logarithmic factor as long as the sample size exceeds  $\widetilde{O}\left(\frac{C_{\star}^{*}H^{4}S(A+B)}{\varepsilon^{2}}\min\left\{\frac{(H\sigma^{+}-1+(1-\sigma^{+})^{H})}{(\sigma^{+})^{2}},\frac{(H\sigma^{-}-1+(1-\sigma^{-})^{H})}{(\sigma^{-})^{2}},H\right\}\right)$  after a burn-in cost independent of  $\varepsilon$ . To the best of our knowledge, this is the first time optimal dependency on state *S* and actions  $\{A, B\}$  has been achieved for offline RTZMGs.
- In addition to the upper bound, we derive information-theoretic lower bounds across various uncertainty levels, independent of the specific distance metric applied. We show that there exists an algorithm requiring at least Ω (C<sup>\*</sup><sub>x</sub>SH<sup>4</sup>(A+B)/ε<sup>2</sup>) samples to find an ε-optimal robust NE policy when the uncertainty level min {σ<sup>+</sup>, σ<sup>-</sup>} ≤ 1/H, and at least Ω (C<sup>\*</sup><sub>x</sub>SH<sup>3</sup>(A+B)/ε<sup>2</sup>min{σ<sup>+</sup>, σ<sup>-</sup>}) samples when min {σ<sup>+</sup>, σ<sup>-</sup>} ≥ 1/H. This indicates that learning RTZMGs is at least as challenging as standard TZMGs (Jin et al., 2022) when the uncertainty is sufficiently small. Besides, we confirm the optimality of RTZ-VI-LCB across different uncertainty levels of the critical parameters, i.e., state S and actions {A, B}, except for the finite-horizon H.

• We design an extended algorithm of RTZ-VI-LCB for robust multi-player general-sum Markov games (named Multi-RTZ-VI-LCB) and achieve an  $\varepsilon$ -optimal robust NE policy in  $\widetilde{O}\left(\frac{C_r^* H^4 S \sum_{i=1}^m A_i}{\varepsilon^2} \min\left\{\left\{\frac{(H\sigma_i - 1 + (1 - \sigma_i)^H)}{(\sigma_i)^2}\right\}_{i=1}^m, H\right\}\right)$  samples with M players and  $A_i$  actions and uncertainty size  $\sigma_i$  per player.

113 114 1.2 RELATED WORK

108

110 111

112

115

116

117

In this section, we review a curated selection of related research, with an emphasis on provably efficient RL algorithms in the tabular setting, as these are the most pertinent to our work.

118 Finite-sample studies of standard TZMGs. Markov games (MGs), or called stochastic games, 119 were first proposed in the early 1950s (Shapley, 1953). Since then, extensive research has been 120 conducted, and MARL has gained significant attention (Oroojlooy & Hajinezhad, 2023), particularly 121 around Nash equilibrium (Littman, 1994; Lee et al., 2020). Numerous MARL algorithms with 122 provable convergence and asymptotic guarantees have been developed (Rashid et al., 2020). More 123 recent work has focused on creating algorithms for standard MARL with non-asymptotic guarantees through finite-sample analysis. In this area, most efforts to compute Nash equilibria are focused on 124 TZMGs. The studies in (Bai & Jin, 2020) and (Xie et al., Jun. 2022) were the first to provide 125 non-asymptotic sample complexity guarantees for model-based (e.g., VI-Explore and VI-ULCB) 126 and model-free algorithms (e.g., OMNI-VI). Further improvements in sample complexity have been 127 explored (Cui et al., 2023; Chen et al., 2022; Liu et al., July 2021; Feng et al., 2023; Li et al., 2024c). 128

129 **Robustness in MARL.** Although progress has been made in standard MARL, existing algorithms 130 may struggle when faced with environmental disturbances or uncertainties, leading to significantly 131 deviated equilibria. Increasing research now focuses on enhancing MARL robustness against 132 uncertainties in different parts of MGs (Vial et al., 2022), including state (Zhou & Liu, 2023), 133 environment (reward and transition dynamics), agent types (Zhang et al., 2021), and other agents' 134 policies (Kannan et al., 2023). A typical method to address robustness against uncertainties of 135 the environment is distributionally robust optimization (DRO), which is a method predominantly 136 explored in supervised learning (Bertsimas et al., 2018; Gao, 2023; Blanchet & Murthy, 2019). The application of DRO to manage model uncertainty in single-agent RL (Iyengar, 2005) has attracted 137 considerable attention. However, when extended to MARL, researchers formulated the problem as 138 robust MGs armed with DRO and developed a relatively understudied field with only a few proven 139 algorithms (Blanchet et al., 2024; Kardeş et al., 2011; Ma et al., 2023; Zhang et al., 2020; Shi 140 et al., 2024b). Thus, relevant algorithms based on partial coverage of datasets while considering the 141 uncertainty level are lacking. 142

143 Single-agent robust RL. In single-agent RL, addressing uncertainties of environments using 144 DRO-such as robust Markov decision processes (MDPs) and distributionally robust dynamic 145 programming—has attracted considerable interest in both theoretical research and practical 146 applications (Badrinath & Kalathil, 2021; Goyal & Grand-Clement, 2023). Recent work has focused 147 on the finite-sample performance of provable robust RL algorithms, exploring different divergence 148 functions for uncertainty sets, various sampling mechanisms, and related challenges (Yang et al., 2023; Blanchet et al., 2024; Shi et al., 2024a). Studies on robust MDPs, particularly relevant 149 here, use uncertainty sets based on TV distance (Liu & Xu, 2024) or Kullback-Leibler (KL) 150 divergence (Shi & Chi, 2024) in tabular settings. It has been shown that addressing robust MDPs 151 does not demand more samples compared with those needed for standard MDPs (Shi et al., 2024a). 152 However, RTZMGs present additional complexities beyond those in robust single-agent RL. 153

154

## 2 PROBLEM FORMULATION

155 156

We focus on offline RTZMGs in this paper, which is a robust version of standard TZMGs taking environmental uncertainties into consideration. RTZMGs form a broader class than standard TZMGs, accommodating various prescribed environmental uncertainty sets. Along with this setting, we investigate an efficient algorithm to achieve robustness and optimal sample complexity on action  $\{A, B\}$  without requiring full coverage of the state-action space. An RTZMG under the finite-horizon setting can be defined as  $\mathcal{MG}_{r} = \{S, \mathcal{A}, \mathcal{B}, \mathcal{U}_{0}^{\sigma^{+}}(P^{0}), \mathcal{U}_{0}^{\sigma^{-}}(P^{0}), r, H\}$ , where

 $\mathcal{S} \coloneqq \{1, \cdots, S\}$  is the state space of size S;  $(\mathcal{A} \coloneqq \{1, \cdots, A\}, \mathcal{B} \coloneqq \{1, \cdots, B\})$  denotes the action spaces of the max-player and the min-player with sizes A and B, respectively; H is the horizon length;  $r = \{r_h\}_{h=1}^{H}$  represents the immediate reward obtained at time step h. Specifically,  $r_h(s, a, b)$  is assumed to be deterministic on a state-action pair (s, a, b) and falls within the range [0, 1]. In RTZMGs, this reward can represent both the gain of the max-player and the loss of the min-player. A crucial difference from standard TZMGs is that, rather than assuming a fixed transition kernel, both players in RTZMGs expect that the transition kernel could be chosen arbitrarily from specified uncertainty sets,  $\mathcal{U}_{\rho}^{\sigma^+}(P^0)$  and  $\mathcal{U}_{\rho}^{\sigma^-}(P^0)$ , respectively. These uncertainty sets are centered on a nominal kernel  $P^0: \mathcal{S} \times \mathcal{A} \times \mathcal{B} \mapsto \Delta(\mathcal{S})$ , with their size and shape defined by a distance metric  $\rho$ and radius parameters  $\sigma^+ > 0$  and  $\sigma^- > 0$ . To accommodate individual robustness preferences, the max-player and min-player can independently define their uncertainty sets  $\mathcal{U}_{\rho}^{\sigma^+}(P^0)$  and  $\mathcal{U}_{\rho}^{\sigma^-}(P^0)$ , selecting different sizes ( $\sigma^+ > 0$  and  $\sigma^- > 0$ ) and potentially different divergence functions ( $\rho$ ) for shaping the sets. In this paper, we consider the same divergence function for both players.

**Uncertainty set with** two-player-wise (s, a, b)-rectangularity. We define the transition kernel uncertainty sets  $\mathcal{U}_{\rho}^{\sigma^+}(P^0)$  and  $\mathcal{U}_{\rho}^{\sigma^-}(P^0)$  for RTZMGs. Inspired by the *rectangularity* condition used in robust single-agent RL (Shi et al., 2024a; Iyengar, 2005), we adapt this concept to a two-player setting, termed two-player-wise (s, a, b)-rectangularity. The adaptation enhances computational tractability and facilitates the robust version of Bellman recursions. It permits each player to select its uncertainty set independently, which can be decomposed for each state-action pair into a product of subsets. Consequently, the uncertainty sets  $\mathcal{U}_{\rho}^{\sigma^+}(P^0)$  and  $\mathcal{U}_{\rho}^{\sigma^-}(P^0)$  for the two players, adhering to two-player-wise (s, a, b)-rectangularity, are mathematically defined as: 

$$\mathcal{U}_{\rho}^{\sigma^{+}}\left(P^{0}\right) := \otimes \mathcal{U}_{\rho}^{\sigma^{+}}\left(P_{h,s,a,b}^{0}\right), \qquad \mathcal{U}_{\rho}^{\sigma^{-}}\left(P^{0}\right) := \otimes \mathcal{U}_{\rho}^{\sigma^{-}}\left(P_{h,s,a,b}^{0}\right), \tag{1}$$

where

$$\mathcal{U}_{\rho}^{\sigma^{+}}\left(P_{h,s,a,b}^{0}\right) := \left\{P_{h,s,a,b} \in \Delta(\mathcal{S}) : \rho\left(P_{h,s,a,b}, P_{h,s,a,b}^{0}\right) \le \sigma^{+}\right\}.$$

Here,  $\otimes$  represents the Cartesian product. The uncertainty set for min-player can be defined similarly. We define a vector of the transition kernel P or  $P^0$  at any state-action pair (s, a, b) as

$$P_{h,s,a,b} := P_h(\cdot \,|\, s, a, b) \in \mathbb{R}^{1 \times S}, \qquad P_{h,s,a,b}^0 := P_h^0(\cdot \,|\, s, a, b) \in \mathbb{R}^{1 \times S}.$$
(2)

Here, the distance function  $\rho$  for each player's uncertainty set can be selected from various options that quantify differences between probability vectors. These include f-divergences (such as KL divergence, TV distance, and chi-square) (Yang et al., 2022), the Wasserstein distance (Xu et al., 2023), and  $\ell_q$  norms (Clavier et al., 2023). 

Robust value functions. In RTZMGs, players seek to optimize their worst-case performance across all possible transition kernels within their respective uncertainty sets  $\mathcal{U}_{0}^{\sigma^{+}}(P^{0})$  and  $\mathcal{U}_{\rho}^{\sigma^{-}}(P^{0})$ . For any product policy  $(\mu \times \nu) \in \Delta(\mathcal{A} \times \mathcal{B})$ , the max-player's worst-case performance at time step h is quantified by the robust value function  $V_h^{\mu,\nu,\sigma^+}$  and the robust Q-function  $Q_h^{\mu,\nu,\sigma^+}$ for all  $(h, s, a, b) \in [H] \times S \times A \times B$ , defined as: 

$$V_{h}^{\mu,\nu,\sigma^{+}}(s) \coloneqq \inf_{P \in \mathcal{U}_{\rho}^{\sigma^{+}}(P^{0})} V_{h}^{\mu,\nu,P}(s) \quad \text{and} \quad Q_{h}^{\mu,\nu,\sigma^{+}}(s,a,b) \coloneqq \inf_{P \in \mathcal{U}_{\rho}^{\sigma^{+}}(P^{0})} Q_{h}^{\mu,\nu,P}; \quad (3)$$

$$V_{h}^{\mu,\nu,\sigma^{-}}(s) \coloneqq \sup_{P \in \mathcal{U}_{o}^{\sigma^{-}}(P^{0})} V_{h}^{\mu,\nu,P}(s) \quad \text{and} \quad Q_{h}^{\mu,\nu,\sigma^{-}}(s,a,b) \coloneqq \sup_{P \in \mathcal{U}_{o}^{\sigma^{-}}(P^{0})} Q_{h}^{\mu,\nu,P}, \quad (4)$$

where

$$V_h^{\mu,\nu,P}(s) \coloneqq \mathbb{E}_{\mu,\nu,P}\left[\sum_{t=h}^H r_t(s_t, a_t, b_t) \mid s_h = s\right];$$

ΓН 

215 
$$Q_h^{\mu,\nu,P}(s,a,b) \coloneqq \mathbb{E}_{\mu,\nu,P}\left[\sum_{t=h}^{\infty} r_t(s_t,a_t,b_t) \mid s_h = s, a_h = a, b_h = b\right].$$

221 222

223

224

225 226 227

228

234

235

241

242

243

244 245

246

247 248 249

261 262

**Offline dataset.** Let  $\mathcal{D}$  be a dataset consisting of K episodes under independence, with each episode produced by implementing a behavior policy  $\{\mu_h^n, \nu_h^n\}_{h=1}^H$  in a nominal MDP  $\mathcal{M}^0 = (\mathcal{S}, \mathcal{A}, \mathcal{B}, H, P^0) := \{P_h^0\}_{h=1}^H, \{r_h\}_{h=1}^H)$ . For  $1 \leq k \leq K$ , the k-th episode  $(s_1^k, a_1^k, b_1^k, \dots, s_H^k, a_H^k, b_H^k, s_{H+1}^k)$  is generated as follows:

$$s_{1}^{k} \sim \varrho^{\mathsf{n}}, \quad a_{h}^{k} \sim \mu_{h}^{\mathsf{n}}(\cdot \,|\, s_{h}^{k}), \quad b_{h}^{k} \sim \nu_{h}^{\mathsf{n}}(\cdot \,|\, s_{h}^{k}), \quad s_{h+1}^{k} \sim P_{h}^{0}(\cdot \,|\, s_{h}^{k}, a_{h}^{k}, b_{h}^{k}), \quad 1 \le h \le H.$$
(5)

Throughout this paper, let  $\varrho^n$  denote the initial distribution related to a historical dataset. We use the short-hand notation for the occupancy distribution w.r.t. the behavior policy  $(\mu^n, \nu^n)$  as:  $\forall (h, s, a, b) \in [H] \times S \times A \times B$ ,

$$d_{h}^{\mathsf{n},P^{0}}(s) = d_{h}^{\mu^{\mathsf{n}},\nu^{\mathsf{n}},P^{0}}(s) \coloneqq \mathbb{P}(s_{h} = s \,|\, s_{1} \sim \varrho^{\mathsf{n}}, \mu^{\mathsf{n}}, \nu^{\mathsf{n}}, P^{0});$$
(6a)

$$d_{h}^{\mathsf{n},P^{0}}(s,a,b) = d_{h}^{\mu^{\mathsf{n}},\nu^{\mathsf{n}},P^{0}}(s,a,b) \coloneqq \mathbb{P}(s_{h}=s \mid s_{1} \sim \varrho^{\mathsf{n}}, \mu^{\mathsf{n}}, \nu^{\mathsf{n}}, P^{0}) \, \mu_{h}^{\mathsf{n}}(a \mid s) \, \nu_{h}^{\mathsf{n}}(b \mid s). \tag{6b}$$

Similarly, for any product policy  $(\mu, \nu)$ , there is,  $\forall (h, s, a, b) \in [H] \times S \times A \times B$ 

$$d_h^{\mu,\nu,P}(s) \coloneqq \mathbb{P}(s_h = s \mid s_1 \sim \varrho, \mu, \nu, P); \tag{7a}$$

$$d_h^{\mu,\nu,P}(s,a,b) \coloneqq \mathbb{P}(s_h = s \mid s_1 \sim \varrho, \mu, \nu, P) \,\mu_h(a \mid s) \,\nu_h(b \mid s). \tag{7b}$$

**Robust Bellman equations.** RTZMGs include a robust version of the Bellman equation, referred to as the *robust Bellman equation*. The robust value functions  $V_h^{\mu,\nu,\sigma^+}(s)$  for max-player in RTZMGs, associated with any product policy  $(\mu, \nu)$ , satisfy:  $\forall (h, s) \in [H] \times S$ ,

$$V_{h}^{\mu,\nu,\sigma^{+}}(s) = \mathbb{E}_{a \sim \mu_{h}(a), b \sim \nu_{h}(a)} \left[ r_{h}(s,a,b) + \inf_{P \in \mathcal{U}_{\rho}^{\sigma^{+}}(P_{h,s,a,b}^{0})} PV_{h+1}^{\mu,\nu,\sigma^{+}} \right].$$
 (8)

 $V_h^{\mu,\nu,\sigma}$  (s) for min-player can be obtained similarly. We highlight that the robust Bellman equations are intrinsically connected to the *two-player-wise* (s, a, b)-*rectangularity* condition (see (1)) applied to the uncertainty set. This condition separates the dependencies of uncertainty subsets among different time steps, the players, and state-action pairs, thus leading to the Bellman recursion.

**Optimal robust policy.** We further define the maximum robust value function with fixed opponent policy for each player as:  $\forall (h, s) \in [H] \times S$ ,

$$V_{h}^{\star,\nu,\sigma^{+}}(s) \coloneqq \max_{\mu:\mathcal{S}\times[H]\mapsto\Delta(\mathcal{A})} V_{h}^{\mu,\nu,\sigma^{+}}(s) = \max_{\mu:\mathcal{S}\times[H]\mapsto\Delta(\mathcal{A})} \inf_{P\in\mathcal{U}_{o}^{\sigma^{+}}(P^{0})} V_{h}^{\mu,\nu,P}(s).$$
(9)

250 Optimal robust policy for min-player can be obtained similarly. As proved by Blanchet et al. (2024), 251 there is at least one policy referred to as  $\mu_h^*(s) : S \times [H] \mapsto \Delta(\mathcal{A})$  (for the max-player) and 252  $\nu_h^*(s) : S \times [H] \mapsto \Delta(\mathcal{B})$  (for the min-player), corresponding to as the *robust best-response policy*. 253 These policies can simultaneously achieve  $V_h^{\star,\nu,\sigma^+}(s)$  (for the max-player) and  $V_h^{\mu,\star,\sigma^-}(s)$  (for the 254 min-player) for all  $s \in S$  and  $h \in [H]$ .

**Robust Nash equilibrium.** In RTZMGs, the dynamics expand beyond traditional TZMGs to involve four participants: two players and two adversaries determining the worst-case transitions. Therefore, finding an equilibrium becomes central in RTZMGs due to potentially conflicting objectives. We introduce the robust variant of standard solution concepts—robust NE for RTZMGs. A product policy  $(\mu, \nu)$  is considered a *robust NE* if

$$\forall (s) \in \mathcal{S}, \quad V_h^{\star,\nu,\sigma^+}(s) = V_h^{\star,\sigma^+}(s); \quad V_h^{\mu,\star,\sigma^-}(s) = V_h^{\star,\sigma^-}(s). \tag{10}$$

A robust NE signifies that given the product policy  $(\mu, \nu)$  of the opponents, no player can enhance their outcome by deviating from their current policy unilaterally when each player accounts for the worst-case scenario within their uncertainty set  $\mathcal{U}_{\rho}^{\sigma^+}(P^0)$  or  $\mathcal{U}_{\rho}^{\sigma^-}(P^0)$ .

Since finding exact robust equilibria can be complex and may not always be feasible, practitioners often seek approximate equilibria. In this context, a product policy  $(\mu \times \nu) \in \Delta(\mathcal{A} \times \mathcal{B})$  can be termed an  $\varepsilon$ -robust NE if

$$\operatorname{Gap}(\mu,\nu) \coloneqq \max\left\{ V_1^{\star,\nu,\sigma^+}(\varrho) - V_1^{\star,\sigma^+}(\varrho), \ V_1^{\star,\sigma^-}(\varrho) - V_1^{\mu,\star,\sigma^-}(\varrho) \right\} \le \varepsilon, \tag{11}$$

where 

$$V_1^{\star,\nu,\sigma^+}(\varrho) = \mathbb{E}_{s \sim \varrho} V_1^{\star,\nu,\sigma^+}(s), \quad \text{and} \quad V_1^{\star,\sigma^+}(\varrho) = \mathbb{E}_{s \sim \varrho} V_1^{\star,\sigma^+}(s).$$

The definitions of  $V_1^{\mu,\star,\sigma^-}(\varrho)$  and  $V_1^{\star,\sigma^-}(\varrho)$  can be obtained similarly. The existence of robust NE has been proved for general divergence functions in the uncertainty set by Blanchet et al. (2024).

**Learning objective** With a dataset collected from the nominal environment, our objective is to find a solution among the  $\varepsilon$ -robust NEs for the RTZMG  $\mathcal{MG}_r$  with respect to a specified uncertainty set  $\mathcal{U}(P^0)$  around the nominal kernel, while minimizing the number of samples required under partial coverage of the state-action space.

#### ALGORITHM DESIGN

In this section, we propose an efficient model-based algorithm RTZ-VI-LCB to achieve robustness and optimal sample complexity on action  $\{A, B\}$ . This algorithm is designed for offline RTZMGs within the finite-horizon setting.

#### 3.1 BUILDING AN EMPIRICAL NOMINAL MDP

According to the empirical frequencies of state transitions, we can naturally construct an empirical estimate  $\widehat{P}^0 = \{\widehat{P}_h^0\}_{h=1}^H$  of  $P^0$ , where

$$\hat{P}_{h}^{0}\left(s'\,|\,s,a,b\right) = \begin{cases} \frac{1}{N_{h}\left(s,a,b\right)} \sum_{i=1}^{N} \mathbb{1}\left\{\left(s_{i},a_{i},b_{i},s'_{i}\right) = \left(s,a,b,s'\right)\right\}, & \text{if } N_{h}\left(s,a,b\right) > 0; \\ \frac{1}{S}, & \text{if } N_{h}\left(s,a,b\right) = 0, \end{cases}$$
(12)

$$\widehat{r}_{h}(s, a, b) = \begin{cases} r_{h}(s, a, b), & \text{if } N_{h}(s, a, b) > 0; \\ 0, & \text{if } N_{h}(s, a, b) = 0, \end{cases}$$
(13)

for any  $(h, s, a, b, s') \in [H] \times S \times A \times B \times S$ . Besides,  $N_h(s, a, b)$  represents the total number of sample transitions from (s, a, b) at step h, and

$$N_h(s, a, b) \coloneqq \sum_{i=1}^N \mathbb{1}\{(s_i, a_i, b_i) = (s, a, b)\}.$$
(14)

Algorithm 1: Two-stage subsampling technique for RTZ-VI-LCB.

- **Input:** Dataset  $\mathcal{D}$ , probability  $\delta$ .
  - <sup>2</sup> Step 1: Data Partitioning. Split  $\mathcal{D}$  into two equal-sized subsets,  $\mathcal{D}^{m}$  and  $\mathcal{D}^{a}$ , each containing K/2 trajectories.
  - <sup>3</sup> Step 2: Defining Transition Bounds. For step h and state s, denote the number of transitions from  $\mathcal{D}^{\mathsf{m}}$  (resp.  $\mathcal{D}^{\mathsf{a}}$ ) as  $N_h^{\mathsf{m}}(s)$  (resp.  $N_h^{\mathsf{a}}(s)$ ). Construct the trimmed count as:

$$N_h^{\mathsf{t}}(s) \coloneqq \max\left\{N_h^{\mathsf{a}}(s) - 10\sqrt{N_h^{\mathsf{a}}(s)\log\frac{HS}{\delta}}, 0\right\};\tag{15}$$

- 4 Step 3: Generating Subsampled Dataset. Randomly sample transitions (quadruples of the form (s, a, b, h, s') from  $\mathcal{D}^{\mathsf{m}}$  uniformly. For each  $(s, h) \in \mathcal{S} \times [H]$ , include  $\min\{N_h^t(s), N_h^m(s)\}$  transitions in the new dataset  $\mathcal{D}^t$ .
- **Output:** Set  $\mathcal{D}_0 = \mathcal{D}^t$ .

Although it is feasible to decompose the historical dataset  $\mathcal{D}$  into sample transitions, the dependencies between transitions within the same episode introduce complexities in our analysis. To address this issue, Li et al. (2024a) introduced a two-fold subsampling method for single-agent RL to preprocess  $\mathcal{D}$ , thereby reducing statistical dependencies and producing a distributionally equivalent dataset  $\mathcal{D}_0$  with independent samples. We adapt this method to TZMGs, as outlined in Algorithm 1. We present the following lemma concerning the dataset  $\mathcal{D}_0$ , which is proved in Appendix C.1.

**Lemma 1** The dataset produced by the two-stage subsampling method is distributionally identical to  $\mathcal{D}_0$  with probability at least  $1-8\delta$ , where  $\{N_h(s,a,b)\}$  are independent of the sample transitions in  $\mathcal{D}^0$  and obey:  $\forall (h, s, a, b) \in [H] \times \mathcal{S} \times \mathcal{A} \times \mathcal{B}$ , 

$$N_h(s,a,b) \ge \frac{Kd_h^{\mathsf{n}}(s,a,b)}{8} - 5\sqrt{Kd_h^{\mathsf{n}}(s,a,b)\log\frac{KH}{\delta}}.$$
(16)

By applying the two-fold sampling method, we can treat the dataset  $\mathcal{D}_0$  as having independent samples, simplifying the analysis significantly as supported by Lemma 1.

AN OPTIMISTIC VARIANT OF ROBUST VI WITH LOWER CONFIDENCE BOUNDS. 3.2

We propose a model-based approach for solving RTZMGs using an approximate  $\hat{P}^0$  for  $P^0$ , which is the nominal transition kernel. Specifically, we introduce VI with lower confidence bounds (LCBs) for RTZMGs (RTZ-VI-LCB) to compute a robust NE for two players, as summarized in Algorithm 2. 

Our algorithm begins at the final time step h = H and proceeds backward through h = H - 1, H - 12,...,1. Drawing from the principle of pessimism in single-agent offline RL (Li et al., 2024a; Jin et al., 2021), we design an optimistic robust Q-value to estimate the robust Q-function at time step  $h \in [H]$  as  $\widehat{Q}_h^+$  and  $\widehat{Q}_h^-$  for all  $(h, s, a, b) \in [H] \times \mathcal{S} \times \mathcal{A} \times \mathcal{B}$ , that is, 

$$\widehat{Q}_{h}^{+}(s,a,b) = \widehat{r}_{h}(s,a,b) + \inf_{P \in \mathcal{U}^{\sigma^{+}}(\widehat{P}_{h,s,a,b}^{0})} P\widehat{V}_{h+1}^{+} + \beta_{h}\left(s,a,b,\widehat{V}_{h+1}^{+}\right);$$
(17a)

$$\widehat{Q}_{h}^{-}(s,a,b) = \widehat{r}_{h}(s,a,b) + \sup_{h=1}^{r \in \mathcal{U}^{+}} P\widehat{V}_{h+1}^{-} - \beta_{h}\left(s,a,b,\widehat{V}_{h+1}^{-}\right).$$
(17b)

$$\widehat{V}_{h}^{-}(s,a,b) = \widehat{r}_{h}(s,a,b) + \sup_{P \in \mathcal{U}^{\sigma^{-}}(\widehat{P}_{h,s,a,b}^{0})} P\widehat{V}_{h+1}^{-} - \beta_{h}\left(s,a,b,\widehat{V}_{h+1}^{-}\right).$$
(17b)

**Dual problem.** Solving (17) directly is computationally intensive because it requires optimizing over an S-dimensional probability simplex, which becomes exponentially more difficult as the state space size S increases. In fortunate, strong duality for TV distance allows us to tackle this problem by solving its dual (Iyengar, 2005):

$$\inf_{P \in \mathcal{U}^{\sigma^+}(\hat{P}^0_{h,s,a,b})} P \hat{V}^+_{h+1} = \max_{\alpha \in [\min_s \hat{V}^+_{h+1}, \max_s \hat{V}^+_{h+1}]} \Big\{ \hat{P}^0_{h,s,a,b} \Big[ \hat{V}^+_{h+1} \Big]_{\alpha} - \sigma^+ \Big( \alpha - \min_{s'} \Big[ \hat{V}^+_{h+1} \Big]_{\alpha} (s') \Big) \Big\}.$$
(18)

where  $\left[\widehat{V}_{h+1}^+\right]_{\alpha}$  denotes the clipped versions of  $\widehat{V}_{h+1}^- \in \mathbb{R}^S$  and  $\widehat{V}_{h+1}^+ \in \mathbb{R}^S$  based on some level  $\alpha \geq 0$ , as follows.  $\sup_{P \in \mathcal{U}^{\sigma^-}(\widehat{P}^0_{h,s,a,b})} P \widehat{V}^-_{h+1}$  can be defined similarly. See Appendix A for details.

$$\left[\widehat{V}_{h+1}^{+}\right]_{\alpha}(s) := \begin{cases} \widehat{V}_{h+1}^{+}(s), & \text{if } \widehat{V}_{h+1}^{+}(s) > \alpha; \\ \alpha, & \text{otherwise;} \end{cases}$$
(19)

**Penalty term.** The optimistic robust Q-function estimate is refined by  $\beta_h(s, a, b, \hat{V})$ , which is a data-driven penalty term and includes the uncertainty in value estimates. We adopt the Bernsteinstyle penalty to better capture the variance structure over time. In particular, for any  $(s, a, b, h) \in$  $\mathcal{S} \times \mathcal{A} \times \mathcal{B} \times [H]$  and  $\delta \in (0, 1)$ , the penalty term  $\beta_h(s, a, b, \widehat{V})$  is defined as:

$$\beta_h\left(s,a,b,\widehat{V}\right) = \min\left\{\max\left\{\sqrt{\frac{C_{\mathsf{n}}\log\frac{KH}{\delta}}{N_h\left(s,a,b\right)}}\mathsf{Var}_{\widehat{P}^0_{h,s,a,b}}(\widehat{V})}, \frac{2C_{\mathsf{n}}H\log\frac{KH}{\delta}}{N_h\left(s,a,b\right)}\right\}, H\right\}, \quad (20)$$

where  $C_n$  is some universal constant, and

ν

$$\operatorname{Var}_{\widehat{P}^{0}_{h,s,a,b}}\left(\widehat{V}\right) \coloneqq \widehat{P}^{0}_{h,s,a,b}\widehat{V}^{2} - (\widehat{P}^{0}_{h,s,a,b}\widehat{V})^{2}.$$
(21)

Note that we choose  $\widehat{P}^0$ , as opposed to  $P^0$  (i.e.,  $\operatorname{Var}_{\widehat{P}^0_{h,s,a,b}}(\widehat{V})$ ) in the variance term, since we have no access to the true transition kernel  $P^0$ . This penalty term is distinct from those used in standard offline TZMGs (Cui et al., 2023; Li et al., 2024a), as it accounts for the unique structure of robust self-play MDPs. Specifically, it provides a tight upper bound on statistical uncertainty, considering the non-linear and implicit dependency introduced by the uncertainty set  $\mathcal{U}(P^0)$ , addressing challenges not present in standard MDP scenarios.

**Policy estimation.** We update the policies using the estimated Q-functions with uncertainty as line 6 in Algorithm 2. Specifically, for any matrix  $\mathbf{N} \in \mathbb{R}^{A \times B}$ , the function ComputNash( $\mathbf{N}$ ) returns a solution  $(\widehat{w}, \widehat{z})$  to the minimax problem  $\max_{w \in \Delta(\mathcal{A})} \min_{z \in \Delta(\mathcal{B})} w^{\top} \mathbf{N} z$ . In other words, for each  $s \in S$ , we compute the NE policies  $(\mu_h^+(s), \nu_h^+(s))$  and  $(\mu_h^-(s), \nu_h^-(s)) \in \Delta(\mathcal{A}) \times \Delta(\mathcal{B})$  for the zero-sum matrix games with payoff matrices  $\widehat{Q}_h^+(s, \cdot, \cdot)$  and  $\widehat{Q}_h^-(s, \cdot, \cdot)$ , respectively. Solving these robust matrix games is generally PPAD-hard due to the potential for players to choose different worst-case transition kernels.

Algorithm 2: Value iteration with lower confidence bounds for RTZMGs (RTZ-VI-LCB).

1 Initialization: Set uncertainty levels  $\sigma^-$  and  $\sigma^+$ ; set  $\widehat{V}_h^-(s) = 0$  and  $\widehat{V}_h^+(s) = H$  for all  $(s,h)\in\mathcal{S}\times[H+1]; \text{set }\widehat{Q}_h^-(s,a,b)=0 \text{ and } \widehat{Q}_h^+(s,a,b)=H \text{ for all } \mathbb{R}$  $(s, a, b, h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{B} \times [H+1].$ <sup>2</sup> Compute the empirical reward function  $\hat{r}$  using (13) and the empirical transition kernel  $P_0$ using (12). 3 for h = H, H - 1, ..., 1 do Update the robust Q-value estimate as  $\hat{Q}_{h}^{+}(s,a,b) = \min\left\{\hat{r}_{h}(s,a,b) + \inf_{P \in \mathcal{U}^{\sigma^{+}}(\hat{P}_{h,s,a,b}^{0})} P\hat{V}_{h+1}^{+} + \beta_{h}\left(s,a,b,\hat{V}_{h+1}^{+}\right), H\right\};$  $\widehat{Q}_{h}^{-}\left(s,a,b\right) = \max\left\{\widehat{r}_{h}\left(s,a,b\right) + \sup_{P \in \mathcal{U}^{\sigma^{-}}\left(\widehat{P}_{h,s,a,b}^{0}\right)} P\widehat{V}_{h+1}^{-} - \beta_{h}\left(s,a,b,\widehat{V}_{h+1}^{-}\right), 0\right\},\$ with  $\beta_h(s, a, b, V) = \min\left\{\max\left\{\sqrt{\frac{C_n \log \frac{KH}{\delta}}{N_h(s, a, b)}} \mathsf{Var}_{\widehat{P}^0_{h, s, a, b}}(V), \frac{2C_n H \log \frac{KH}{\delta}}{N_h(s, a, b)}\right\}, H\right\}.$ **Compute** Nash policy for each  $s \in S$  as  $(\mu_h^+(s), \nu_h^+(s)) = \text{ComputNash}\left(\widehat{Q}_h^+(s, \cdot, \cdot)\right);$  $\left(\mu_{h}^{-}\left(s\right),\nu_{h}^{-}\left(s\right)\right) = \mathsf{ComputNash}\left(\widehat{Q}_{h}^{-}\left(s,\cdot,\cdot\right)\right),$ **Update** the robust value estimate for each  $s \in S$  as  $\hat{V}_{h}^{-}(s) = \mathbb{E}_{a \sim \mu_{k}^{-}(s), b \sim \nu_{k}^{-}(s)} \left[ \hat{Q}_{h}^{-}(s, a, b) \right], \qquad \hat{V}_{h}^{+}(s) = \mathbb{E}_{a \sim \mu_{k}^{+}(s), b \sim \nu_{k}^{+}(s)} \left[ \hat{Q}_{h}^{+}(s, a, b) \right].$ **s Output**: The policy pair  $(\hat{\mu}, \hat{\nu})$ , where  $\hat{\mu} = {\{\mu_h^-\}}_{h=1}^H$  and  $\hat{\nu} = {\{\nu_h^+\}}_{h=1}^H$ . 

4 PERFORMANCE GUARANTEES

**Robust unilateral clipped concentrability.** To assess the effectiveness of the historical dataset for achieving the desired goal, it is essential to measure the distributional discrepancy between the historical data and the target data. Drawing on the *single-policy clipped concentrability* assumption in the single-agent RL (Li et al., 2024a), we propose a novel assumption for RTZMGs as:

Assumption 1 (Robust unilateral clipped concentrability) The behavior policies of the historical dataset D satisfies

for some quantity  $C_r^{\star} \in \left[\frac{1}{S(A+B)}, \infty\right]$ . We define  $C_r^{\star}$  as the smallest value that satisfies (22), referring to it as the robust unilateral clipped concentrability coefficient. For consistency, we adopt the convention 0/0 = 0.

Notably, if  $d_h^{\mu,\nu^*,P}(s,a,b)$  or  $d_h^{\mu^*,\nu,P}(s,a,b)$  is larger than  $\frac{1}{S(A+B)}$ , the robust unilateral clipped concentrability assumption above do not require the data distribution  $d_h^{n,P^0}(s,a,b)$  to scale with  $d_h^{\mu,\nu^*,P}(s,a,b)$  or  $d_h^{\mu^*,\nu,P}(s,a,b)$  proportionally. We here outline the principal theoretical findings concerning the sample complexity of learning robust NE in RTZMGs, including an upper bound for the RTZ-VI-LCB algorithm (Algorithm 2) and an information-theoretic lower bound. Initially, we present the finite-sample guarantee for RTZ-VI-LCB, with detailed proof provided in Appendix B.

**Theorem 1 (Upper bound for RTZ-VI-LCB)** Under the TV uncertainty set  $\mathcal{U}^{\sigma^+}(\cdot)$  and  $\mathcal{U}^{\sigma^-}(\cdot)$ defined in (2) with  $\sigma^+$ ,  $\sigma^- \in (0, 1]$ . Define  $d^n_m = \min_{h,s,a,b} \{d^n_h(s, a, b) : d^n_h(s, a, b) > 0\}$ . Define  $f(\sigma^+, \sigma^-) = \min \left\{ \frac{(H\sigma^+ - 1 + (1-\sigma^+)^H)}{(\sigma^+)^2}, \frac{(H\sigma^- - 1 + (1-\sigma^-)^H)}{(\sigma^-)^2}, H \right\}$ . Consider any  $\delta \in (0, 1)$  and any *RTZMG*  $\mathcal{MG}_r = \{S, \mathcal{A}, \mathcal{B}, \mathcal{U}^{\sigma^+}(P^0), \mathcal{U}^{\sigma^-}(P^0), r, H\}$ . For sufficient large constants  $c_0, c_1 > 0$ , with probability at least  $1 - \delta$ , we can achieve

$$\operatorname{Gap}(\widehat{\mu}, \widehat{\nu}) \le c_1 \sqrt{\frac{C_{\mathsf{r}}^{\star} H^3 S(A+B) \log \frac{KH}{\delta}}{K}} f(\sigma^+, \sigma^-, H),$$
(23)

with the total number of samples T exceeding

$$T = KH \ge c_0 \frac{H^2 S(A+B)}{d_{\mathsf{m}}^{\mathsf{n}}} \log \frac{KH}{\delta} f(\sigma^+, \sigma^-, H).$$
(24)

Now, we introduce a lower bound of sample complexity in RTZMGs, whose proof is in Appendix D.

**Theorem 2 (Lower bound for solving robust MGs)** Consider any tuple  $\mathcal{MG}_{\mathsf{r}} = \{S, \mathcal{A}, \mathcal{B}, \mathcal{U}^{\sigma^+}(P^0), \mathcal{U}^{\sigma^-}(P^0), r, H\}$  obeying  $H > 16 \log 2$  and  $\sigma^+, \sigma^- \in (0, 1 - c_0]$  with any small efficiently positive constant  $0 < c_0 \leq \frac{1}{4}$ . Let

$$\varepsilon \leq \begin{cases} \frac{c_2}{H}, & \text{if } \max\left\{\sigma^+, \sigma^-\right\} \leq \frac{c_2}{2H}, \\ 1 & \text{otherwise} \end{cases}$$
(25)

for any  $c_2 \leq \frac{1}{4}$ . With an initial state distribution  $\rho$ , we can construct a set of RTZMGs  $\left\{ \mathcal{M}_f^{\phi} | f \in \mathcal{F} = \{0, 1, \cdots, SA - 1\}, \phi = [\phi_h]_{1 \leq h \leq H} \in \Phi \subseteq \{0, 1\}^H \right\}$  such that for any dataset with K independent samples trajectories and H lengths per trajectories satisfying  $C \leq C_r^* \leq 2C$ , such that

$$\inf_{\widehat{\mu},\widehat{\nu}} \max_{(f,\phi)\in\mathcal{F}\times\Phi} \left\{ \mathbb{P}_{\phi} \left( \operatorname{Gap}(\widehat{\mu},\widehat{\nu}) > \varepsilon \right) \right\} \ge \frac{1}{8},$$
(26)

474 provided that

$$T = KH \le \frac{c_2 C_{\rm r}^{\star} H^3 S(A+B) \min\{\frac{1}{\min\{\sigma^+, \sigma^-\}}, H\}}{\varepsilon^2}.$$
(27)

Here,  $c_2$  denotes an efficiently small constant. The infimum is obtained over all estimators  $(\widehat{\mu}, \widehat{\nu})$ .

<sup>480</sup> Moreover, our algorithm can be extended to multi-player general-sum Markov games with m players <sup>481</sup> and  $A_i$  actions and uncertainty size  $\sigma_i$  per player with details provided in Appendix F, i.e., Multi-<sup>482</sup> RTZ-VI-LCB. Specifically, we obtain the following theoretical guarantee of Multi-RTZ-VI-LCB:

**Theorem 3 (Upper bound for Multi-RTZ-VI-LCB)** Consider any  $\delta \in (0,1)$  and any robust multi-player general-sum MGs  $\mathcal{MG}_{\mathsf{r}} = \mathcal{M}(\mathcal{S}, \{\mathcal{A}_i\}_{i=1}^m, H, \{\mathcal{U}_{\rho^i}^{\sigma_i}(P^0)\}_{i=1}^m, \{r_i\}_{i=1}^m)$ . Under the TV uncertainty set  $\mathcal{U}^{\sigma_i}(\cdot)$  defined in (2) with  $\sigma_i \in (0,1]$  for  $i = 1, 2, \cdots, m$ . Define  $d_{\mathsf{m}}^{\mathsf{n}} =$   $\min_{h,s,\boldsymbol{a}} \{d_h^{\mathsf{n}}(s,\boldsymbol{a}) : d_h^{\mathsf{n}}(s,\boldsymbol{a}) > 0\}, \text{ and } f(\{\sigma_i\}_{i=1}^m, H) = \min\left\{\left\{\frac{(H\sigma_i - 1 + (1-\sigma_i)^H)}{(\sigma_i)^2}\right\}_{i=1}^m, H\right\}.$   $For sufficient large constants c_0, c_1 > 0, \text{ with probability of at least } 1 - \delta, \text{ we can achieve}$ 

$$\operatorname{Gap}(\widehat{\pi}) \le c_1 \sqrt{\frac{C_{\mathsf{r}}^{\star} H^3 S \sum_{i=1}^m A_i \log \frac{KH}{\delta}}{K}} f(\{\sigma_i\}_{i=1}^m, H),$$
(28)

with the total number of samples T exceeding

489 490 491

496 497

498

499 500 501

504

505

506

507

510

511

512

513

514

515 516

517

518

519

521 522

523 524

$$T = KH \ge c_0 \frac{H^2 S \sum_{i=1}^{m} A_i}{d_{\mathsf{m}}^{\mathsf{n}}} \log \frac{KH}{\delta} f(\{\sigma_i\}_{i=1}^{m}, H).$$
(29)

Here are the key implications of these theorems:

 Theorem 1 demonstrates that the proposed RTZ-VI-LCB algorithm can attain an ε-robust NE solution when the total sample size exceeds:

$$\widetilde{O}\left(\frac{C_{\mathbf{r}}^{\star}H^{4}S(A+B)}{\varepsilon^{2}}\min\left\{\frac{(H\sigma^{+}-1+(1-\sigma^{+})^{H})}{(\sigma^{+})^{2}},\frac{(H\sigma^{-}-1+(1-\sigma^{-})^{H})}{(\sigma^{-})^{2}},H\right\}\right)$$

suggesting that the sample efficiency for robust offline TZMGs is strongly influenced by the dataset quality (quantified by  $C_r^*$ ) and the problem structure of RTZMGs (reflected in the occupancy distributions  $d_m^n$ ). If  $C_r^*$  is as small as  $\frac{1}{S(A+B)}$ , the upper bound of the sample complexity exhibits a weaker dependency on actions  $\{A, B\}$  and state S. Combining this upper bound with the lower bound in Theorem 2 shows that RTZ-VI-LCB's sample complexity is optimal w.r.t. key factors S, A, B and  $\varepsilon$ . This is the first optimal sample complexity upper bound for offline RTZMGs, regarding state S and actions  $\{A, B\}$ .

• Theorem 2 conveys two important points. When the uncertainty level is small (i.e.,  $\min\{\sigma^+, \sigma^-\} \lesssim \frac{1}{H}$ ), no algorithm can find an  $\varepsilon$ -optimal robust policy with fewer than  $\Omega\left(\frac{C_r^*SH^4(A+B)}{\varepsilon^2}\right)$  samples, matching the complexity requirement for non-robust offline TZMGs (Jin et al., 2022). This implies that robust TZMGs are at least as challenging as standard TZMGs for low uncertainty. When the uncertainty level satisfies  $\min\{\sigma^+, \sigma^-\} \gtrsim \frac{1}{H}$ , no algorithm can find an  $\varepsilon$ -optimal robust policy with the numbers of samples fewer than  $\Omega\left(\frac{C_r^*SH^3(A+B)}{\varepsilon^2\min\{\sigma^+,\sigma^-\}}\right)$ . Thus, RTZ-VI-LCB is the first provably near-optimal algorithm on *S* and {*A*, *B*} for RTZMGs without requiring full coverage assumptions.

 Theorem 3 demonstrates that the proposed Multi-RTZ-VI-LCB algorithm can attain an εrobust NE solution when the total sample size exceeds:

$$\widetilde{O}\left(\frac{C_{\mathbf{r}}^{\star}H^4S\sum_{i=1}^{m}A_i}{\varepsilon^2}\min\left\{\left\{\frac{(H\sigma_i-1+(1-\sigma_i)^H)}{(\sigma_i)^2}\right\}_{i=1}^m,H\right\}\right),$$

suggesting that the algorithm can break the curse of multiagency.

#### 5 CONCLUSION

To balance model robustness with sample efficiency, we design an efficient robust model-based algorithm for offline RTZMGs, which is value iteration with lower confidence bounds for RTZMGs (RTZ-VI-LCB). Our algorithm integrates robust VI with the principle of pessimism. By imposing a tailored and mild assumption (robust unilateral clipped concentrability) on the historical dataset to account for the distribution shift, we do not require full state-action space coverage. We address robustness against the distribution shifts in the worse-case scenario of the shared environment, analyze the finite-sample complexity of the proposed RTZ-VI-LCB algorithm, and establish an information-theoretic lower bound to evaluate its optimality across various uncertainty levels.

To the best of our knowledge, this is the first provably optimal algorithm for offline RTZMGs that addresses the dependency on states S and actions  $\{A, B\}$ , while accounting for model perturbations and partial coverage. Furthermore, we extend RTZ-VI-LCB to multi-agent general-sum MGs, demonstrating a breakthrough in breaking the curse of multiagency. Our algorithm opens up several intriguing questions, such as designing efficient model-free algorithms for robust offline TZMGs with partial coverage and exploring ways to adjust the size and metric of the uncertainty set to complete the algorithmic design.

# 540 REFERENCES

579

592

- Alekh Agarwal, Sham Kakade, and Lin F Yang. Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pp. 67–83. PMLR, 2020.
- Kishan Panaganti Badrinath and Dileep Kalathil. Robust reinforcement learning using least squares policy iteration with provable performance guarantees. In *International Conference on Machine Learning*, pp. 511–520. PMLR, 2021.
- Yu Bai and Chi Jin. Provable self-play algorithms for competitive reinforcement learning. In International conference on machine learning, pp. 551–560. PMLR, 2020.
- Dimitris Bertsimas, Vishal Gupta, and Nathan Kallus. Data-driven robust optimization. *Mathematical Programming*, 167:235–292, 2018.
- Sushrut Bhalla, Sriram Ganapathi Subramanian, and Mark Crowley. Deep multi agent reinforcement
   learning for autonomous driving. In *Canadian Conference on Artificial Intelligence*, pp. 67–78.
   Springer, 2020.
- Jose Blanchet and Karthyek Murthy. Quantifying distributional model risk via optimal transport.
   *Mathematics of Operations Research*, 44(2):565–600, 2019.
- Jose Blanchet, Miao Lu, Tong Zhang, and Han Zhong. Double pessimism is provably efficient for distributionally robust offline reinforcement learning: Generic algorithm and robust partial coverage. *Advances in Neural Information Processing Systems*, 36, 2024.
- Zilong Cao, Pan Zhou, Ruixuan Li, Siqi Huang, and Dapeng Wu. Multiagent deep reinforcement learning for joint multichannel access and task offloading of mobile-edge computing in industry 4.0. *IEEE Internet of Things Journal*, 7(7):6201–6213, 2020.
- Yuxin Chen, Yuejie Chi, Jianqing Fan, Cong Ma, et al. Spectral methods for data science: A statistical perspective. *Foundations and Trends*® *in Machine Learning*, 14(5):566–806, 2021.
- Zixiang Chen, Dongruo Zhou, and Quanquan Gu. Almost optimal algorithms for two-player zero sum linear mixture Markov games. In *International Conference on Algorithmic Learning Theory*,
   pp. 227–261. PMLR, 2022.
- Pierre Clavier, Erwan Le Pennec, and Matthieu Geist. Towards minimax optimality of model-based robust reinforcement learning. *arXiv preprint arXiv:2302.05372*, 2023.
- Qiwen Cui, Kaiqing Zhang, and Simon Du. Breaking the curse of multiagents in a large state space:
  Rl in markov games with independent linear function approximation. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 2651–2652. PMLR, 2023.
- John C Duchi. Introductory lectures on stochastic optimization. *The mathematics of data*, 25: 99–186, 2018.
- Songtao Feng, Ming Yin, Yu-Xiang Wang, Jing Yang, and Yingbin Liang. Model-free algorithm with improved sample efficiency for zero-sum markov games, 2023.
- Rui Gao. Finite-sample guarantees for Wasserstein distributionally robust optimization: Breaking
   the curse of dimensionality. *Operations Research*, 71(6):2291–2306, 2023.
- Edgar N Gilbert. A comparison of signalling alphabets. *The Bell system technical journal*, 31(3): 504–522, 1952.
- Vineet Goyal and Julien Grand-Clement. Robust markov decision processes: Beyond rectangularity.
   *Mathematics of Operations Research*, 48(1):203–226, 2023.
- Garud N Iyengar. Robust dynamic programming. *Mathematics of Operations Research*, 30(2): 257–280, 2005.

607

609

627

631

640

- 594 Chi Jin, Qinghua Liu, Yuanhao Wang, and Tiancheng Yu. V-learning-a simple, efficient, 595 decentralized algorithm for multiagent RL. In ICLR 2022 Workshop on Gamification and 596 Multiagent Solutions, 2022. 597
- Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline RL? In 598 International Conference on Machine Learning, pp. 5084–5096. PMLR, 2021.
- 600 Shyam Sundar Kannan, Vishnunandan LN Venkatesh, and Byung-Cheol Min. Smart-Ilm: Smart 601 multi-agent robot task planning using large language models. arXiv preprint arXiv:2309.10062, 602 2023. 603
- Erim Kardeş, Fernando Ordóñez, and Randolph W Hall. Discounted robust stochastic games and an 604 application to queueing control. Operations research, 59(2):365–382, 2011. 605
- Nathan Lambert, Markus Wulfmeier, William Whitney, Arunkumar Byravan, Michael Bloesch, Vibhavari Dasagi, Tim Hertweck, and Martin Riedmiller. The challenges of exploration for offline 608 reinforcement learning. arXiv preprint arXiv:2201.11861, 2022.
- Chung-Wei Lee, Haipeng Luo, Chen-Yu Wei, and Mengxiao Zhang. Linear last-iterate convergence 610 for matrix games and stochastic games. arXiv preprint arXiv:2006.09517, 2020. 611
- 612 Gen Li, Laixi Shi, Yuxin Chen, Yuejie Chi, and Yuting Wei. Settling the sample complexity of 613 model-based offline reinforcement learning. The Annals of Statistics, 52(1):233–260, 2024a. 614
- Gen Li, Yuting Wei, Yuejie Chi, and Yuxin Chen. Breaking the sample size barrier in model-based 615 reinforcement learning with a generative model. Operations Research, 72(1):203-221, 2024b. 616
- 617 Na Li, Yuchen Jiao, Hangguan Shan, and Shefeng Yan. Provable memory efficient self-play 618 algorithm for model-free reinforcement learning. In The Twelfth International Conference on 619 Learning Representations, 2024c. 620
- Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In 621 Machine learning proceedings 1994, pp. 157–163. Elsevier, 1994. 622
- 623 Qinghua Liu, Tiancheng Yu, Yu Bai, and Chi Jin. A sharp analysis of model-based reinforcement 624 learning with self-play. In Marina Meila and Tong Zhang (eds.), Proceedings of the 38th 625 International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning 626 Research, pp. 7001–7010. PMLR, July 2021.
- Zhishuai Liu and Pan Xu. Distributionally robust off-dynamics reinforcement learning: Provable 628 efficiency with linear function approximation. In International Conference on Artificial 629 Intelligence and Statistics, pp. 2719–2727. PMLR, 2024. 630
- Shaocong Ma, Ziyi Chen, Shaofeng Zou, and Yi Zhou. Decentralized robust v-learning for solving 632 markov games with model uncertainty. Journal of Machine Learning Research, 24(371):1-40, 633 2023. 634
- Afshin Oroojlooy and Davood Hajinezhad. A review of cooperative multi-agent deep reinforcement 635 learning. Applied Intelligence, 53(11):13677–13722, 2023. 636
- 637 Tabish Rashid, Gregory Farquhar, Bei Peng, and Shimon Whiteson. Weighted qmix: Expanding 638 monotonic value function factorisation for deep multi-agent reinforcement learning. Advances in 639 neural information processing systems, 33:10199–10210, 2020.
- Lloyd S Shapley. Stochastic games. Proceedings of the National Academy of Sciences, 39(10): 641 1095–1100, 1953. 642
- 643 Laixi Shi and Yuejie Chi. Distributionally robust model-based offline reinforcement learning with 644 near-optimal sample complexity. Journal of Machine Learning Research, 25(200):1–91, 2024. 645
- Laixi Shi, Gen Li, Yuting Wei, Yuxin Chen, Matthieu Geist, and Yuejie Chi. The curious price of 646 distributional robustness in reinforcement learning with a generative model. Advances in Neural 647 Information Processing Systems, 36, 2024a.

- 648 Laixi Shi, Eric Mazumdar, Yuejie Chi, and Adam Wierman. Sample-efficient robust multi-agent 649 reinforcement learning in the face of environmental uncertainty. In Forty-first International 650 Conference on Machine Learning, 2024b. 651 David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, 652 Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go 653 without human knowledge. Nature, 550(7676):354-359, 2017. 654 655 Alexandre B. Tsybakov. Introduction to Nonparametric Estimation. Springer Publishing Company, 656 Incorporated, 1st edition, 2008. ISBN 0387790519. 657 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, 658 volume 47. Cambridge university press, 2018. 659 Daniel Vial, Sanjay Shakkottai, and R Srikant. Robust multi-agent bandits over undirected graphs. 661 Proceedings of the ACM on Measurement and Analysis of Computing Systems, 6(3):1–57, 2022. 662 Qiaomin Xie, Yudong Chen, Zhaoran Wang, and Zhuoran Yang. Learning zero-sum simultaneous-663 move Markov games using function approximation and correlated equilibrium. Math. Oper. Res., 664 48(1):433–462, Jun. 2022. ISSN 0364-765X. doi: 10.1287/moor.2022.1268. 665 666 Zaiyan Xu, Kishan Panaganti, and Dileep Kalathil. Improved sample complexity bounds for distributionally robust reinforcement learning. In International Conference on Artificial 667 Intelligence and Statistics, pp. 9728–9754. PMLR, 2023. 668 669 Yuling Yan, Gen Li, Yuxin Chen, and Jianqing Fan. Model-based reinforcement learning for offline 670 zero-sum markov games. Operations Research, 2024. 671 Wenhao Yang, Liangyu Zhang, and Zhihua Zhang. Toward theoretical understandings of robust 672 markov decision processes: Sample complexity and asymptotics. The Annals of Statistics, 50(6): 673 3223-3248, 2022. 674 675 Wenhao Yang, Han Wang, Tadashi Kozuno, Scott M Jordan, and Zhihua Zhang. Avoiding 676 model estimation in robust markov decision processes with a generative model. arXiv preprint 677 arXiv:2302.01248, 5, 2023. 678 Christopher Yeh, Chenlin Meng, Sherrie Wang, Anne Driscoll, Erik Rozi, Patrick Liu, Jihyeon Lee, 679 Marshall Burke, David B Lobell, and Stefano Ermon. Sustainbench: Benchmarks for monitoring 680 the sustainable development goals with machine learning. arXiv preprint arXiv:2111.04724, 681 2021. 682 683 Resilience enhancement of multi-agent Lanting Zeng, Dawei Qiu, and Mingyang Sun. 684 reinforcement learning-based demand response against adversarial attacks. Applied Energy, 324: 685 119688, 2022. 686 Huan Zhang, Hongge Chen, Duane S Boning, and Cho-Jui Hsieh. Robust reinforcement learning 687 on state observations with learned optimal adversary. In International Conference on Learning 688 Representations, 2021. 689 690 Kaiqing Zhang, Tao Sun, Yunzhe Tao, Sahika Genc, Sunil Mallya, and Tamer Basar. Robust multiagent reinforcement learning with model uncertainty. Advances in neural information processing 691 systems, 33:10571-10583, 2020. 692 693 Ziyuan Zhou and Guanjun Liu. Robustness testing for multi-agent reinforcement learning: State 694 perturbations on critical agents. arXiv preprint arXiv:2306.06136, 2023. 696 697 699
- 700
- 701

#### PRELIMINARIES А

**Dual equivalence of robust Bellman.** We can compute the robust Bellman operator by solving its dual formulation rather than the original form, as long as the predefined uncertainty set is in a benign form (e.g., utilizing TV distance as the divergence function) (Iyengar, 2005; Shi et al., 2024a). Taking TV distance as an example, we describe the equivalence under strong duality between the robust Bellman operator and its dual form as Lemma 2.

**Lemma 2** Consider any TV uncertainty set  $\mathcal{U}^{\sigma^+}(P)$  and  $\mathcal{U}^{\sigma^-}(P)$  associated with fixed uncertainty levels  $\sigma^+$  and  $\sigma^- \in (0,1]$  and any probability vector  $P \in \Delta(S)$ , respectively. For any vector  $V \in \mathbb{R}^S$  obeying V > 0, one has

$$\inf_{P \in \mathcal{U}^{\sigma^+}(P)} PV = \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P[V]_{\alpha} - \sigma^+ \left( \alpha - \min_{s'} [V]_{\alpha} \left( s' \right) \right) \right\};$$
(30a)

$$\sup_{P \in \mathcal{U}^{\sigma^-}(P)} PV = \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P[V]_{\alpha} - \sigma^- \left(\alpha - \min_{s'} [V]_{\alpha} \left(s'\right)\right) \right\},$$
(30b)

where  $[V]_{\alpha}$  is defined in (19)

ŀ

The proof of Lemma 2 is similar to Iyengar (2005, Lemma 4.3). Therefore, comparing the standard Bellman operator, the lemma above guarantees that no more computation cost is required when applying the robust Bellman operator, ignoring some logarithmic factors (Iyengar, 2005).

Facts of RTZMGs and empirical RTZMGs. Recall the definition of any RTZMG  $MG_r$  =  $\{\mathcal{S}, \mathcal{A}, \mathcal{B}, \mathcal{U}_{\rho}^{\sigma^+}(P^0), \mathcal{U}_{\rho}^{\sigma^-}(P^0), r, H\}$ . According to robust Bellman equations in (8), one has: for any product policy  $(\mu, \nu)$  and any  $(h, s, a, b) \in [H] \times S \times A \times B$ ,

$$Q_{h}^{\mu,\nu,\sigma^{+}}(s,a,b) = r_{h}(s,a,b) + \inf_{P \in \mathcal{U}_{\rho}^{\sigma^{+}}(P_{h,s,a,b}^{0})} PV_{h+1}^{\mu,\nu,\sigma^{+}};$$
(31a)

$$Q_{h}^{\mu,\nu,\sigma^{-}}(s,a,b) = r_{h}(s,a,b) + \sup_{P \in \mathcal{U}_{\rho}^{\sigma^{-}}(P_{h,s,a,b}^{0})} PV_{h+1}^{\mu,\nu,\sigma^{-}},$$
(31b)

where

$$V_h^{\mu,\nu,\sigma^+}(s) = \mathbb{E}_{a \sim \mu_h(s), b \sim \nu_h(s)} \left[ Q_h^{\mu,\nu,\sigma^+}(s,a,b) \right];$$
$$V_h^{\mu,\nu,\sigma^-}(s) = \mathbb{E}_{a \sim \mu_h(s), b \sim \nu_h(s)} \left[ Q_h^{\mu,\nu,\sigma^-}(s,a,b) \right].$$

> Considering the offline setting, we use  $\widehat{\mathcal{MG}}_{\mathsf{r}} = \{\mathcal{S}, \mathcal{A}, \mathcal{B}, \mathcal{U}_{\rho}^{\sigma^+}(\widehat{P}^0), \mathcal{U}_{\rho}^{\sigma^-}(\widehat{P}^0), r, H\}$  to represent the empirical RTZMG, which is establishing along with the estimated nominal distribution  $\widehat{P}^0$ in (12). Therefore, for any product policy  $(\mu, \nu)$ , we define the empirical robust value function (resp. empirical robust Q-function) in  $\widehat{\mathcal{MG}}_r$  as  $\widehat{V}_h^{\mu,\nu,\sigma^+}$  and  $\widehat{V}_h^{\mu,\nu,\sigma^-}$  (resp.  $\widehat{Q}_h^{\mu,\nu,\sigma^+}$  and  $\widehat{Q}_h^{\mu,\nu,\sigma^-}$ ), which are analogous to (4). Moreover, we can similarly define the optimal of the empirical robust value function for both player over  $\widehat{\mathcal{M}}\widehat{\mathcal{G}}_r$ , which is: for  $\forall s \in \mathcal{S}$ ,

$$\widehat{V}_{h}^{\star,\nu,\sigma^{+}}(s) = \widehat{V}_{h}^{\mu^{\star},\nu,\sigma^{+}}(s) \coloneqq \max_{\mu:\mathcal{S}\times[H]\to\Delta(\mathcal{A})} \widehat{V}_{h}^{\mu,\nu,\sigma^{+}}(s) = \max_{\mu:\mathcal{S}\times[H]\to\Delta(\mathcal{A})} \inf_{P\in\mathcal{U}^{\sigma^{+}}(\widehat{P}^{0})} \widehat{V}_{h}^{\mu,\nu,P}(s);$$
(32a)

750  
751  
752  
753  

$$\widehat{V}_{h}^{\mu,\star,\sigma^{-}}(s) = \widehat{V}_{h}^{\mu,\nu^{\star},\sigma^{-}}(s) \coloneqq \max_{\nu:\mathcal{S}\times[H]\to\Delta(\mathcal{B})} \widehat{V}_{h}^{\mu,\nu,\sigma^{-}}(s) = \max_{\nu:\mathcal{S}\times[H]\to\Delta(\mathcal{B})} \inf_{P\in\mathcal{U}^{\sigma^{-}}(\widehat{P}^{0})} \widehat{V}_{h}^{\mu,\nu,P}(s).$$
(32b)

Notably, for all  $s \in S$ , there exists at least one *robust best-response* policy that can achieve  $\widehat{V}_{h}^{\star,\nu,\sigma^{+}}(s)$  and  $\widehat{V}_{h}^{\mu,\star,\sigma^{-}}(s)$ , as proved by Blanchet et al. (2024).

Therefore, we can obtain the empirical robust Bellman equation similar to (8) as: for any product policy  $(\mu, \nu)$ ,

761 762 763

769 770

771 772

773 774

775

776 777

778

779

785 786 787  $\widehat{Q}_{h}^{\mu,\nu,\sigma^{+}}(s,a,b) = r_{h}(s,a,b) + \inf_{P \in \mathcal{U}_{\rho}^{\sigma^{+}}(\widehat{P}_{h,s,a,b}^{0})} P\widehat{V}_{h+1}^{\mu,\nu,\sigma^{+}};$ (33a)

$$\widehat{Q}_{h}^{\mu,\nu,\sigma^{-}}(s,a,b) = r_{h}(s,a,b) + \sup_{P \in \mathcal{U}_{\rho}^{\sigma^{-}}(\widehat{P}_{h,s,a,b}^{0})} P\widehat{V}_{h+1}^{\mu,\nu,\sigma^{-}},$$
(33b)

where

$$\widehat{V}_{h}^{\mu,\nu,\sigma^{+}}(s) = \mathbb{E}_{a \sim \mu_{h}(s), b \sim \nu_{h}(s)}[\widehat{Q}_{h}^{\mu,\nu,\sigma^{+}}(s,a,b)];$$
$$\widehat{V}_{h}^{\mu,\nu,\sigma^{-}}(s) = \mathbb{E}_{a \sim \mu_{h}(s), b \sim \nu_{h}(s)}[\widehat{Q}_{h}^{\mu,\nu,\sigma^{-}}(s,a,b)].$$

### B PROOF OF THEOREM 1

The proof of Theorem 1 can be separated into three steps, as outlined below.

#### B.1 STEP 1: DECOUPLING STATISTICAL DEPENDENCY

Before bounding  $\operatorname{Gap}(\widehat{\mu}, \widehat{\nu})$ , we introduce an important lemma, quantifying the difference between  $\widehat{P}$  and P when projected in the direction of the value function.

**Lemma 3** Instate the assumptions in Theorem 1. Consider any vector  $V \in \mathbb{R}^S$  with  $||V||_{\infty} \leq H$ for all  $(h, s, a, b) \in [H] \times S \times A \times B$  satisfying  $N_h(s, a, b) > 0$ . With probability at least  $1 - \delta$ , one has

$$\left|\inf_{P\in\mathcal{U}^{\sigma^+}(\widehat{P}^0_{h,s,a,b})}PV - \inf_{P\in\mathcal{U}^{\sigma^+}(P^0_{h,s,a,b})}PV\right| \leq C_4\sqrt{\frac{1}{N_h(s,a,b)}}\mathsf{Var}_{\widehat{P}^0_{h,s,a,b}}(V)\log\frac{KH}{\delta} + C_4\frac{H\log\frac{KH}{\delta}}{N_h(s,a,b)}$$
(34)

for some sufficiently large constant  $C_4 > 0$ , and

$$\operatorname{Var}_{\widehat{P}_{h,s,a,b}^{0}}\left(V\right) \leq 2\operatorname{Var}_{P_{h,s,a,b}^{0}}\left(V\right) + O\left(\frac{H^{2}}{N_{h}\left(s,a,b\right)}\log\frac{KH}{\delta}\right).$$
(35)

Proof can be found in Appendix C.3.

In simple terms, (34) provides a Bernstein-type concentration bound, while (35) ensures that the empirical variance estimate (i.e., the plug-in estimate) closely matches the true variance. Notably, Lemma 3 does not require V to be statistically independent of  $\hat{P}^0_{h,s,a,b}$ , which is essential given the complex statistical dependencies in our iterative algorithm. Under the leave-one-out analysis (see, e.g., Agarwal et al. (2020); Chen et al. (2021); Li et al. (2024a;b)), we prove Lemma 3 to decouple statistical dependencies, as illustrated in Appendix C.3. With Lemma 3, we can now have

$$\left|\inf_{\mathcal{P}\in\mathcal{U}^{\sigma^+}(\widehat{P}^0_{h,s,a,b})}PV - \inf_{\mathcal{P}\in\mathcal{U}^{\sigma^+}(P^0_{h,s,a,b})}PV\right| \le \beta_h\left(s,a,b,V\right)$$
(36)

798 799

802 803

804 805

796 797

Therefore, we conclude that  $\widehat{Q}_{h}^{+}(s, a, b)$  is an optimistic estimation of  $\widehat{Q}_{h}^{\mu,\nu,\sigma^{+}}(s, a, b)$ , which is summarized below.

**Lemma 4** With probability exceeding  $1 - \delta$ , it holds that

for any  $(h, s, a, b) \in [H] \times S \times A \times B$  satisfying  $N_h(s, a, b) \geq 1$ .

$$\widehat{Q}_{h}^{+}(s,a,b) \ge Q_{h}^{\star,\widehat{\nu},\sigma^{+}}(s,a,b) \quad and \quad \widehat{V}_{h}^{+}(s) \ge V_{h}^{\star,\widehat{\nu},\sigma^{+}}(s);$$
(37)

806 807 See Appendix C.4 for detail proofs.

Besides, we introduce another key lemma highlighting the difference between RTZMGs and standard TZMGs from the same idea by Shi et al. (2024b, Lemma 3). The range of the robust value function narrows as the uncertainty level  $\sigma^+$  of its uncertainty set increases, as shown below. **Lemma 5** Consider the uncertainty set  $\mathcal{U}^{\sigma^+}(\cdot)$  with TV distance and any RTZMG  $\mathcal{MG}_r = \{\mathcal{S}, \mathcal{A}, \mathcal{B}, \mathcal{U}^{\sigma^+}(P), \mathcal{U}^{\sigma^-}(P), r, H\}$ . The optimistic robust value function estimate  $\widehat{V}_h^+$ :

$$\forall h \in [H]: \quad \max_{s \in \mathcal{S}} \widehat{V}_h^+ - \min_{s \in \mathcal{S}} \widehat{V}_h^+ \le \min\left\{\frac{(H+1)\left(1 - (1 - \sigma^+)^{H-h}\right)}{\sigma^+}, H\right\}$$

See Appendix C.5 for detail proofs.

 **B.2** Step 2: decomposing the error  $\text{Gap}(\hat{\mu}, \hat{\nu})$ 

821 The goal of our algorithm is to output an  $\varepsilon$ -robust NE policy  $(\hat{\mu}, \hat{\nu})$  satisfying  $\operatorname{Gap}(\hat{\mu}, \hat{\nu})$  in (11), 822 i.e.,

$$\operatorname{Gap}(\widehat{\mu},\widehat{\nu}) \coloneqq \max\left\{ V_1^{\star,\widehat{\nu},\sigma^+}(\varrho) - V_1^{\star,\sigma^+}(\varrho), \ V_1^{\star,\sigma^-}(\varrho) - V_1^{\widehat{\mu},\star,\sigma^-}(\varrho) \right\} \le \varepsilon.$$

Due to the symmetry between max-player and min-player, we assume without loss of generality that  $V_1^{\star,\hat{\nu},\sigma^+}(\varrho) - V_1^{\star,\sigma^+}(\varrho)$  is larger than  $V_1^{\star,\sigma^-}(\varrho) - V_1^{\hat{\mu},\star,\sigma^-}(\varrho)$ , leading to  $\operatorname{Gap}(\hat{\mu},\hat{\nu}) \leq \{V_1^{\star,\hat{\nu},\sigma^+}(\varrho) - V_1^{\star,\sigma^+}(\varrho)\}$ .

According to the relationship in Lemma 4, we obtain

$$V_{h}^{\star,\widehat{\nu},\sigma^{+}}(s) \leq \widehat{V}_{h}^{+}(s) = \max_{\mu \in \Delta(\mathcal{A})} \min_{\nu \in \Delta(\mathcal{B})} \mathbb{E}_{(a,b) \sim (\mu(s),\nu(s))} \left[ \widehat{Q}_{h}^{+}(s,a,b) \right]$$
$$\leq \max_{\mu \in \Delta(\mathcal{A})} \mathbb{E}_{(a,b) \sim (\mu(s),\nu^{\star}(s))} \left[ Q_{h}^{+}(s,a,b) \right], \tag{38}$$

where the first equality comes from line 6 in Algorithm 2. Therefore, there exists a deterministic policy  $\mu^{d} : S \leftarrow \Delta(A)$  satisfying that for any  $s \in S$ 

$$\mu^{\mathsf{d}}(s) \coloneqq \arg \max_{\mu \in \Delta(\mathcal{A})} \mathbb{E}_{(a,b) \sim (\mu(s),\nu^{\star}(s))} \left[ Q_h^+(s,a,b) \right].$$
(39)

Before starting, we introduce several useful notations:

 The state-action space covered by the behavior policy (μ<sup>n</sup>, ν<sup>n</sup>) in the nominal transition kernel P<sup>0</sup> is denoted as

$$\mathcal{C}^{\mathsf{n}} = \{(h, s, a, b) : d_h^{\mathsf{n}}(s, a, b) > 0\}.$$
(40)

The set of potential state occupancy distributions w.r.t. the policy (μ<sup>d</sup>(s), ν\*(s)) in a model within the uncertainty set P ∈ U<sup>σ+</sup> (P<sup>0</sup>) for any time step h ∈ [H] is denoted as

$$\mathcal{D}_{h}^{\mathbf{p}} \coloneqq \left\{ \left[ d_{h}^{\mu^{\mathsf{d}}(s),\nu^{\star}(s),P}(s) \right]_{s\in\mathcal{S}} : P \in \mathcal{U}^{\sigma^{+}}\left(P^{0}\right) \right\};$$

$$(41)$$

$$\mathcal{D}_{h}^{\mathsf{pa}} \coloneqq \left\{ \left[ d_{h}^{\mu^{\mathsf{d}}(s),\nu^{\star}(s),P}(s,a,b) \right]_{(s,a,b)\in\mathcal{S}\times\mathcal{A}\times\mathcal{B}} : P \in \mathcal{U}^{\sigma^{+}}\left(P^{0}\right) \right\}.$$
(42)

• For convenience and without ambiguity, we introduce an additional notation for  $h \in [H]$  as

$$\beta_h^{\mu^{\mathsf{d}},\nu^{\star}}(s) = \mathbb{E}_{(a,b)\sim(\mu^{\mathsf{d}}(s),\nu^{\star}(s))}\beta_h\left(s,a,b,\widehat{V}_{h+1}^+\right)$$

In particular, the vector  $\beta_h^{\mu^{\mathsf{d}},\nu^{\star}} \in \mathbb{R}^S$  is defined with its *s*-th item given by  $\beta_h^{\mu^{\mathsf{d}},\nu^{\star}}(s)$ .

• Similarly, we can define the notation related to rewards for  $h \in [H]$  as

$$\widehat{r}_{h}^{\mu^{\mathsf{d}},\nu^{\star}}(s) = \mathbb{E}_{(a,b)\sim(\mu^{\mathsf{d}}(s),\nu^{\star}(s))}\widehat{r}_{h}\left(s,a,b\right)$$

According to the update rule in line 4 in Algorithm 2 and robust Bellman equality (31), we derive  $V_{h}^{\star,\hat{\nu},\sigma^{+}}(s) - V_{h}^{\star,\sigma^{+}}(s)$   $\leq \widehat{V}_{h}^{+}(s) - V_{h}^{\mu^{d},\nu^{\star},\sigma^{+}}(s)$   $\leq \mathbb{E}_{(a,b)\sim(\mu^{d}(s),\nu^{\star}(s))} \inf_{P \in \mathcal{U}^{\sigma^{+}}(\hat{P}_{h,s,a,b}^{0})} P \widehat{V}_{h+1}^{+} + \beta_{h}^{\mu^{d},\nu^{\star}}(s)$   $- \mathbb{E}_{(a,b)\sim(\mu^{d}(s),\nu^{\star}(s))} \left[ \inf_{P \in \mathcal{U}^{\sigma^{+}}(P_{h,s,a,b}^{0})} P \widehat{V}_{h+1}^{+} - \inf_{P \in \mathcal{U}^{\sigma^{+}}(P_{h,s,a,b}^{0})} P V_{h+1}^{\mu^{d},\nu^{\star},\sigma^{+}} + \left| \inf_{P \in \mathcal{U}^{\sigma^{+}}(P_{h,s,a,b}^{0})} P \widehat{V}_{h+1}^{+} - \inf_{P \in \mathcal{U}^{\sigma^{+}}(P_{h,s,a,b}^{0})} P \widehat{V}_{h+1}^{+} \right| \right] + \beta_{h}^{\mu^{d},\nu^{\star},\sigma^{+}} \left| \sup_{e \in \mathcal{U}^{\sigma^{+}}(P_{h,s,a,b}^{0})} P \widehat{V}_{h+1}^{+} - \inf_{P \in \mathcal{U}^{\sigma^{+}}(P_{h,s,a,b}^{0})} P V_{h+1}^{\mu^{d},\nu^{\star},\sigma^{+}} \right| + 2\beta_{h}^{\mu^{d},\nu^{\star}}(s)$ 

Here, (ii) is valid under the notation

$$P_{h,s,a,b}^{\inf,V} \coloneqq \operatorname{argmin}_{P \in \mathcal{U}^{\sigma^+} \left(P_{h,s,a,b}^0\right)} P V_{h+1}^{\mu^{\mathsf{d}},\nu^{\star},\sigma^+}$$
(44)

(43)

(46)

and consequently,

$$\inf_{P \in \mathcal{U}^{\sigma^+}(P_{h,s,a,b}^0)} PV_{h+1}^{\mu^{d},\nu^{\star},\sigma^+} = P_{h,s,a,b}^{\inf,V} V_{h+1}^{\mu^{d},\nu^{\star},\sigma^+}, \text{ and } \inf_{P \in \mathcal{U}^{\sigma^+}(P_{h,s,a,b}^0)} P\widehat{V}_{h+1}^+ \le P_{h,s,a,b}^{\inf,V} \widehat{V}_{h+1}^+.$$

Besides, (i) in (43) exists due to (36) in Lemma 3 for  $N_h(s, a, b) > 0$  and

 $\overset{(\mathrm{ii})}{\leq} \mathbb{E}_{(a,b)\sim(\mu^{\mathsf{d}}(s),\nu^{\star}(s))} \left[ P_{h,s,a,b}^{\inf,V} \left( \widehat{V}_{h+1}^{+} - V_{h+1}^{\mu^{\mathsf{d}},\nu^{\star},\sigma^{+}} \right) \right] + 2\beta_{h}^{\mu^{\mathsf{d}},\nu^{\star}}(s).$ 

$$\left| \inf_{P \in \mathcal{U}^{\sigma^+}(P^0_{h,s,a,b})} P \widehat{V}^+_{h+1} - \inf_{P \in \mathcal{U}^{\sigma^+}(\widehat{P}^0_{h,s,a,b})} P \widehat{V}^+_{h+1} \right| \le H = \beta_h^{\mu^d,\nu^\star}(s)$$
(45)

for  $N_h(s, a, b) = 0$ .

For ease of proof, we introduce a notation as  $\widetilde{P}_{h,s}^{\inf,V} := \mathbb{E}_{(a,b)\sim(\mu^d(s),\nu^\star(s))}P_{h,s,a,b}^{\inf,V}$ . Furthermore, we define a sequence of matrices  $\widetilde{P}_h^{\inf,V} \in \mathbb{R}^{S\times S}$ . We can utilizing (43) recursively over the time steps  $h, h+1, \cdots, H$  and derive

$$V_{h}^{\star,\hat{\nu},\sigma^{+}}(s) - V_{h}^{\star,\sigma^{+}}(s) \leq \widehat{V}_{h}^{+}(s) - V_{h}^{\mu^{d},\nu^{\star},\sigma^{+}}(s) \\ \leq \widetilde{P}_{h}^{\inf,V}\left(\widehat{V}_{h+1}^{+} - V_{h+1}^{\mu^{d},\nu^{\star},\sigma^{+}}\right) + 2\beta_{h}^{\mu^{d},\nu^{\star}}(s) \\ \leq \widetilde{P}_{h}^{\inf,V}\widetilde{P}_{h+1}^{\inf,V}\left(\widehat{V}_{h+2}^{+} - V_{h+2}^{\mu^{d},\nu^{\star},\sigma^{+}}\right) + 2\widetilde{P}_{h}^{\inf,V}\beta_{h+1}^{\mu^{d},\nu^{\star}} + 2\beta_{h}^{\mu^{d},\nu^{\star}}(s)$$

 $\leq \cdots \leq 2 \sum_{i=h}^{H} \left( \prod_{i=h}^{i-1} \widetilde{P}_{j}^{\inf, V} \right) \beta_{i}^{\mu^{\mathsf{d}}, \nu^{\star}},$ 

where we define  $\left(\prod_{j=h}^{i-1} \widetilde{P}_{j}^{\inf,V}\right) = I$  for convenience.

For any  $d_h^{\mu^d,\nu^\star} \in \mathcal{D}_h^p$  (cf. (41)), taking inner product with (46) yields 

$$\left\langle d_{h}^{\mu^{\mathsf{d}},\nu^{\star}}, V_{h}^{\star,\widehat{\nu},\sigma^{+}}(s) - V_{h}^{\star,\sigma^{+}}(s) \right\rangle \leq \left\langle d_{h}^{\mu^{\mathsf{d}},\nu^{\star}}, 2\sum_{i=h}^{H} \left( \prod_{j=h}^{i-1} \widetilde{P}_{j}^{\mathrm{inf},V} \right) \beta_{i}^{\mu^{\mathsf{d}},\nu^{\star}} \right\rangle$$
$$= 2\sum_{i=h}^{H} \left\langle d_{i}^{\mathsf{p},\mu^{\mathsf{d}},\nu^{\star}}, \beta_{i}^{\mu^{\mathsf{d}},\nu^{\star}} \right\rangle, \tag{47}$$

where

$$d_{i}^{\mathbf{p},\mu^{\mathsf{d}},\nu^{\star}} \coloneqq \left[ \left( d_{h}^{\mu^{\mathsf{d}},\nu^{\star}} \right)^{\top} \left( \prod_{j=h}^{i-1} \widetilde{P}_{j}^{\mathrm{inf},V} \right) \right]^{\top} \in \mathcal{D}_{i}^{\mathsf{p}}$$
(48)

by the definition of  $\mathcal{D}_i^p$  (cf. (41)) for all  $i = h + 1, \dots, H$ . 

Next, we control  $\langle d_i^{\mathbf{p},\mu^{\mathbf{d}},\nu^{\star}}, \beta_i^{\mu^{\mathbf{d}},\nu^{\star}} \rangle$  utilizing concentrability. First of all, according to (20) in Lemma 3, we demonstrate that the pessimistic penalty satisfies

$$\beta_{i}(s, a, b, \hat{V}) \leq \max\left\{ \sqrt{\frac{C_{n} \log \frac{KH}{\delta}}{N_{i}(s, a, b)}} \operatorname{Var}_{\hat{P}_{i, s, a, b}^{0}}(\hat{V}), \frac{2C_{n}H \log \frac{KH}{\delta}}{N_{i}(s, a, b)} \right\}$$

$$\leq \sqrt{\frac{C_{n} \log \frac{KH}{\delta}}{N_{i}(s, a, b)}} \operatorname{Var}_{\hat{P}_{i, s, a, b}^{0}}(\hat{V}) + \frac{2C_{n}H \log \frac{KH}{\delta}}{N_{i}(s, a, b)}$$

$$\stackrel{(i)}{\leq} \sqrt{\frac{C_{n} \log \frac{KH}{\delta}}{N_{i}(s, a, b)}} \left( 2\operatorname{Var}_{P_{i, s, a, b}^{0}}(\hat{V}) + \frac{C_{0}H^{2}}{N_{i}(s, a, b)} \log \frac{KH}{\delta} \right)} + \frac{2C_{n}H \log \frac{KH}{\delta}}{N_{i}(s, a, b)}$$

$$\stackrel{(ii)}{\leq} \sqrt{\frac{2C_{n} \log \frac{KH}{\delta}}{N_{i}(s, a, b)}} \operatorname{Var}_{P_{i, s, a, b}^{0}}(\hat{V}) + \frac{\left(2C_{n} + \sqrt{C_{n}C_{0}}\right) H \log \frac{KH}{\delta}}{N_{i}(s, a, b)}}$$

$$(49)$$

where (i) holds by applying (35) for some sufficiently large  $C_0$  and (ii) exists follows from the Cauchy-Schwarz inequality. Therefore, combining the definition of  $\beta_i^{\mu^{\alpha},\nu^{\star}}(s)$ , we obtain

$$\langle d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}, \beta_{i}^{\mu^{d},\nu^{\star}} \rangle = \sum_{s \in \mathcal{S}} d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}(s) \beta_{i}^{\mu^{d},\nu^{\star}}(s)$$

$$= \sum_{s \in \mathcal{S}} d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}(s) \mathbb{E}_{(a,b)\sim(\mu^{d}(s),\nu^{\star}(s))} \beta_{i}(s,a,b,\hat{V})$$

$$= \sum_{(s,a,b)\in\mathcal{S}\times\mathcal{A}\times\mathcal{B}} d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}(s) \mathbb{1}\{a = \mu^{d}(s)\}\nu^{\star}(b|s)\beta_{i}(s,a,b,\hat{V})$$

$$= \sum_{(s,b)\in\mathcal{S}\times\mathcal{B}} d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}(s,\mu^{d}(s),b)\beta_{i}(s,\mu^{d}(s),b,\hat{V}),$$
(50)

where the last equation holds due to the definition in (7b). Then, we observe  $d_{b}^{\mathsf{p},\mu^{\mathsf{d}},\nu^{\star}}(s,a,b) \in \mathcal{D}_{b}^{\mathsf{pa}}$ (cf. (42)). Thereafter, we divide the bound (50) into two cases.

For the first case, i.e.,  $s \in S$  where  $\max_{P \in \mathcal{U}^{\sigma^+}(P^0)} d_i^{\mu^{\mathsf{d}},\nu^{\star},P}(s,\mu^{\mathsf{d}}(s),b) = 0$ , it follows from the definition (cf. (41)) that for any  $d_i^{\mathsf{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b) \in \mathcal{D}_i^{\mathsf{pa}}$ , it satisfies that  $d_{z}^{\mathbf{p},\mu^{\mathsf{d}},\nu^{\star}}$ 

$$\sum_{i}^{p,\mu^{a},\nu^{*}}(s,\mu^{d}(s),b) = 0.$$
(51)

For the second case, i.e.,  $s \in S$  where  $\max_{P \in \mathcal{U}^{\sigma^+}(P^0)} d_i^{\mu^{\mathsf{d}},\nu^{\star},P}(s,\mu^{\mathsf{d}}(s),b) > 0$ , by the assumption in (22) 

970  
971 
$$\max_{P \in \mathcal{U}^{\sigma^+}(P^0)} \frac{\min\left\{d_i^{\mu^{\mathsf{d}},\nu^{\star},P}\left(s,\mu^{\mathsf{d}}(s),b\right),\frac{1}{S(A+B)}\right\}}{d_i^{\mathsf{n}}\left(s,\mu^{\mathsf{d}}(s),b\right)} \le C_{\mathsf{r}}^{\star} < \infty.$$

972 It implies that

$$d_i^{\mathsf{n}}(s, \mu^{\mathsf{d}}(s), b) > 0 \quad \text{and} \quad \left(i, s, \mu^{\mathsf{d}}(s), b\right) \in \mathcal{C}^{\mathsf{n}}.$$
(52)

Lemma 1 tells that with probability at least  $1 - 8\delta$ ,

$$N_{i}(s, \mu^{\mathsf{d}}(s), b) \geq \frac{Kd_{i}^{\mathsf{n}}(s, \mu^{\mathsf{d}}(s), b)}{8} - 5\sqrt{Kd_{i}^{\mathsf{n}}(s, \mu^{\mathsf{d}}(s), b)\log\frac{KH}{\delta}}$$

$$\stackrel{(i)}{\geq} \frac{Kd_{i}^{\mathsf{n}}(s, \mu^{\mathsf{d}}(s), b)}{16}$$

$$\stackrel{(ii)}{\geq} \frac{K\max_{P\in\mathcal{U}^{\sigma}(P^{0})}\min\left\{d_{i}^{\mu^{\mathsf{d}},\nu^{\star},P}\left(s, \mu^{\mathsf{d}}(s), b\right), \frac{1}{S(A+B)}\right\}}{16C_{\mathsf{r}}^{\star}}$$

$$\geq \frac{K\min\left\{d_{i}^{\mathsf{p},\mu^{\mathsf{d}},\nu^{\star}}\left(s, \mu^{\mathsf{d}}(s), b\right), \frac{1}{S(A+B)}\right\}}{16C_{\mathsf{r}}^{\star}}, \tag{53}$$

where (ii) comes from Assumption 1 and (i) holds due to

$$Kd_{i}^{\mathsf{n}}(s,\mu^{\mathsf{d}}(s),b) \geq c_{0}\frac{HS(A+B)}{d_{\mathsf{m}}^{\mathsf{n}}}\log\frac{KH}{\delta}f(\sigma^{+},\sigma^{-},H)d_{i}^{\mathsf{n}}(s,\mu^{\mathsf{d}}(s),b)$$
$$\geq c_{0}HS(A+B)\log\frac{KH}{\delta}f(\sigma^{+},\sigma^{-},H) \geq 1600\log\frac{KH}{\delta},\tag{54}$$

where  $f(\sigma^+, \sigma^-, H) = \min\left\{\frac{H\sigma^+ + 1 - (1 - \sigma^+)^H}{(\sigma^+)^2}, \frac{H\sigma^- + 1 - (1 - \sigma^-)^H}{(\sigma^-)^2}, H\right\}$ , the first inequality follows from condition (24), and the second inequality follows from

$$d_{\mathsf{m}}^{\mathsf{n}} = \min_{h,s,\mu^{\mathsf{d}}(s),b} \left\{ d_{h}^{\mathsf{n}}(s,\mu^{\mathsf{d}}(s),b) : d_{h}^{\mathsf{n}}(s,\mu^{\mathsf{d}}(s),b) > 0 \right\} \le d_{i}^{\mathsf{n}}\left(s,\mu^{\mathsf{d}}(s),b\right).$$
(55)

1002 Combining the results in (49) and (50), we arrive at

$$\begin{aligned}
&\left\langle d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}},\beta_{i}^{\mu^{d},\nu^{\star}}\right\rangle \\
&= \sum_{(s,b)\in\mathcal{S}\times\mathcal{B}} d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}(s,\mu^{d}(s),b)\beta_{i}(s,\mu^{d}(s),b,\hat{V}) \\
&\leq \sum_{(s,b)\in\mathcal{S}\times\mathcal{B}} d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}(s,\mu^{d}(s),b)\sqrt{\frac{2C_{n}\log\frac{KH}{\delta}}{N_{i}(s,\mu^{d}(s),b)}} \mathsf{Var}_{P_{i,s,\mu^{d}(s),b}^{0}}(\hat{V}) \\
&+ \sum_{(s,b)\in\mathcal{S}\times\mathcal{B}} d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}(s,\mu^{d}(s),b)\frac{\left(2C_{n}+\sqrt{C_{n}C_{0}}\right)H\log\frac{KH}{\delta}}{N_{i}(s,\mu^{d}(s),b)} \\
&\stackrel{(\mathrm{i})}{\leq} \sum_{(s,b)\in\mathcal{S}\times\mathcal{B}} d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}(s,\mu^{d}(s),b)\sqrt{\frac{32C_{r}^{\star}C_{n}\log\frac{KH}{\delta}}{K\min\left\{d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}(s,\mu^{d}(s),b),\frac{1}{S(A+B)}\right\}}}\mathsf{Var}_{P_{i,s,\mu^{d}(s),b}^{0}}(\hat{V}) \\
&+ \sum_{(s,b)\in\mathcal{S}\times\mathcal{B}} d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}(s,\mu^{d}(s),b)\frac{16C_{r}^{\star}\left(2C_{n}+\sqrt{C_{n}C_{0}}\right)H\log\frac{KH}{\delta}}{K\min\left\{d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}(s,\mu^{d}(s),b),\frac{1}{S(A+B)}\right\}}}.
\end{aligned}$$
(56)

1025 Therefore, according to (47), we just need to bound  $\sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_i^{p,\mu^d,\nu^*}(s,\mu^d(s),b)B_1$  and  $\sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_i^{p,\mu^d,\nu^*}(s,\mu^d(s),b)B_2$ , which is introduced as follows. 
$$\begin{aligned} & \text{Part 1: Bounding } \sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b)B_{1} \quad \text{Combining the result in (54) with} \\ & \sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b)B_{1} \quad \text{yields} \end{aligned} \\ & \sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b)B_{1} \\ & = \sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b)\sqrt{\frac{32C_{i}^{*}C_{n}\log\frac{KH}{\delta}}{K\min\left\{d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b\right\}} \text{Var}_{P_{i,s,\mu^{d}}(s),b}(\widehat{V})} \\ & \sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b) \sqrt{\frac{32C_{i}^{*}C_{n}\log\frac{KH}{\delta}}{K\min\left\{d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b\right\}} \text{Var}_{P_{i,s,\mu^{d}}(s),b}(\widehat{V})} \\ & \max\left\{\sqrt{\frac{32C_{i}^{*}C_{n}\log\frac{KH}{\delta}}{Kd_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b)} \text{Var}_{P_{i,s,\mu^{d}}(s),b}(\widehat{V})}, \sqrt{\frac{32C_{i}^{*}C_{n}S(A+B)\log\frac{KH}{\delta}}{K}} \text{Var}_{P_{i,s,\mu^{d}}(s),b}(\widehat{V})} \right\} \\ & \leq \sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b) \sqrt{\frac{32C_{i}^{*}C_{n}S(A+B)\log\frac{KH}{\delta}}{K}} \text{Var}_{P_{i,s,\mu^{d}}(s),b}(\widehat{V})} \\ & + \sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b) \sqrt{\frac{32C_{i}^{*}C_{n}S(A+B)\log\frac{KH}{\delta}}{K}} \text{Var}_{P_{i,s,\mu^{d}}(s),b}(\widehat{V})} \\ & \leq \sqrt{\frac{32C_{i}^{*}C_{n}S(A+B)\log\frac{KH}{\delta}}{K}} \left( \sqrt{H\sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}}(s,\mu^{d}(s),b) \text{Var}_{P_{i,s,\mu^{d}}(s),b}(\widehat{V})} \\ & + \sqrt{\sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b) \text{Var}_{P_{i,s,\mu^{d}}(s),b}(\widehat{V})} \times \sqrt{\sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b)}} \\ & = \sqrt{\frac{128C_{i}^{*}C_{n}HS(A+B)\log\frac{KH}{\delta}}{K}} \sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b) \text{Var}_{P_{i,s,\mu^{d}}(s),b}(\widehat{V})} \\ & = \sqrt{\frac{128C_{i}^{*}C_{n}HS(A+B)\log\frac{KH}{\delta}}{K}} \sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b) \text{Var}_{P_{i,s,\mu^{d}}(s),b}(\widehat{V})} \\ & = \sqrt{\frac{128C_{i}^{*}C_{n}HS(A+B)\log\frac{KH}{\delta}}{K}} \sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b) \text{Var}_{P_{i,s,\mu^{d}}(s),b}(\widehat{V})} \\ & = \sqrt{\frac{128C_{i}^{*}C_{n}HS(A+B)\log\frac{KH}{\delta}}{K}} \sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_{i}^{p,\mu^{d},\nu^{*}}(s,\mu^{d}(s),b)}$$

where the last inequality follows from the Cauchy-Schwarz inequality. Then, we introduce the following lemma about  $\sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_i^{\mathsf{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b) \operatorname{Var}_{P_{i,s,\mu^{\mathsf{d}}(s),b}}(\widehat{V})$ , whose proof is postponed to Appendix C.6. 

**Lemma 6** Considering  $\forall \delta \in (0, 1)$ , with probability at least  $1 - \delta$ , one has: for any product policy  $(\widehat{\mu}, \widehat{\nu}),$ 

Armed with Lemma 6, (57) can be further bounded as 

 $\sum_{i=1}^{n} \sum_{\substack{(a,b) \in S \times B}} d_i^{\mathfrak{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b) B_1$ 

 $\sum_{i=1}^{n} \sum_{\substack{(s,b) \in S \times B}} d_i^{\mathbf{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b) B_2$ 

 $\stackrel{(i)}{\leq} \frac{32C_{\mathsf{r}}^{\star} \left(2C_{\mathsf{n}} + \sqrt{C_{\mathsf{n}}C_{3}}\right) H^{2}S(A+B) \log \frac{KH}{\delta}}{K},$ 

 $\sum \frac{d_i^{\mathbf{p},\mu^{\mathbf{d}},\nu^{\star}}(s,\mu^{\mathbf{d}}(s),b)}{\frac{1}{2} \left(\frac{1}{2} \right)\right)\right)\right)\right)}\right)\right)\right)}\right)}$ 

**Part 2: Bounding**  $\sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_i^{\mathbf{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b)B_2$  Combining the result in (53) with  $\sum_{i=1}^{H} \sum_{(s,b) \in \mathcal{S} \times \mathcal{B}} \overline{d_i^{\mathsf{p},\mu^\mathsf{d},\nu^\star}}(s,\mu^\mathsf{d}(s),b) B_2 \text{ yields}$ 

 $=\sum_{i=1}^{H}\sum_{(s,b)\in\mathcal{S}\times\mathcal{B}}d_{i}^{\mathsf{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b)\frac{16C_{\mathsf{r}}^{\star}\left(2C_{\mathsf{n}}+\sqrt{C_{\mathsf{n}}C_{3}}\right)H\log\frac{KH}{\delta}}{K\min\left\{d_{i}^{\mathsf{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b),\frac{1}{S(A+B)}\right\}}$ 

 $\leq \sqrt{\frac{128 C_{\mathsf{r}}^{\star} C_{\mathsf{n}} HS(A+B) \log \frac{KH}{\delta}}{K}} \sqrt{H \min\left\{\frac{2(H\sigma^{+}-1+(1-\sigma^{+})^{H})}{(\sigma^{+})^{2}}, H\right\}}$ 

 $\times \sqrt{\left(4\sum_{i=1}^{H}\sum_{(s,b)\in\mathcal{S}\times\mathcal{B}}d_{i}^{\mathsf{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b)\beta_{i}(s,\mu^{\mathsf{d}}(s),b,\widehat{V})+(H+3)\right)}.$ 

(59)

(60)

where the inequality holds by the trivial fact

$$\leq \sum_{(s,b)\in\mathcal{S}\times\mathcal{B}} d_i^{\mathsf{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b) \left(\frac{1}{d_i^{\mathsf{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b)} + \frac{1}{1/S(A+B)}\right)$$

$$= \sum_{(s,b)\in\mathcal{S}\times\mathcal{B}} 1 + S(A+B) \sum_{(s,\mu^{\mathsf{d}},\mu^{\mathsf{d}},\nu^{\star})} d_i^{\mathsf{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b) \leq 2S(A+B)$$

$$=\sum_{(s,b)\in\mathcal{S}\times\mathcal{B}}1+S(A+B)\sum_{(s,b)\in\mathcal{S}\times\mathcal{B}}d_i^{\mathbf{p},\mu^{\mathbf{d}},\nu^{\star}}(s,\mu^{\mathbf{d}}(s),b)\leq 2S(A+B).$$
(61)

**Putting all together** Combining the results (59) and (60) in Part 1 and Part 2, we obtain

$$\sum_{i=1}^{H} \sum_{(s,b)\in S\times B} d_{i}^{\mathfrak{p},\mu^{\mathfrak{d}},\nu^{\star}}(s,\mu^{\mathfrak{d}}(s),b)\beta_{i}(s,\mu^{\mathfrak{d}}(s),b,\widehat{V}) \\
\leq \sqrt{\frac{128C_{\mathsf{r}}^{\star}C_{\mathsf{n}}H^{2}S(A+B)\log\frac{KH}{\delta}}{K}} \min\left\{\frac{2(H\sigma^{+}-1+(1-\sigma^{+})^{H})}{(\sigma^{+})^{2}},H\right\}}{\sqrt{\left(4\sum_{i=1}^{H} \sum_{(s,b)\in S\times B} d_{i}^{\mathfrak{p},\mu^{\mathfrak{d}},\nu^{\star}}(s,\mu^{\mathfrak{d}}(s),b)\beta_{i}(s,\mu^{\mathfrak{d}}(s),b,\widehat{V})+(H+3)\right)}} \\
+ \frac{32C_{\mathsf{r}}^{\star}\left(2C_{\mathsf{n}}+\sqrt{C_{\mathsf{n}}C_{3}}\right)H^{2}S(A+B)\log\frac{KH}{\delta}}{K}}{K}, \qquad (62)$$

$$\begin{aligned} & \text{which can further bound as} \\ & \text{int} \\ & \text{int} \\ & \sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_i^{p,\mu^d,\nu^*}(s,\mu^d(s),b)\beta_i(s,\mu^d(s),b,\hat{V}) \\ & \text{int} \\ & \frac{128C_r^*C_nH^2(H+3)S(A+B)\log\frac{KH}{\delta}}{K}\min\left\{\frac{2(H\sigma^+-1+(1-\sigma^+)^H)}{(\sigma^+)^2},H\right\} \\ & + \frac{32C_r^*\left(2C_n+\sqrt{C_nC_3}\right)H^2S(A+B)\log\frac{KH}{\delta}}{K} + \sqrt{\frac{512C_r^*C_nH^2S(A+B)\log\frac{KH}{\delta}}{K}} \\ & \times \sqrt{\min\left\{\frac{2(H\sigma^+-1+(1-\sigma^+)^H)}{(\sigma^+)^2},H\right\}}\sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_i^{p,\mu^d,\nu^*}(s,\mu^d(s),b)\beta_i(s,\mu^d(s),b,\hat{V})} \\ & \frac{\sqrt{\frac{128C_r^*C_nH^2(H+3)S(A+B)\log\frac{KH}{\delta}}{K}}}{K}\min\left\{\frac{2(H\sigma^+-1+(1-\sigma^+)^H)}{(\sigma^+)^2},H\right\}} \\ & + \frac{32C_r^*\left(2C_n+\sqrt{C_nC_3}\right)H^2S(A+B)\log\frac{KH}{\delta}}{K}}{K} \\ & + \frac{256C_r^*C_nH^2S(A+B)\log\frac{KH}{\delta}}{K}}{K}\min\left\{\frac{2(H\sigma^+-1+(1-\sigma^+)^H)}{(\sigma^+)^2},H\right\}} \\ & + \frac{1}{2}\sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_i^{p,\mu^d,\nu^*}(s,\mu^d(s),b)\beta_i(s,\mu^d(s),b,\hat{V}), \\ & \text{there the last ration follows from the AM GM inequality. Rearranging terms, it follows that } \end{aligned}$$

where the last relation follows from the AM-GM inequality. Rearranging terms, it follows that

$$\begin{split} &\sum_{i=1}^{H} \sum_{(s,b)\in\mathcal{S}\times\mathcal{B}} d_{i}^{\mathbf{p},\mu^{\mathbf{d}},\nu^{\star}}(s,\mu^{\mathbf{d}}(s),b)\beta_{i}(s,\mu^{\mathbf{d}}(s),b,\widehat{V}) \\ &\leq \sqrt{\frac{512C_{\mathbf{r}}^{\star}C_{\mathbf{n}}H^{2}(H+3)S(A+B)\log\frac{KH}{\delta}}{K}} \min\left\{\frac{2(H\sigma^{+}-1+(1-\sigma^{+})^{H})}{(\sigma^{+})^{2}},H\right\} \\ &+ \frac{64C_{\mathbf{r}}^{\star}\left(2C_{\mathbf{n}}+\sqrt{C_{\mathbf{n}}}C_{3}\right)H^{2}S(A+B)\log\frac{KH}{\delta}}{K}}{K} \\ &+ \frac{512C_{\mathbf{r}}^{\star}C_{\mathbf{n}}H^{2}S(A+B)\log\frac{KH}{\delta}}{K}\min\left\{\frac{2(H\sigma^{+}-1+(1-\sigma^{+})^{H})}{(\sigma^{+})^{2}},H\right\} \\ &\leq \sqrt{\frac{512C_{\mathbf{r}}^{\star}C_{\mathbf{n}}H^{2}(H+3)S(A+B)\log\frac{KH}{\delta}}{K}}\min\left\{\frac{2(H\sigma^{+}-1+(1-\sigma^{+})^{H})}{(\sigma^{+})^{2}},H\right\}}{H^{2}} \\ &+ \frac{C_{\mathbf{r}}^{\star}C_{2}H^{2}S(A+B)\log\frac{KH}{\delta}}{K}\min\left\{\frac{2(H\sigma^{+}-1+(1-\sigma^{+})^{H})}{(\sigma^{+})^{2}},H\right\}. \end{split}$$
(64)

Along with the above result, we are ready to bound  $V_1^{\star,\sigma^+}(\varrho) - V_1^{\widehat{\mu},\star,\sigma^+}(\varrho)$ . There exists some sufficiently large constants  $C_1, C_2, C_3 > 0$ , and 

$$\begin{aligned} & 1179 \\ & 1179 \\ & 1180 \\ & 1181 \\ & 1182 \\ & 1181 \\ & 1182 \\ & 1182 \\ & 1183 \\ & 1182 \\ & 1183 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & 1184 \\ & 1185 \\ & & 1184 \\ & 1185 \\ & & \\ &$$

where the last inequality follows from condition (24).

# 1188 B.3 STEP 3: SUMMING UP THE RESULTS

1190 Consequently, we obtain the upper bound of  $V_1^{\star,\hat{\nu},\sigma^+}(\varrho) - V_1^{\hat{\mu},\hat{\nu},\sigma^+}(\varrho)$  in (65). Similarly, 1191

$$V_{1}^{\star,\sigma^{-}}(\varrho) - V_{1}^{\widehat{\mu},\star,\sigma^{-}}(\varrho) \\ \leq \sqrt{\frac{C_{r}^{\star}C_{3}H^{2}S(A+B)\log\frac{KH}{\delta}}{K}} \min\left\{\frac{(H+1)(H\sigma^{-}-1+(1-\sigma^{-})^{H})}{(\sigma^{-})^{2}},H\right\}},$$
 (66)

1196 which directly leads to

1197 1198 1199

1201 1202

1205 1206

1219

1221

1223

1228

1233 1234

1192 1193 1194

1195

$$\begin{aligned}
\operatorname{Gap}(\widehat{\mu},\widehat{\nu}) &\leq c_1 \sqrt{\frac{C_{\mathsf{r}}^{\star} H^2 S(A+B) \log \frac{KH}{\delta}}{K}} \\
&\times \sqrt{\min\left\{\frac{2(H\sigma^+ - 1 + (1-\sigma^+)^H)}{(\sigma^+)^2}, \frac{2(H\sigma^- - 1 + (1-\sigma^-)^H)}{(\sigma^-)^2}, H\right\}}, \quad (67)
\end{aligned}$$

1203 for some sufficiently large  $c_1$  and

$$K \ge HS(A+B)\log\frac{KH}{\delta}\min\left\{\frac{2(H\sigma^+ - 1 + (1-\sigma^+)^H)}{(\sigma^+)^2}, \frac{2(H\sigma^- - 1 + (1-\sigma^-)^H)}{(\sigma^-)^2}, H\right\}.$$

1207 **Discussion of (67).** For the term  $T = \min(f(\sigma^+, \sigma^-), H)$ , considering the symmetry between 1208  $\sigma^+$  and  $\sigma^-$ , we define  $g(\sigma^+, H) = H\sigma^+ - H(1 - \sigma^+)^H - (\sigma^+)^2 H$ . For  $H \ge 2$ , we derive the first derivative as  $\frac{\partial g(\sigma^+, H)}{\partial \sigma^+} = H + H^2 (1 - \sigma^+)^{H-1} - 2H\sigma^+$ . Further, the second derivative is 1209 1210 given by  $\frac{\partial^2 g(\sigma^+, H)}{\partial (\sigma^+)^2} = -H^2(H-1)(1-\sigma^+)^{H-2} - 2H < 0$ , indicating that  $g(\sigma^+, H)$  is concave. 1211 By evaluating the first derivative at the boundaries, we find  $\frac{\partial g(\sigma^+, H)}{\partial \sigma^+}|_{\sigma^+ \to 0} \to H^2 + H > 0$ 1212 1213 and  $\frac{\partial g(\sigma^+, H)}{\partial \sigma^+}|_{\sigma^+=1} = -H < 0$ , which shows that  $g(\sigma^+, H)$  first increases monotonically, 1214 reaches a maximum at some point  $\sigma^*$ , and then decreases monotonically. Furthermore, since 1215  $\begin{array}{l} g(\sigma^{+} \to 0, H) \to -H < 0 \text{ and } g(\sigma^{+} = 1, H) = 0, \text{ there exists } 0 < \sigma^{0} < 1 \text{ such that } g(\sigma^{0}, H) = 0. \text{ Thus, when } \sigma^{0} \lesssim \min\{\sigma^{+}, \sigma^{-}\} \lesssim 1, \text{ we have } T = H. \text{ Otherwise, } T = \min\left\{\frac{(H\sigma^{+} - 1 + (1 - \sigma^{+})^{H})}{(\sigma^{+})^{2}}, \frac{(H\sigma^{-} - 1 + (1 - \sigma^{-})^{H})}{(\sigma^{-})^{2}}\right\}. \end{array}$ 1216 1217 1218

#### 1220 C AUXILIARY LEMMAS FOR THEOREM 1

1222 C.1 PROOF OF LEMMA 1

In this part, we prove Lemma 1 produced in Algorithm 1.

1225 Before next proof, we clarify the independent property. Let us examine two distinct data-generation 1226 mechanisms, where a sample transition quadruple (s, a, b, h, s') represents a transition from state s 1227 with actions (a, b) to state s' at step h.

**Step 1:** Augmenting  $\mathcal{D}^{t}$  to Create  $\mathcal{D}^{t,a}$ . To construct the augmented dataset  $\mathcal{D}^{t,a}$ , for each  $(s,h) \in \mathcal{S} \times [H]$ , we proceed as follows: (i). Include in  $\mathcal{D}^{t,a}$  all  $N_{h}^{t}(s)$  sample transitions in  $\mathcal{D}^{t}$  originating from state *s* at step *h*. (ii). If  $N_{h}^{t}(s) > N_{h}^{m}(s)$ , supplement  $\mathcal{D}^{t,a}$  with an additional  $N_{h}^{t}(s) - N_{h}^{m}(s)$  independent sample transitions  $\{(s, a_{h,s}^{(i)}, b_{h,s}^{(i)}, h, s_{h,s}^{'(i)})\}$ , generated as follows:

$$a_{h,s}^{(i)} \stackrel{\text{i.i.d.}}{\sim} \mu_{h}^{\mathbf{b}}(\cdot|s), \quad b_{h,s}^{(i)} \stackrel{\text{i.i.d.}}{\sim} \nu_{h}^{\mathbf{b}}(\cdot|s), \quad s_{h,s}^{\prime(i)} \stackrel{\text{i.i.d.}}{\sim} P_{h}\big(\cdot|s, a_{h,s}^{(i)}, b_{h,s}^{(i)}\big), \quad N_{h}^{\mathbf{m}}(s) < i \le N_{h}^{\mathbf{t}}(s).$$

1235 1236 Step 2: Constructing  $\mathcal{D}^{\text{iid}}$ . For each  $(s,h) \in \mathcal{S} \times [H]$ , generate  $N_h^t(s)$  independent sample 1237 transitions  $\{(s, a_{h,s}^{(i)}, b_{h,s}^{(i)}, h, s'_{h,s}^{(i)})\}$  as follows:

1241

$$a_{h,s}^{(i) \text{ i.i.d. }} \overset{\text{i.i.d. }}{\sim} \mu_h^{\text{b}}(\cdot|s), \quad b_{h,s}^{(i) \text{ i.i.d. }} \overset{\text{i.i.d. }}{\sim} \nu_h^{\text{b}}(\cdot|s), \quad s_{h,s}^{\prime(i) \text{ i.i.d. }} \overset{\text{i.i.d. }}{\sim} P_h\big(\cdot|s,a,b\big), \quad 1 \le i \le N_h^{\text{t}}(s)$$

1240 The resulting dataset is defined as:

$$\mathcal{D}^{\text{iid}} \coloneqq \Big\{ \left( s, a_{h,s}^{(i)}, b_{h,s}^{(i)}, h, s_{h,s}^{\prime(i)} \right) \mid s \in \mathcal{S}, 1 \le h \le H, 1 \le i \le N_h^{\text{t}}(s) \Big\}.$$

**Establishing independent property.** The dataset  $\mathcal{D}^{t,a}$  deviates from  $\mathcal{D}^{t}$  only when  $N_{h}^{t}(s) > N_{h}^{m}(s)$  holds. This augmentation ensures that  $\mathcal{D}^{t,a}$  contains precisely  $N_{h}^{t}(s)$  sample transitions from state s at step h. Both  $\mathcal{D}^{t,a}$  and  $\mathcal{D}^{iid}$  comprise exactly  $N_h^t(s)$  sample transitions from state s at step h, with  $\{N_h^t(s)\}$  being statistically independent of the randomness in sample generation. Consequently, given  $\{N_h^t(s)\}$ , the sample transitions in  $\mathcal{D}^{t,a}$  across different steps are statistically independent. As a result, both  $\mathcal{D}^t$  and  $\mathcal{D}^{iid}$  can be regarded as collections of independent samples. 

Next, we begin to prove  $N_h^t(s) \leq N_h^m(s)$ . Since  $\mathcal{D}^a$  is generated by half of the sample trajectories in line 2 in Algorithm 1, there is 

$$N_h^{\mathsf{a}}(s) = \sum_{k=K/2+1}^{K} \mathbbm{1}\left\{s_h^k = s\right\}$$

for each  $s \in S$  and  $1 \le h \le H$ . Thus, we can view  $N_h^a(s)$  as the sum of K/2 independent Bernoulli random variables with mean  $d_h^{\mu^n,\nu^n}(s)$ . According to the Bernstein inequality and the union bound, we derive 

$$\mathbb{P}\left\{\exists (s,h) \in \mathcal{S} \times [H] : \left| N_h^{\mathsf{a}}(s) - \frac{K}{2} d_h^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s) \right| \ge N_0 \right\}$$
$$\le \sum_{s \in \mathcal{S}, h \in [H]} \mathbb{P}\left\{ \left| N_h^{\mathsf{a}}(s) - \frac{K}{2} d_h^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s) \right| \ge N_0 \right\}$$

$$\leq \sum_{s \in \mathcal{S}, h \in [H]} \mathbb{P}^{\cdot}$$

$$\leq 2HS \exp\left(-\frac{N_0^2/2}{N_{h,s}+N_0/3}\right), \quad \forall N_0 \ge 0,$$

where

$$N_{h,s} \coloneqq \frac{K}{2} \mathsf{Var} \left( \mathbbm{1}\{s_h^t = s\} \right) = \frac{K d_h^{\mu^n,\nu^n}(s) \left( 1 - d_h^{\mu^n,\nu^n}(s) \right)}{2} \le \frac{K d_h^{\mu^n,\nu^n}(s)}{2}$$

Therefore, with probability at least  $1 - 2\delta$ , we yield that:  $\forall s \in S$  and  $\forall 1 \le h \le H$ , 

$$\begin{aligned} & | 1269 \\ 1270 \\ 1271 \\ 1272 \end{aligned} \quad \left| N_h^{a}(s) - \frac{K}{2} d_h^{\mu^{n},\nu^{n}}(s) \right| \le \sqrt{4N_{h,s} \log \frac{HS}{\delta}} + \frac{2}{3} \log \frac{HS}{\delta} \le \sqrt{2K d_h^{\mu^{n},\nu^{n}}(s) \log \frac{HS}{\delta}} + \log \frac{HS}{\delta}. \end{aligned}$$
(68)

As generated by the same way between  $\mathcal{D}^{m}$  and  $\mathcal{D}^{a}$ , we similarly obtain that with probability exceeding  $1 - 2\delta$ ,  $\forall s \in S$  and  $\forall 1 \leq h \leq H$ , 

$$\left|N_{h}^{\mathsf{m}}(s) - \frac{K}{2} d_{h}^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s)\right| \leq \sqrt{2K d_{h}^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s) \log \frac{HS}{\delta}} + \log \frac{HS}{\delta}.$$
(69)

Combining (68) and (69), there is 

$$|N_h^{\mathsf{m}}(s) - N_h^{\mathsf{a}}(s)| \le 2\sqrt{2Kd_h^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s)\log\frac{HS}{\delta}} + 2\log\frac{HS}{\delta}$$
(70)

for all  $s \in S$  and  $1 \le h \le H$ . 

Now, we complete the proof of  $N_h^t(s) \le N_h^m(s)$ , which can be divided into two cases. 

The first case is  $N_h^a(s) \leq 100 \log \frac{HS}{\delta}$ . According to the definition in (15), we obtain 

$$N_h^{\mathsf{t}}(s) = \max\left\{N_h^{\mathsf{a}}(s) - 10\sqrt{N_h^{\mathsf{a}}(s)\log\frac{HS}{\delta}}, 0\right\} = 0 \le N_h^{\mathsf{m}}(s).$$
(71)

The second case is  $N_h^a(s) > 100 \log \frac{HS}{\delta}$ . Followed by (68), we obtain 

$$\frac{K}{2}d_h^{\mu^{\mathfrak{n}},\nu^{\mathfrak{n}}}(s) + \sqrt{2Kd_h^{\mu^{\mathfrak{n}},\nu^{\mathfrak{n}}}(s)\log\frac{HS}{\delta}} + \log\frac{HS}{\delta} \ge N_h^{\mathfrak{a}}(s),$$

leading to 

$$Kd_h^{\mu^n,\nu^n}(s) \ge (9\sqrt{2})^2 \log \frac{HS}{\delta} \ge 100 \log \frac{HS}{\delta}.$$
(72)

1296 Thus, we take (72) back to (68) and derive

$$N_{h}^{a}(s) \geq \frac{K}{2} d_{h}^{\mu^{n},\nu^{n}}(s) - \sqrt{2K d_{h}^{\mu^{n},\nu^{n}}(s) \log \frac{HS}{\delta}} - \log \frac{HS}{\delta} \geq \frac{K}{4} d_{h}^{\mu^{n},\nu^{n}}(s).$$
(73)

1300 Consequently, in the case of  $N_h^a(s) > 100 \log \frac{HS}{\delta}$ , we have

$$N_{h}^{\mathsf{t}}(s) = \max\left\{N_{h}^{\mathsf{a}}(s) - 10\sqrt{N_{h}^{\mathsf{a}}(s)\log\frac{HS}{\delta}}, 0\right\}$$

$$= N_h^{\mathsf{a}}(s) - 10\sqrt{N_h^{\mathsf{a}}(s)\log\frac{HS}{\delta}}$$

$$\stackrel{(\mathrm{i})}{\leq} N_h^{\mathrm{a}}(s) - 5\sqrt{Kd_h^{\mu^{\mathrm{n}},\nu^{\mathrm{n}}}(s)\log\frac{HS}{\delta}}$$

$$\stackrel{\text{(ii)}}{\leq} N_h^{\mathsf{a}}(s) - \left\{ 2\sqrt{2Kd_h^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s)\log\frac{HS}{\delta}} + 2\log\frac{HS}{\delta} \right\} \stackrel{\text{(iii)}}{\leq} N_h^{\mathsf{m}}(s), \tag{74}$$

1311 (1312) where (i) holds under condition (73), (ii) exists under the condition (72), and (iii) comes from the inequality (70) with probability at least  $1 - 2\delta$ .

1314 Combining the results in (71) and (74) together, we establish  $N_h^t(s) \le N_h^m(s)$ .

Now, we claim the following bound with proof in Appendix C.2:  $\forall (s, a, b, h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{B} \times [H]$ , with probability exceeding  $1 - 2\delta$ ,

$$N_{h}^{\mathsf{t}}(s,a,b) \ge N_{h}^{\mathsf{t}}(s)\mu_{h}^{\mathsf{n}}(a\,|\,s)\nu_{h}^{\mathsf{n}}(b\,|\,s) - \sqrt{4N_{h}^{\mathsf{t}}(s)\mu_{h}^{\mathsf{n}}(a\,|\,s)\nu_{h}^{\mathsf{n}}(b\,|\,s)\log\frac{KH}{\delta} - \log\frac{KH}{\delta}}.$$
 (75)

Armed with the fact  $N_h^t(s) \le N_h^m(s)$  and claim (75), we start to prove (16). In the following, we discuss two cases, i.e.,  $Kd_h^{\mu^n,\nu^n}(s,a,b) \le 1600 \log \frac{KH}{\delta}$  and  $Kd_h^{\mu^n,\nu^n}(s,a,b) > 1600 \log \frac{KH}{\delta}$ .

For the first case of  $Kd_h^{\mu^n,\nu^n}(s,a) \le 1600 \log \frac{KH}{\delta}$ , we can easily classified that

$$\frac{K}{8}d_{h}^{\mu^{n},\nu^{n}}(s,a) - 5\sqrt{Kd_{h}^{\mu^{n},\nu^{n}}(s,a)\log\frac{KH}{\delta}} \le 0 \le N_{h}^{t}(s,a).$$
(76)

For the second case of  $Kd_h^{\mu^n,\nu^n}(s,a,b) = Kd_h^{\mu^n,\nu^n}(s)\mu_h^n(a \mid s)\nu_h^n(b \mid s) > 1600 \log \frac{KH}{\delta}$ , we obtain KH

$$N_h^{\mathsf{a}}(s) \ge \frac{K}{4} d_h^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s) \ge 400 \log \frac{KH}{\delta},$$
(77)

which is derived by the same line of (73) with slight modification. The property in (77) and the definition of  $N_h^t(s)$  together yield

$$N_h^{\mathsf{t}}(s) \ge N_h^{\mathsf{a}}(s) - 10\sqrt{N_h^{\mathsf{a}}(s)\log\frac{KH}{\delta}}$$
$$> \frac{K_{d^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s)}}{10\sqrt{K_{d^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s)\log\frac{KH}{\delta}}}$$

$$\geq \frac{K}{4} d_h^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s) - 10\sqrt{\frac{K}{4}} d_h^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s) \log \frac{KH}{\delta} \geq \frac{K}{8} d_h^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s)$$

1338 As a consequent,

$$N_{h}^{\mathsf{t}}(s)\mu_{h}^{\mathsf{n}}(a\,|\,s)\nu_{h}^{\mathsf{n}}(b\,|\,s) \ge \frac{K}{8}d_{h}^{\mu^{\mathsf{n}},\nu^{\mathsf{n}}}(s)\mu_{h}^{\mathsf{n}}(a\,|\,s)\nu_{h}^{\mathsf{n}}(b\,|\,s)$$
(78)

$$=\frac{K}{8}d_{h}^{\mu^{n},\nu^{n}}(s,a,b) \ge 200\log\frac{KH}{\delta},$$
(79)

where the last inequality holds due to the assumption of the second case. Taking the lower bound (78) with (75) together, there is

1346  
1347 
$$N_{h}^{t}(s, a, b) \geq \frac{K}{8} d_{h}^{\mu^{n}, \nu^{n}}(s, a, b) - \sqrt{\frac{K}{2}} d_{h}^{\mu^{n}, \nu^{n}}(s, a, b) \log \frac{KH}{\delta} - \log \frac{KH}{\delta}$$

1348  
1349 
$$\geq \frac{K}{8} d_h^{\mu^n,\nu^n}(s,a,b) - 2\sqrt{K d_h^{\mu^n,\nu^n}(s,a,b) \log \frac{KH}{\delta}}.$$

Putting the result above and (76) together, according to the claim (75), we can finally complete the proof of Lemma 1.

1353 1354

1355

C.2 PROOF OF CLAIM (75).

To prove claim (75), we analyze two cases, i.e.,  $N_h^t(s)\mu_h^n(a \mid s)\nu_h^n(b \mid s) \leq 4\log\frac{KH}{\delta}$  and  $N_h^t(s)\mu_h^n(a \mid s)\nu_h^n(b \mid s) > 4\log\frac{KH}{\delta}$ .

For the first case of  $N_h^t(s)\mu_h^n(a \mid s)\nu_h^n(b \mid s) \le 4\log\frac{KH}{\delta}$ , we conclude the right-hand side of (75) is negative, leading to the claim (75).

For the second case of  $N_h^t(s)\mu_h^n(a \mid s)\nu_h^n(b \mid s) > 4\log\frac{KH}{\delta}$ , we compose a special set  $\mathcal{D}^l$  as

$$\mathcal{D}^{\mathsf{I}} \coloneqq \left\{ (s, a, b, h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{B} \times [H] \mid N_h^{\mathsf{t}}(s) \mu_h^{\mathsf{n}}(a \mid s) \nu_h^{\mathsf{n}}(b \mid s) > 4 \log \frac{KH}{\delta} \right\}.$$
(80)

With the fact of

$$\begin{split} \sum_{(s,a,b,h)\in\mathcal{S}\times\mathcal{A}\times\mathcal{B}\times[H]} N_h^{\mathsf{t}}(s)\mu_h^{\mathsf{n}}(a\,|\,s)\nu_h^{\mathsf{n}}(b\,|\,s) &= \sum_{(s,h)\in\mathcal{S}\times[H]} N_h^{\mathsf{t}}(s) \sum_{(a,b)\in\mathcal{A}\times\mathcal{B}} \mu_h^{\mathsf{n}}(a\,|\,s)\nu_h^{\mathsf{n}}(b\,|\,s) \\ &= \sum_{(s,h)\in\mathcal{S}\times[H]} N_h^{\mathsf{t}}(s) \leq \sum_{(s,h)\in\mathcal{S}\times[H]} N_h^{\mathsf{a}}(s) = \frac{KH}{2}, \end{split}$$

1376 1377 the cardinality of  $\mathcal{D}^{\mathsf{I}}$  can be bounded as:

$$\left|\mathcal{D}^{\mathsf{I}}\right| < \frac{\sum_{(s,a,b,h)} N_h^{\mathsf{t}}(s)\mu_h^{\mathsf{n}}(a\,|\,s)\nu_h^{\mathsf{n}}(b\,|\,s)}{4\log\frac{KH}{\delta}} \le KH/2.$$
(81)

Besides, we can view  $N_h^t(s, a)$  as the sum of  $N_h^t(s)$  independent Bernoulli random variables with mean  $\mu_h^n(a \mid s)\nu_h^n(b \mid s)$ , which holds due to  $N_h^t(s) \leq N_h^m(s)$  with high probability and condition on  $N_h^t(s)$ ,  $N_h^m(s)$ . Analogous to (68) based on the condition  $N_h^t(s) \leq N_h^m(s)$ , we can repeat the Bernstein-type argument and obtain that for any fixed triple (s, a, b, h), with probability at least  $1-2\delta/(KH)$ ,

1390 1391

1392 1393

$$N_{h}^{\mathsf{t}}(s,a,b) \ge N_{h}^{\mathsf{t}}(s)\mu_{h}^{\mathsf{n}}(a\,|\,s)\nu_{h}^{\mathsf{n}}(b\,|\,s) - \sqrt{4N_{h}^{\mathsf{t}}(s)\mu_{h}^{\mathsf{n}}(a\,|\,s)\nu_{h}^{\mathsf{n}}(b\,|\,s)\log\frac{KH}{\delta}} - \log\frac{KH}{\delta}.$$
(82)

Therefore, with probability exceeding  $1 - \delta$ , (82) holds for all  $(s, a, b, h) \in \mathcal{D}^{\mathsf{I}}$  by utilizing the union bound of (81) over all  $(s, a, b, h) \in \mathcal{D}^{\mathsf{I}}$ .

Consequently, combining the results above under two cases, we derive that the property (75) holds for all  $(s, a, b, h) \in S \times A \times B \times [H]$  with probability at least  $1 - \delta$ .

1399 1400

1401 C.3 PROOF OF LEMMA 3

- 1402
- 1403 We prove Lemma 3 similar to the proof of claim 1 by Yan et al. (2024), which is separated into two parts as follows.

Т

Part 1: proof of inequality (34). According to the definition in (18), for any fixed value vector Vindependent from  $\hat{P}_{h,s,a,b}^0$ , we have

$$\begin{vmatrix} \inf_{P \in \mathcal{U}^{\sigma^{+}}(\hat{P}_{h,s,a,b}^{0})} PV - \inf_{P \in \mathcal{U}^{\sigma^{+}}(P_{h,s,a,b}^{0})} PV \end{vmatrix}$$
$$= \begin{vmatrix} \max_{\alpha \in [\min_{s} V(s), \max_{s} V(s)]} \left\{ \hat{P}_{h,s,a,b}^{0} [V]_{\alpha} - \sigma^{+} \left(\alpha - \min_{s'} [V]_{\alpha} (s')\right) \right\} \\ - \max_{\alpha \in [\min_{s} V(s), \max_{s} V(s)]} \left\{ P_{h,s,a,b}^{0} [V]_{\alpha} - \sigma^{+} \left(\alpha - \min_{s'} [V]_{\alpha} (s')\right) \right\} \end{vmatrix}$$
$$\leq \max_{\alpha \in [\min_{s} V(s), \max_{s} V(s)]} \left| \hat{P}_{h,s,a,b}^{0} [V]_{\alpha} - P_{h,s,a,b}^{0} [V]_{\alpha} \end{vmatrix}$$
$$\leq \max_{\alpha \in [0,H]} \left| \hat{P}_{h,s,a,b}^{0} [V]_{\alpha} - P_{h,s,a,b}^{0} [V]_{\alpha} \right|,$$
(83)

Т

where the last inequality exists due to the fact that the maximum operator is 1-Lipschitz.

According to the definition of empirical transition kernel  $\widehat{P}^0_{h,s,a,b}$ , we get

$$(\widehat{P}_{h,s,a,b}^{0} - P_{h,s,a,b}^{0}) [V]_{\alpha}$$

$$= \sum_{s' \in \mathcal{S}} \underbrace{ [V(s')]_{\alpha} \left[ \frac{\sum_{i=1}^{N} \mathbb{1} \{h_i = h, s_i = s, a_i = a, b_i = b, s'_i = s'\}}{N_h(s, a, b)} - P_h^0(s' \mid s, a, b) \right]}_{=:X_{s'}}$$

as a sum of independent random variables. Based on the relationship between  $P_{h,s,a,b}^0$  and  $\hat{P}_{h,s,a,b}^0$ , we verify  $\mathbb{E}[X_{s'}] = 0$  and  $|X_{s'}| \leq H$  for all  $s' \in S$ . Therefore, with probability exceeding  $1 - \delta$ and for some universal constant  $C_4 > 0$ , under the Bernstein inequality (Vershynin, 2018, Theorem 2.8.4), we have

$$\left( \widehat{P}_{h,s,a,b}^{0} - P_{h,s,a,b}^{0} \right) [V]_{\alpha} \leq C_{4} \sqrt{\frac{1}{N_{h}\left(s,a,b\right)}} \mathsf{Var}_{P_{h,s,a,b}^{0}}\left([V]_{\alpha}\right) \log \frac{KH}{\delta} + \frac{C_{4}H \log \frac{KH}{\delta}}{N_{h}\left(s,a,b\right)}$$

$$\leq C_{4} \sqrt{\frac{1}{N_{h}\left(s,a,b\right)}} \mathsf{Var}_{P_{h,s,a,b}^{0}}\left(V\right) \log \frac{KH}{\delta} + \frac{C_{4}H \log \frac{KH}{\delta}}{N_{h}\left(s,a,b\right)}, \quad (84)$$

where the last inequality comes from the definition of  $[V]_{\alpha}$  in (19).

$$\begin{aligned} \mathsf{Var}_{P^{0}_{h,s,a,b}}\left(V\right) &= P^{0}_{h,s,a,b}\left(\overline{V}\circ\overline{V}\right) \\ &= \widehat{P}^{0}_{h,s,a,b}\left(\overline{V}\circ\overline{V}\right) + \left(P^{0}_{h,s,a,b} - \widehat{P}^{0}_{h,s,a,b}\right)\left(\overline{V}\circ\overline{V}\right) \\ &= \mathsf{Var}_{\widehat{P}^{0}_{h,s,a,b}}\left(V\right) + \left[\left(P^{0}_{h,s,a,b} - \widehat{P}^{0}_{h,s,a,b}\right)V\right]^{2} + \left(P^{0}_{h,s,a,b} - \widehat{P}^{0}_{h,s,a,b}\right)\left(\overline{V}\circ\overline{V}\right), \end{aligned} \tag{85}$$

#### where the last equation holds since

Let  $\overline{V} := V - \left(P_{h,s,a,b}^0 V\right) 1$ , we have

$$\begin{aligned} \widehat{P}_{h,s,a,b}^{0}(\overline{V} \circ \overline{V}) &= \widehat{P}_{h,s,a,b}^{0}\left(\left[V - \left(P_{h,s,a,b}^{0}V\right)1\right] \circ \left[V - \left(P_{h,s,a,b}^{0}V\right)1\right]\right) \\ &= \widehat{P}_{h,s,a,b}^{0}\left(V \circ V\right) - 2\left(P_{h,s,a,b}^{0}V\right)\left(\widehat{P}_{h,s,a,b}^{0}V\right) + \left(P_{h,s,a,b}^{0}V\right)^{2} \\ &= \widehat{P}_{h,s,a,b}^{0}\left(\left[V - \left(\widehat{P}_{h,s,a,b}^{0}V\right)1\right] \circ \left[V - \left(\widehat{P}_{h,s,a,b}^{0}V\right)1\right]\right) + \left(\widehat{P}_{h,s,a,b}^{0}V\right) \end{aligned}$$

1453  
1454 
$$= \widehat{P}_{h,s,a,b}^{0} \left( \left[ V - \left( \widehat{P}_{h,s,a,b}^{0} V \right) 1 \right] \circ \left[ V - \left( \widehat{P}_{h,s,a,b}^{0} V \right) 1 \right] \right) + \left( \widehat{P}_{h,s,a,b}^{0} V \right)^{2}$$
1455

1455  
1456 
$$-2\left(P_{h,s,a,b}^{0}V\right)\left(\widehat{P}_{h,s,a,b}^{0}V\right) + \left(P_{h,s,a,b}^{0}V\right)^{2}$$

1457 
$$= \operatorname{Var}_{\widehat{P}^{0}_{h,s,a,b}}(V) + \left[ \left( P^{0}_{h,s,a,b} - \widehat{P}^{0}_{h,s,a,b} \right) V \right]^{2}.$$

Analogous to (84), with probability exceeding  $1 - \delta$ , there is 

$$\begin{aligned} &|(\widehat{P}_{h,s,a,b}^{0} - P_{h,s,a,b}^{0})(\overline{V} \circ \overline{V})| \leq C_{4}\sqrt{\frac{1}{N_{h}(s,a,b)}}\mathsf{Var}_{P_{h,s,a,b}^{0}}(\overline{V} \circ \overline{V})\log\frac{KH}{\delta} + \frac{C_{4}H^{2}\log\frac{KH}{\delta}}{N_{h}(s,a,b)} \\ &|\mathsf{1462}|\\ &|\mathsf{1463}|\\ &|\mathsf{1464}|\\ &|\mathsf{1464}|\\ &|\mathsf{1465}|\\ &|\mathsf{1465}|\\ &|\mathsf{1466}|\\ &|\mathsf{1466}|\\$$

where the last inequation comes from the fact that

$$\mathsf{Var}_{P^0_{h,s,a,b}}\left(\overline{V}\circ\overline{V}\right) \leq P^0_{h,s,a,b}\left(\overline{V}\circ\overline{V}\circ\overline{V}\circ\overline{V}\right) \leq H^2 P^0_{h,s,a,b}\left(\overline{V}\circ\overline{V}\right) = H^2 \mathsf{Var}_{P^0_{h,s,a,b}}\left(V\right).$$

Under the result in (86), we bound (85) further as: 

$$\begin{split} \mathsf{Var}_{P_{h,s,a,b}^{0}}\left(V\right) \leq &\mathsf{Var}_{\widehat{P}_{h,s,a,b}^{0}}\left(V\right) + \left[\left(P_{h,s,a,b}^{0} - \widehat{P}_{h,s,a,b}^{0}\right)V\right]^{2} \\ &+ C_{4}\sqrt{\frac{H^{2}\log\frac{KH}{\delta}}{N_{h}\left(s,a,b\right)}}\mathsf{Var}_{P_{h,s,a,b}^{0}}\left(V\right)} + \frac{C_{4}H^{2}\log\frac{KH}{\delta}}{N_{h}\left(s,a,b\right)} \\ \leq &\mathsf{Var}_{\widehat{P}_{h,s,a,b}^{0}}\left(V\right) + \left[\left(P_{h,s,a,b}^{0} - \widehat{P}_{h,s,a,b}^{0}\right)V\right]^{2} + \frac{C_{4}H^{2}\log\frac{KH}{\delta}}{N_{h}\left(s,a,b\right)} \\ &+ \frac{1}{2}\mathsf{Var}_{P_{h,s,a,b}^{0}}\left(V\right) + \frac{C_{4}^{2}H^{2}\log\frac{KH}{\delta}}{2N_{h}\left(s,a,b\right)}, \end{split}$$

where the last relation holds due to the AM-GM inequality. Therefore, we obtain 

$$\mathsf{Var}_{P_{h,s,a,b}^{0}}\left(V\right) \leq 2\mathsf{Var}_{\widehat{P}_{h,s,a,b}^{0}}\left(V\right) + 2\left[\left(P_{h,s,a,b}^{0} - \widehat{P}_{h,s,a,b}^{0}\right)V\right]^{2} + \frac{\left(C_{4}^{2} + 2C_{4}\right)H^{2}\log\frac{KH}{\delta}}{N_{h}\left(s,a,b\right)}.$$
 (87)

Combining (87) and (84), we derive

$$\begin{aligned} &|(\hat{P}_{h,s,a,b}^{0} - P_{h,s,a,b}^{0})V| \leq \sqrt{\frac{2C_{4}^{2}}{N_{h}\left(s,a,b\right)}} \mathsf{Var}_{\hat{P}_{h,s,a,b}^{0}}\left(V\right) \log \frac{KH}{\delta} + \frac{\sqrt{C_{4}^{2}\left(C_{4}^{2} + 2C_{4}\right)H \log \frac{KH}{\delta}}}{N_{h}\left(s,a,b\right)} \\ &+ \sqrt{\frac{2C_{4}^{2}}{N_{h}\left(s,a,b\right)} \log \frac{KH}{\delta}} \left|(\hat{P}_{h,s,a,b}^{0} - P_{h,s,a,b}^{0})V\right| + \frac{C_{4}H \log \frac{KH}{\delta}}{N_{h}\left(s,a,b\right)}}. \end{aligned}$$

$$\begin{aligned} & (88) \end{aligned}$$

In the following, we consider two cases, i.e.,  $N_h(s, a, b) \leq \frac{1}{8C_4^2} \log \frac{KH}{\delta}$  and  $N_h(s, a, b) > 0$  $\frac{1}{8C_4^2}\log\frac{KH}{\delta}.$ 

For the first case of  $N_h(s, a, b) \leq \frac{1}{8C_4^2} \log \frac{KH}{\delta}$ , (34) is valid since 

$$\left|\inf_{P\in\mathcal{U}^{\sigma^{+}}(\widehat{P}^{0}_{h,s,a,b})}PV - \inf_{P\in\mathcal{U}^{\sigma^{+}}(P^{0}_{h,s,a,b})}PV\right| \leq \max_{\alpha\in[\min_{s}V(s),\max_{s}V(s)]}\left|\left(\widehat{P}^{0}_{h,s,a,b} - P^{0}_{h,s,a,b}\right)V\right|$$
$$\leq 2H = O\left(\frac{H\log\frac{KH}{\delta}}{N_{h}\left(s,a,b\right)}\right).$$
(89)

For the second case of  $N_h(s, a, b) > \frac{1}{8C_4^2} \log \frac{KH}{\delta}$ , we observe from (88) that 

$$\left| \left( \hat{P}_{h,s,a,b}^{0} - P_{h,s,a,b}^{0} \right) V \right| \leq \frac{1}{2} \left| \left( \hat{P}_{h,s,a,b}^{0} - P_{h,s,a,b}^{0} \right) V \right| + \sqrt{\frac{2C_{4}^{2}}{N_{h}\left(s,a,b\right)}} \mathsf{Var}_{\hat{P}_{h,s,a,b}^{0}}\left( V \right) \log \frac{KH}{\delta}$$

1510  
1511 
$$+ \frac{C_4 + \sqrt{C_4^2 (C_4^2 + 2C_4)}}{N_h (s, a, b)} H \log \frac{KH}{\delta}.$$

1512 Rearrange terms above and yield1513

$$\left| \left( \widehat{P}_{h,s,a,b}^{0} - P_{h,s,a,b}^{0} \right) V \right| \leq \sqrt{\frac{8C_{4}^{2}}{N_{h}\left(s,a,b\right)}} \operatorname{Var}_{\widehat{P}_{h,s,a,b}^{0}}\left(V\right) \log \frac{KH}{\delta} + 2H \frac{C_{4} + \sqrt{C_{4}^{2}\left(C_{4}^{2} + 2C_{4}\right)}}{N_{h}\left(s,a,b\right)} \log \frac{KH}{\delta}.$$
(90)

Putting (83) and (90) together, we get

$$\left| \inf_{P \in \mathcal{U}^{\sigma^+}(\hat{P}^0_{h,s,a,b})} PV - \inf_{P \in \mathcal{U}^{\sigma^+}(P^0_{h,s,a,b})} PV \right| \leq \sqrt{\frac{8C_4^2}{N_h(s,a,b)}} \mathsf{Var}_{\hat{P}^0_{h,s,a,b}}(V) \log \frac{KH}{\delta} + 2H \frac{C_4 + \sqrt{C_4^2(C_4^2 + 2C_4)}}{N_h(s,a,b)} \log \frac{KH}{\delta}.$$
 (91)

Putting the above bounds for two cases together, we conclude the proof of (34).

**Part 2: proof of inequality (35).** In the process of proving inequality (35), we just divide the problem into two cases, i.e.,  $N_h(s, a, b) < 16C_4^2 \log \frac{KH}{\delta}$  and  $N_h(s, a, b) \ge 16C_4^2 \log \frac{KH}{\delta}$ .

For the first case of  $N_h(s, a, b) < 16C_4^2 \log \frac{KH}{\delta}$ , the result (35) is valid since

$$\operatorname{Var}_{\widehat{P}_{h,s,a,b}^{0}}\left(V\right) \leq H^{2} = O\left(\frac{H^{2}\log\frac{KH}{\delta}}{N_{h}\left(s,a,b\right)}\right).$$

1537 For the second case of  $N_h(s, a, b) \ge 16C_4^2 \log \frac{KH}{\delta}$ , there is

$$\begin{aligned} \mathsf{Var}_{\widehat{P}_{h,s,a,b}^{0}}\left(V\right) \stackrel{\text{(i)}}{=} \mathsf{Var}_{P_{h,s,a,b}^{0}}\left(V\right) - \left[\left(P_{h,s,a,b}^{0} - \widehat{P}_{h,s,a,b}^{0}\right)V\right]^{2} - \left(P_{h,s,a,b}^{0} - \widehat{P}_{h,s,a,b}^{0}\right)\left(\overline{V} \circ \overline{V}\right) \\ \stackrel{\text{(ii)}}{\leq} \mathsf{Var}_{P_{h,s,a,b}^{0}}\left(V\right) + C_{4}\sqrt{\frac{H^{2}}{N_{h}\left(s,a,b\right)}}\mathsf{Var}_{P_{h,s,a,b}^{0}}\left(V\right)\log\frac{KH}{\delta}} + \frac{C_{4}H^{2}\log\frac{KH}{\delta}}{N_{h}\left(s,a,b\right)} \\ \stackrel{\text{(iii)}}{\leq} 2\mathsf{Var}_{P_{h,s,a,b}^{0}}\left(V\right) + \frac{\left(C_{4}^{2}/4 + C_{4}\right)H^{2}\log\frac{KH}{\delta}}{N_{h}\left(s,a,b\right)} \end{aligned}$$

 $= 2 \mathsf{Var}_{P^0_{h,s,a,b}}\left(V\right) + O\left(\frac{H^2 \log \frac{KH}{\delta}}{N_h\left(s,a,b\right)}\right),$ 

where (i) comes from (85), (ii) holds due to (86), and (iii) exists under the AM-GM inequality.
Putting the two cases together, we complete the proof of (35). Thus, Lemma 3 is finally proven.

1552 C.4 Proof of Lemma 4 

1554 Assuming that  $\widehat{Q}_{h}^{+}(s, a, b) \ge Q_{h}^{\star, \widehat{\nu}, \sigma^{+}}(s, a, b)$  holds, then we can easily obtain  $\widehat{V}_{h}^{+}(s) \ge V_{h}^{\star, \nu, \sigma^{+}}(s)$ , 1556 since

$$\widehat{V}_{h}^{+}(s) = \mathbb{E}_{a \sim \mu_{h}^{+}(s), b \sim \nu_{h}^{+}(s)} \left[ \widehat{Q}_{h}^{+}(s, a, b) \right]$$

$$\stackrel{(i)}{\geq} \mathbb{E}_{a \sim \mu^{\star}(s), b \sim \widehat{\nu}(s)} \left[ \widehat{Q}_{h}^{+}(s, a, b) \right] \geq \mathbb{E}_{a \sim \mu^{\star}(s), b \sim \widehat{\nu}(s)} \left[ Q_{h}^{\star, \widehat{\nu}, \sigma^{+}}(s, a, b) \right] = V_{h}^{\star, \widehat{\nu}, \sigma^{+}}(s),$$

where (i) holds due to the fact that  $\hat{\nu} = \nu_h^+$  and  $(\mu_h^+, \nu_h^+)$  is the Nash equilibrium of  $\hat{Q}_h^+(s, a, b)$ . Consequently, we just need to verify

$$\widehat{Q}_h^+(s,a,b) \ge Q_h^{\star,\widehat{\nu},\sigma^+}(s,a,b),\tag{92}$$

which is obtained by induction.

1566 It can be easily verified that (92) holds at the base case when h = H + 1 under the trivial fact 1567  $\widehat{Q}_{H+1}^+(s,a,b) = Q_{H+1}^{\star,\sigma^+}(s,a,b) = 0.$ 1568 Suppose that (92) holds for all  $(s, a, b) \in S \times A \times B$  at some time step  $h \in [H]$  next. 1569 1570 According to the update rule in line 4 in Algorithm 2, (92) exists if  $\hat{Q}_h^+(s,a,b) = H$  because 1571  $\widehat{Q}_h^+(s,a,b) = H \ge Q_h^{\star,\widehat{\nu},\sigma^+}(s,a,b).$ 1572 1573 Besides, in the case of  $N_h(s, a, b) = 0$ , we have  $\beta_h\left(s, a, b, \widehat{V}_{h+1}^+\right) = H$ , leading to  $\widehat{Q}_h^+(s, a, b) = 0$ 1574  $H \ge Q_h^{\star,\widehat{\nu},\sigma^+}(s,a,b)$ . Otherwise, for  $N_h(s,a,b) > 0$ ,  $\widehat{Q}_h^+(s,a,b)$  is updated as 1575 1576  $\widehat{Q}_{h}^{+}\left(s,a,b\right) = \widehat{r}\left(s,a,b\right) + \inf_{P \in \mathcal{U}^{\sigma^{+}}\left(\widehat{P}_{h,s,a,b}^{0}\right)} P\widehat{V}_{h+1}^{+} + \beta_{h}\left(s,a,b,\widehat{V}_{h+1}^{+}\right)$  $\geq \widehat{r}\left(s,a,b\right) + \inf_{P \in \mathcal{U}^{\sigma^+}\left(P_{h,s,a,b}^0\right)} P\widehat{V}_{h+1}^+ + \beta_h\left(s,a,b,\widehat{V}_{h+1}^+\right)$ 1579 1581  $-\left|\inf_{P\in\mathcal{U}^{\sigma^+}(\widehat{P}^0_{h,s,a,b})}P\widehat{V}^+_{h+1}-\inf_{P\in\mathcal{U}^{\sigma^+}(P^0_{h,s,a,b})}P\widehat{V}^+_{h+1}\right|$  $\geq \widehat{r}(s,a,b) + \inf_{\substack{P \in \mathcal{U}^{\sigma^+}(P_{h,s,a,b}^0)}} P\widehat{V}_{h+1}^+ + 0$  $\geq \widehat{r}(s,a,b) + \inf_{\substack{P \in \mathcal{U}^{\sigma^+}(P_{h,s,a,b}^0)}} P\widehat{V}_{h+1}^{\star,\widehat{\nu},\sigma^+} + 0 = Q_h^{\star,\widehat{\nu},\sigma^+}(s,a,b),$ 1585 (93)1587 1588 where the second inequality holds due to (36) in Lemma 3 and the last equality comes from the 1589 empirical robust Bellman equation (33). 1590 1591 Armed with the case of h = H + 1, we complete prove Lemma 4 by induction. 1592

#### 1593 C.5 PROOF OF LEMMA 5 1594

1597 1598

1604 1605

1609 1610 1611

Following the proof by Shi et al. (2024b, Lemma 3), we bound  $\min_{s \in S} \hat{V}_h^+(s)$  and  $\max_{s \in S} \hat{V}_h^+(s)$ , respectively. Specifically, we have

$$\min_{s \in \mathcal{S}} \widehat{V}_{h}^{+}(s) = \min_{s \in \mathcal{S}} \mathbb{E}_{(a,b) \sim \mu_{h}^{+} \times \nu_{h}^{+}} \left[ \widehat{Q}_{h}^{+}(s,a,b) \right] \\
= \min_{s \in \mathcal{S}} \mathbb{E}_{(a,b) \sim \mu_{h}^{+} \times \nu_{h}^{+}} \left[ \widehat{r}_{h}(s,a,b) + \inf_{P \in \mathcal{U}^{\sigma^{+}}(\widehat{P}_{h,s,a,b}^{0})} P\widehat{V}_{h+1}^{+} + \beta_{h} \left( s,a,b, \widehat{V}_{h+1} \right) \right] \\
\geq 0 + \min_{s \in \mathcal{S}} \widehat{V}_{h+1}^{+}(s) + 0,$$
(94)

where the middle equality is valid due to the update rule in line 4 in Algorithm 2. Similarly, there is

$$\max_{s \in \mathcal{S}} \widehat{V}_{h}^{+} = \max_{s \in \mathcal{S}} \mathbb{E}_{(a,b) \sim \mu_{h}^{+} \times \nu_{h}^{+}} \left[ \widehat{Q}_{h}^{+}(s,a,b) \right]$$

$$= \max_{s \in \mathcal{S}} \mathbb{E}_{(a,b) \sim \mu_{h}^{+} \times \nu_{h}^{+}} \left[ \widehat{r}_{h}(s,a,b) + \inf_{P \in \mathcal{U}^{\sigma^{+}}(\widehat{P}_{h,s,a,b}^{0})} P\widehat{V}_{h+1}^{+} + \beta_{h} \left( s,a,b, \widehat{V}_{h+1} \right) \right]$$

$$\leq 1 + \max_{(s,a,b) \in \mathcal{S} \times \mathcal{A} \times \mathcal{B}} \inf_{P \in \mathcal{U}^{\sigma^{+}}(\widehat{P}_{h,s,a,b})} P\widehat{V}_{h+1}^{+} + H.$$
(95)

In order to prove Lemma 5, we here introduce several useful notations. For any  $h \in [H]$ , there exists at least one state  $s_h^*$  that satisfies  $\widehat{V}_h^+(s_h^*) = \min_{s \in S} \widehat{V}_h^+(s)$ .

Furthermore, for any accessible uncertainty set  $\sigma^+ > 0$  and  $(s, a, b) \in \mathcal{S} \times \mathcal{A} \times \mathcal{B}$ , we define an auxiliary vector  $\hat{P}'_{h,s,a,b} \in \mathbb{R}^S$  by reducing the values of several elements of  $\hat{P}^0_{h,s,a,b}$  strictly, namely,  $0 \leq \hat{P}'_{h,s,a,b} \leq \hat{P}^0_{h,s,a,b}$  and  $\sum \hat{P}^0_{h,s,a,b}(s') - \hat{P}'_{h,s,a,b}(s') = \left\| \hat{P}'_{h,s,a,b} - \hat{P}^0_{h,s,a,b} \right\|_{*} = \sigma^+$ .

1618 
$$0 \le P'_{h,s,a,b} \le P^0_{h,s,a,b}$$
 and  $\sum_{s' \in S} P^0_{h,s,a,b}(s') - P'_{h,s,a,b}(s') = \left\| P'_{h,s,a,b} - P^0_{h,s,a,b} \right\|_1$   
1619

(96)

We use  $l_{s_h^*}$  to represent a S-dimensional standard basis under  $s_h^*$ , we can derive that

$$\frac{1}{2} \left\| \widehat{P}_{h,s,a,b}^{\prime} + \sigma^{+} \left[ l_{s_{h}^{\star}} \right]^{\top} - \widehat{P}_{h,s,a,b}^{0} \right\|_{1} \leq \frac{1}{2} \left\| \widehat{P}_{h,s,a,b}^{\prime} - \widehat{P}_{h,s,a,b}^{0} \right\|_{1} + \frac{1}{2} \left\| \sigma^{+} \left[ l_{s_{h}^{\star}} \right]^{\top} \right\|_{1} \leq \sigma^{+}, \quad (97)$$

where the first inequality is valid since the 'distance' function (e.g., TV distance) satisfies the triangle inequality.

1627 Therefore, we can conclude that  $\hat{P}'_{h,s,a,b} + \sigma^+ [l_{s_h^\star}]^\top \in \mathcal{U}^{\sigma^+}(\hat{P}^0_{h,s,a,b})$  and  $\hat{P}'_{h,s,a,b} + \sigma^+ [l_{s_h^\star}]^\top$  is a 1628 distribution vector based on (97), leading to 

$$\inf_{P \in \mathcal{U}^{\sigma^+}(\hat{P}^0_{h,s,a,b})} P \hat{V}^+_{h+1} \leq \left( \hat{P}'_{h,s,a,b} + \sigma^+ [l_{s^*_{i,h}}]^+ \right) \hat{V}^+_{h+1} \\
\leq \| \hat{P}'_{h,s,a,b} \|_1 \| \hat{V}^+_{h+1} \|_\infty + \sigma^+ \hat{V}^+_{h+1}(s^*_{h+1}) \\
\leq (1 - \sigma^+) \max_{s \in \mathcal{S}} \hat{V}^+_{h+1}(s) + \sigma^+ \min_{s \in \mathcal{S}} \hat{V}^+_{h+1}(s),$$
(98)

where the last inequality holds since

$$\begin{aligned} & \|P_{h,s,a,b}'\|_1 = \sum_{s'} P_{h,s,a,b}'(s') = -\sum_{s'} \left( P_{h,s,a,b}^0(s') - P_{h,s,a,b}'(s') \right) + \sum_{s'} P_{h,s,a,b}^0(s') = 1 - \sigma^+. \end{aligned}$$

$$\begin{aligned} & \text{(99)} \end{aligned}$$

## 1641 Putting (98) and (95) together shows

$$\max_{s \in \mathcal{S}} \widehat{V}_{h}^{+}(s) \leq 1 + \max_{(s,a,b) \in \mathcal{S} \times \mathcal{A} \times \mathcal{B}} \inf_{P \in \mathcal{U}^{\sigma^{+}}(P_{h,s,a,b}^{0})} P \widehat{V}_{h+1}^{+} + H$$
$$\leq H + 1 + (1 - \sigma^{+}) \max_{s \in \mathcal{S}} \widehat{V}_{h+1}^{+}(s) + \sigma^{+} \min_{s \in \mathcal{S}} \widehat{V}_{h+1}^{+}(s).$$
(100)

Taking the result (100) with (94), we obtain

$$\max_{s \in S} \widehat{V}_{h}^{+} - \min_{s \in S} \widehat{V}_{h}^{+} \\
\leq H + 1 + (1 - \sigma^{+}) \max_{s \in S} \widehat{V}_{h+1}^{+}(s) + \sigma^{+} \min_{s \in S} \widehat{V}_{h+1}^{+}(s) - \min_{s \in S} V_{h+1}^{+}(s) \\
= H + 1 + (1 - \sigma^{+}) \left( \max_{s \in S} \widehat{V}_{h+1}^{+}(s) - \min_{s \in S} \widehat{V}_{h+1}^{+}(s) \right) \\
\leq H + 1 + (1 - \sigma^{+}) \left[ H + 1 + (1 - \sigma^{+}) \left( \max_{s \in S} \widehat{V}_{h+2}^{+}(s) - \min_{s \in S} \widehat{V}_{h+2}^{+}(s) \right) \right] \\
\leq \dots \leq \frac{(H + 1) \left( 1 - (1 - \sigma^{+})^{H - h} \right)}{\sigma^{+}}.$$
(101)

Combining this result with  $\max_{s \in S} \widehat{V}_h^+(s) - \min_{s \in S} \widehat{V}_h^+(s) \le H$ , we complete the proof.

## 1663 C.6 PROOF OF LEMMA 6

First of all, we introduce some auxiliary values and reward functions to control  $\sum_{i=1}^{H} \sum_{(s,b) \in S \times B} d_i^{p,\mu^d,\nu^*}(s,\mu^d(s),b) \operatorname{Var}_{P_{i,s,\mu^d(s),b}^0}(\widehat{V}) \text{ as below: for any time step } i$   $\widehat{V}_i^{\mathsf{m}} \coloneqq \min_{s \in S} \widehat{V}_i^{+}(s) \colon \text{the minimum value of all the entries in vector } \widehat{V}_i^{+}.$   $\widehat{V}_i' \coloneqq \widehat{V}_i^{+} - \widehat{V}_i^{\mathsf{m}} 1 \colon \text{truncated value function.}$   $\widehat{r}_i^{\mu^d,\nu^*}(s) \equiv \mathbb{E}_{(a,b) \sim (\mu^d(s),\nu^*(s))} \widehat{r}_i(s,a,b) \colon \text{average reward function.}$   $\widehat{r}_i^{\mathsf{m}} = r_i^{\mu^d,\nu^*} + (\widehat{V}_{i+1}^{\mathsf{m}} - \widehat{V}_i^{\mathsf{m}}) 1 \colon \text{truncated reward function.}$  Then applying the robust Bellman's consistency equation in (33) gives 

$$\widehat{V}'_i = \widehat{V}^+_i - \widehat{V}^\mathsf{m}_i 1$$

1677  
1678 
$$\stackrel{(i)}{\leq} \hat{r}_{i}^{\mu^{d},\nu^{\star}} + \tilde{P}_{i}^{\inf,\hat{V}}\hat{V}_{i+1}^{+} + 2\beta_{i}^{\mu^{d},\nu^{\star}} - \hat{V}_{i}^{\mathsf{m}}$$

$$\begin{aligned}
\stackrel{(1)}{=} \widehat{r}_{i}^{\mu^{d},\nu^{\star}} &+ \widetilde{P}_{i}^{\inf,\widehat{V}}\widehat{V}_{i+1}^{+} + 2\beta_{i}^{\mu^{d},\nu^{\star}} - \widehat{V}_{i}^{\mathsf{m}}\mathbf{1} \\
\stackrel{(67)}{=} \widehat{r}_{i}^{\mu^{d},\nu^{\star}} &+ \widetilde{P}_{i}^{\inf,\widehat{V}}\widehat{V}_{i+1}^{+} + \left(\widehat{V}_{i+1}^{\mathsf{m}}\mathbf{1} - \widehat{V}_{i}^{\mathsf{m}}\mathbf{1}\right) - \widehat{V}_{i+1}^{\mathsf{m}}\mathbf{1} + 2\beta_{i}^{\mu^{d},\nu^{\star}} \\
\stackrel{(68)}{=} \widehat{r}_{i}^{\mathsf{m}} + \widetilde{P}_{i}^{\inf,\widehat{V}}\widehat{V}_{i+1}^{+} - \widehat{V}_{i+1}^{\mathsf{m}}\mathbf{1} + 2\beta_{i}^{\mu^{d},\nu^{\star}}
\end{aligned}$$

1682  
1683 
$$= \hat{r}_{i}^{\mathsf{m}} + \tilde{P}_{i}^{\inf,\hat{V}}\hat{V}_{i+1}' + 2\beta_{i}^{\mu^{\mathsf{d}},\nu^{\star}}, \qquad (102)$$

where (i) follows from the fact that

Here, (ii) is valid under the notation 

$$P_{i,s,a,b}^{\inf,\widehat{V}} \coloneqq \operatorname{argmin}_{P \in \mathcal{U}^{\sigma^+}\left(P_{i,s,a,b}^0\right)} P\widehat{V}_{i+1}^+$$
(104)

and (iii) holds under the notation as  $\widetilde{P}_{i,s}^{\inf,\widehat{V}} \coloneqq \mathbb{E}_{(a,b)\sim(\mu^{d}(s),\nu^{\star}(s))}P_{i,s,a,b}^{\inf,\widehat{V}}$  and the sequence as  $\widetilde{P}_{i}^{\inf,V} \in \mathbb{R}^{S\times S}$  Besides, (i) in (103) exists due to (36) in Lemma 3 for  $N_{i}(s, a, b) > 0$  and 

$$\left| \inf_{P \in \mathcal{U}^{\sigma^+}(P^0_{i,s,a,b})} P \widehat{V}^+_{i+1} - \inf_{P \in \mathcal{U}^{\sigma^+}(\widehat{P}^0_{i,s,a,b})} P \widehat{V}^+_{i+1} \right| \le H = \beta_i^{\mu^d,\nu^*}(s)$$
(105)

for  $N_i(s, a, b) = 0$ . 

The above fact leads to 

where (i) follows from the fact that  $\operatorname{Var}_{P_{i,s}^{\inf,V}}(V-b1) = \operatorname{Var}_{P_{i,s}^{\inf,V}}(V)$  for any value vector  $V \in \mathbb{R}^S$ and scalar *b*, (ii) holds with the fact

$$\begin{split} & \mathbb{E}_{(a,b)\sim(\mu^{\mathsf{d}}(s),\nu^{\star}(s))}\left[\left(P_{i,s,a,b}^{\inf,\widehat{V}}\widehat{V}'_{i+1}\right)\circ\left(P_{i,s,a,b}^{\inf,\widehat{V}}\widehat{V}'_{i+1}\right)\right] \\ \geq & \mathbb{E}_{(a,b)\sim(\mu^{\mathsf{d}}(s),\nu^{\star}(s))}\left[\left(P_{i,s,a,b}^{\inf,\widehat{V}}\widehat{V}'_{i+1}\right)\right]\circ\mathbb{E}_{(a,b)\sim(\mu^{\mathsf{d}}(s),\nu^{\star}(s))}\left[\left(P_{i,s,a,b}^{\inf,\widehat{V}}\widehat{V}'_{i+1}\right)\right] \end{split}$$

1740 (iv) arises from  $\hat{r}_i^m \le r_i \le 1$  due to  $\hat{V}_{i+1}^m - \hat{V}_i^m \le 0$  by definition, and (iii) comes from (102). 1741 Consequently, combining (48), we arrive at 

$$\begin{aligned}
& 1744 \\
& 1745 \\
& \sum_{(s,b)\in S\times B} d_i^{p,\mu^d,\nu^*}(s,\mu^d(s),b) \operatorname{Var}_{P_{i,s,a,b}^{\operatorname{inf},V}}(\hat{V}_{i+1}^+) \\
& 1747 \\
& = \sum_{s\in S} d_i^{p,\mu^d,\nu^*}(s) \mathbb{E}_{(a,b)\sim(\mu^d(s),\nu^*(s))} \operatorname{Var}_{P_{i,s,a,b}^{\operatorname{inf},V}}(\hat{V}_{i+1}^+) \\
& 1749 \\
& \leq \sum_{s\in S} d_i^{p,\mu^d,\nu^*}(s) \left( \widetilde{P}_{i,s}^{\operatorname{inf},V}\left( \widehat{V}_{i+1}' \circ \widehat{V}_{i+1}' \right) - \widehat{V}_i'(s) \circ \widehat{V}_i'(s) + \left( \left\| \widehat{V}_i' \right\|_{\infty} + \left\| \widehat{V}_{i+1}' \right\|_{\infty} \right) \left( 2\beta_i^{\mu^d,\nu^*}(s) + 1 \right) \right) \\
& 1751 \\
& \leq \sum_{s\in S} d_i^{p,\mu^d,\nu^*}(s) \left( \widetilde{P}_{i,s}^{\operatorname{inf},V}\left( \widehat{V}_{i+1}' \circ \widehat{V}_{i+1}' \right) - \widehat{V}_i'(s) \circ \widehat{V}_i'(s) \right) + \left( \left\| \widehat{V}_i' \right\|_{\infty} + \left\| \widehat{V}_{i+1}' \right\|_{\infty} \right) \\
& 1752 \\
& \leq \sum_{s\in S} d_i^{p,\mu^d,\nu^*}(s) \left( \widetilde{P}_{i,s}^{\operatorname{inf},V}\left( \widehat{V}_{i+1}' \circ \widehat{V}_{i+1}' \right) - \widehat{V}_i'(s) \circ \widehat{V}_i'(s) \right) + \left( \left\| \widehat{V}_i' \right\|_{\infty} + \left\| \widehat{V}_{i+1}' \right\|_{\infty} \right) \\
& 1754 \\
& + 2 \left( \left\| \widehat{V}_i' \right\|_{\infty} + \left\| \widehat{V}_{i+1}' \right\|_{\infty} \right) \sum_{s\in S} d_i^{p,\mu^d,\nu^*}(s) \widehat{V}_i'(s) \circ \widehat{V}_i'(s) \right) + \left( \left\| \widehat{V}_i' \right\|_{\infty} + \left\| \widehat{V}_{i+1}' \right\|_{\infty} \right) \\
& 1759 \\
& + 2 \left( \left\| \widehat{V}_i' \right\|_{\infty} + \left\| \widehat{V}_{i+1}' \right\|_{\infty} \right) \sum_{s\in S} d_i^{p,\mu^d,\nu^*}(s) \widehat{V}_i'(s) \circ \widehat{V}_i'(s) \right) + \left( \left\| \widehat{V}_i' \right\|_{\infty} + \left\| \widehat{V}_{i+1}' \right\|_{\infty} \right) \\
& 1761 \\
& 1762 \\
& = \sum_{s\in S} \left( d_{i+1}^{p,\mu^d,\nu^*}(s) \left( \widehat{V}_{i+1}'(s) \circ \widehat{V}_{i+1}'(s) \right) - d_i^{p,\mu^d,\nu^*}(s) \widehat{V}_i'(s) \circ \widehat{V}_i'(s) \right) + \left( \left\| \widehat{V}_i' \right\|_{\infty} + \left\| \widehat{V}_{i+1}' \right\|_{\infty} \right) \\
& 1764 \\
& + 2 \left( \left\| \widehat{V}_i' \right\|_{\infty} + \left\| \widehat{V}_{i+1}' \right\|_{\infty} \right) \sum_{(s,b)\in S\times B} d_i^{p,\mu^d,\nu^*}(s,\mu^d(s),b) \beta_i(s,\mu^d(s),b,\widehat{V}). \quad (107)
\end{aligned}$$

Besides, under TV distance, we have

$$\begin{aligned} \left| \mathsf{Var}_{P_{i,s,a,b}^{0}}(\widehat{V}_{i+1}^{+}) - \mathsf{Var}_{P_{i,s,a,b}^{\inf,V}}(\widehat{V}_{i+1}^{+}) \right| &= \left| \mathsf{Var}_{P_{i,s,a,b}^{0}}(\widehat{V}_{i+1}') - \mathsf{Var}_{P_{i,s,a,b}^{\inf,V}}(\widehat{V}_{i+1}') \right| \\ &\leq \left\| P_{i,s,a,b}^{0} - P_{i,s,a,b}^{\inf,V} \right\|_{1} \left\| \widehat{V}_{i+1}' \right\|_{\infty}^{2} \\ &\leq \sigma^{+} \left\| \widehat{V}_{i+1}' \right\|_{\infty}^{2} \leq (H+1) \left\| \widehat{V}_{i+1}' \right\|_{\infty}^{2}, \end{aligned}$$
(108)

where the last inequality comes from Lemma 5.

Therefore, we derive  $\sum_{i=1}^{\infty} \sum_{\substack{(a,b) \in \mathcal{S} \times \mathcal{R}}} d_i^{\mathfrak{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b) \mathsf{Var}_{P_{i,s,a,b}^0}(\widehat{V}_{i+1}^+)$  $\leq \sum_{i=1}^{n} \sum_{(a,b) \in S \times B} d_i^{\mathbf{p},\mu^{\mathbf{d}},\nu^{\star}}(s,\mu^{\mathbf{d}}(s),b) \mathsf{Var}_{P_{i,s,a,b}^{\mathrm{inf},V}}(\widehat{V}_{i+1}^+)$  $+\sum_{i=1}^{H}\sum_{(l,l)\in\mathcal{C}_{i},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b)\left|\mathsf{Var}_{P_{i,s,a,b}^{0}}(\widehat{V}_{i+1}^{+})-\mathsf{Var}_{P_{i,s,a,b}^{\mathrm{inf},V}}(\widehat{V}_{i+1}^{+})\right|$  $\leq \sum_{i=1}^{H} 2\left(\left\|\widehat{V}_{i}'\right\|_{\infty} + \left\|\widehat{V}_{i+1}'\right\|_{\infty}\right) \sum_{i=1}^{L} d_{i}^{\mathbf{p},\mu^{d},\nu^{\star}}(s)\beta_{i}^{\mu^{d},\nu^{\star}}(s) + \sum_{i=1}^{H} \left(\left\|\widehat{V}_{i}'\right\|_{\infty} + (H+2)\left\|\widehat{V}_{i+1}'\right\|_{\infty}\right)$  $+\sum_{a} d_{H+1}^{\mathbf{p},\mu^{d},\nu^{\star}}(s) \widehat{V}_{H+1}'(s) \circ \widehat{V}_{H+1}'(s)$  $\leq 4 \sum_{i=1}^{H} \min\left\{\frac{(H+1)\left(1-(1-\sigma^{+})^{H-i}\right)}{\sigma^{+}}, H\right\} \sum_{(i,j)\in\mathcal{C}_{*},\mathcal{P}} d_{i}^{\mathbf{p},\mu^{\mathsf{d}},\nu^{\star}}(s,\mu^{\mathsf{d}}(s),b)\beta_{i}(s,\mu^{\mathsf{d}}(s),b,\widehat{V})$  $+(H+3)\sum_{i=1}^{H}\min\left\{\frac{(H+1)\left(1-(1-\sigma^{+})^{H-i}\right)}{\sigma^{+}},H\right\}$  $\stackrel{(i)}{\leq} 4 \sum_{i=1}^{H} \min\left\{\frac{(H+1)\left(1-(1-\sigma^{+})^{H-i}\right)}{\sigma^{+}}, H\right\} \sum_{i=1}^{H} \sum_{(s,h) \in S \times \mathcal{B}} d_{i}^{\mathbf{p}, \mu^{\mathsf{d}}, \nu^{\star}}(s, \mu^{\mathsf{d}}(s), b) \beta_{i}(s, \mu^{\mathsf{d}}(s), b, \widehat{V})$  $+(H+3)\sum_{i=1}^{H}\min\left\{\frac{(H+1)\left(1-(1-\sigma^{+})^{H-i}\right)}{\sigma^{+}},H\right\}$  $\overset{\text{(ii)}}{\leq} 4H \min\left\{\frac{2(H\sigma^+ - 1 + (1 - \sigma^+)^H)}{(\sigma^+)^2}, H\right\} \sum_{i=1}^H \sum_{(i,j) \in \mathcal{Q}} d_i^{\mathsf{p},\mu^\mathsf{d},\nu^\star}(s,\mu^\mathsf{d}(s),b)\beta_i(s,\mu^\mathsf{d}(s),b,\widehat{V})$ +  $(H+3)H\min\left\{\frac{2(H\sigma^+ - 1 + (1-\sigma^+)^H)}{(\sigma^+)^2}, H\right\},\$ (109)

where (i) comes from Cauchy-Schwarz inequality and the (ii) holds since

$$\begin{split} \sum_{i=1}^{H} \frac{(H+1)\left(1-(1-\sigma^{+})^{H-i}\right)}{\sigma^{+}} = & \frac{H(H+1)}{\sigma^{+}} - \sum_{i=0}^{H-1} \frac{(H+1)(1-\sigma^{+})^{i}}{\sigma^{+}} \\ &= & \frac{H(H+1)}{\sigma^{+}} - \frac{(H+1)(1-(1-\sigma^{+})^{H})}{(\sigma^{+})^{2}} \\ &= & \frac{(H+1)(H\sigma^{+}-1+(1-\sigma^{+})^{H})}{(\sigma^{+})^{2}} \\ &\leq & \frac{2H(H\sigma^{+}-1+(1-\sigma^{+})^{H})}{(\sigma^{+})^{2}}. \end{split}$$

#### **PROOF OF THEOREM 2** D

In this section, we focus on a simpler class of RTZMGs: robust Markov decision processes (MDPs), which are single-agent versions of RTZMGs. 

Before proceeding, we briefly define a Robust MDP (RMDP) in the finite-horizon episodic Recall that an RTZMG with an uncertainty set is represented as  $\mathcal{MG}$ setting. =  $\{\mathcal{S}, \mathcal{A}, \mathcal{B}, \mathcal{U}^{\sigma^+}(P^0), \mathcal{U}^{\sigma^-}(P^0), r, H\}$ . For simplicity, we assume  $\mathcal{A} \geq \mathcal{B}$ , and set  $|\mathcal{B}| = 1$ , meaning the min-player's actions do not affect transitions or rewards. Thus, finding a robust NE in such RTZMGs reduces to finding the max-player's optimal policy in a corresponding RMDP  $\mathcal{M}_{r} = \{S, \mathcal{A}, \mathcal{U}^{\sigma^{+}}(P^{0}), r, H\}.$ 

Thus, in this section, we construct the lower bound for finding the optimal policy in RTZMGs, which also implies a lower bound for finding robust NE in RTZMGs. We first highlight a useful property about KL divergence from Tsybakov (2008, Lemma 2.7), which can be helpful in this section.

**Lemma 7** For any  $p, q \in (0, 1)$ , it holds that

 $\mathsf{KL}(p \parallel q) \le \frac{(p-q)^2}{q(1-q)}.$ (110)

### 1847 D.1 STEP 1: CONSTRUCTING A FAMILY OF HARD MARKOV GAME INSTANCES

1849 The hard instances developed here differ from standard MDP since we need to consider that the transition kernel can be perturbed in robust MDPs.

**Constructing hard robust MDP instances.** To begin with, we first introduce an auxiliary collection  $\Phi \subseteq \{0,1\}^H$ , consisting of *H*-dimensional vectors. In addition, resorting to the Gilbert-Varshamov lemma (Gilbert, 1952), we notice that there exists a set  $\Phi \subseteq \{0,1\}^H$  such that:

for any 
$$\phi, \widetilde{\phi} \in \Phi$$
 obeying  $\phi \neq \widetilde{\phi}$ :  $\|\phi - \widetilde{\phi}\|_1 \ge \frac{H}{8}$  and  $|\Phi| \ge e^{H/8}$ . (111)

With this in mind, we construct a set of RMDPs as below:

$$\mathcal{M}(\mathcal{F}, \Phi) \coloneqq \left\{ \mathcal{M}_{f}^{\phi} = \left( \mathcal{S}, \mathcal{A}, \mathcal{U}^{\sigma^{+}}(P^{f, \phi}), r, H \right) | f \in \mathcal{F} = \{0, 1, \cdots, SA - 1\}, \\ \phi = [\phi_{h}]_{1 \le h \le H} \in \Phi \right\},$$
(112)

65 where

1844

1845 1846

1851

1855 1856 1857

1859 1860

1862 1863 1864

1866 1867

$$S = \{0, 1, \dots, S - 1\},$$
 and  $A = \{0, 1, \dots, A - 1\},$ 

1868 and  $\sigma^+$  will be introduced momentarily.

In simple terms, the collection  $\mathcal{M}(\mathcal{F}, \Phi)$  consists of SA subsets, each containing  $|\Phi|$  different RMDPs associated with some  $f \in \mathcal{F}$ . The state space for each RMDP  $\mathcal{M}_{f}^{\phi} \in \mathcal{M}(\mathcal{F}, \Phi)$ , denoted as  $S_{\text{one}}$ , includes two types of states:  $\mathcal{M} = \{m_i \mid i \in \mathcal{F}\}$  and  $\mathcal{N} = \{n_i \mid i \in \mathcal{F}\}$ . Each state in  $\mathcal{M}$  and  $\mathcal{N}$  has two possible actions,  $\mathcal{A}_{\text{one}} = \{0, 1\}$ . Thus, there are a total of 2SA states and 4SAstate-action pairs.

1875 With these notations, we define the transition kernels for  $\mathcal{M}(\mathcal{F}, \Phi)$ . For any RMDP  $\mathcal{M}_{f}^{\phi} \in \mathcal{M}(\mathcal{F}, \Phi)$ , the transition kernel  $P^{f,\phi} = \{P_{h}^{f,\phi}\}_{h=1}^{H}$  is defined as follows, for any  $(s, a, s', h) \in \mathcal{S}_{one} \times \mathcal{A}_{one} \times \mathcal{S}_{one} \times [H]$ , 1878

$$P_{h}^{f,\phi}(s' \mid s, a) = \begin{cases} p\mathbb{1}(s' = n_{f}) + (1-p)\mathbb{1}(s' = s) & \text{if } s = m_{f}, a = \phi_{h} \\ q\mathbb{1}(s' = n_{f}) + (1-q)\mathbb{1}(s' = s) & \text{if } s = m_{f}, a = 1 - \phi_{h} \\ \mathbb{1}(s' = s) & \text{otherwise} \end{cases}$$
(113)

1882 where p and q follow  $p > q \ge \frac{1}{2}$ .

In addition, the reward function is defined as

$$\forall (h, s, a) \in [H] \times \mathcal{S}_{one} \times \mathcal{A}_{one} : \quad r_h(s, a) = \begin{cases} 1 & \text{if } s \in \mathcal{N} \\ 0 & \text{otherwise.} \end{cases}$$
(114)

1888 Uncertainty set of the transition kernels. Denote the transition kernel vector as

18

1879 1880 1881

1885

(115)

Recalling the uncertainty set defined in (1), we know that  $\mathcal{U}^{\sigma^+}(P^{f,\phi})$  represents:

1892  
1892 
$$\mathcal{U}^{\sigma^+}(P^{f,\phi}) \coloneqq \otimes \mathcal{U}^{\sigma^+}(P^{f,\phi}_{h,s,a}), \quad \mathcal{U}^{\sigma^+}(P^{f,\phi}_{h,s,a}) \coloneqq \left\{ \widetilde{P}^{f,\phi}_{h,s,a} \in \Delta(\mathcal{S}) : \rho\left(\widetilde{P}^{f,\phi}_{h,s,a} - P^{f,\phi}_{h,s,a}\right) \le \sigma^+ \right\},$$
1894 where  $\odot$  represents the Casterior and dust even  $(h, s, a) \in [U] \hookrightarrow \mathcal{S}$  and  $A$ 

where  $\otimes$  represents the Cartesian product over  $(h, s, a) \in [H] \times S_{one} \times A_{one}$ .

For the convenience of the subsequent proof, we analyze the TV distance as an uncertainty set for example, which means

$$\mathcal{U}^{\sigma^+}(P^{f,\phi}_{h,s,a}) \coloneqq \left\{ \widetilde{P}^{f,\phi}_{h,s,a} \in \Delta(\mathcal{S}) : \frac{1}{2} \left\| \widetilde{P}^{f,\phi}_{h,s,a} - P^{f,\phi}_{h,s,a} \right\| \le \sigma^+ \right\}.$$
(116)

1901 1902 Next, we introduce useful notations and facts for this section. For any RMDP  $\mathcal{M}_{f}^{\phi} \in \mathcal{M}(\mathcal{F}, \Phi)$  and 1903 any  $(h, s, a, s') \in [H] \times \mathcal{S}_{one} \times \mathcal{A}_{one} \times \mathcal{S}_{one}$ , we define the minimum transition probability from 1904 (s, a) to s', determined by any perturbed transition kernel  $P_{h,s,a} \in \mathcal{U}^{\sigma^+}(P_{h,s,a}^{f,\phi})$ , as:

$$P_{h}^{\inf,f,\phi}(s' \mid s,a) \coloneqq \inf_{P_{h,s,a} \in \mathcal{U}^{\sigma^{+}}(P_{h,s,a}^{f,\phi})} P_{h}(s' \mid s,a) = \max\{P_{h}(s' \mid s,a) - \sigma^{+}, 0\},$$
(117)

where the last equation follows directly from the definition of  $\mathcal{U}^{\sigma^+}(\cdot)$  in (116), with the remaining probability distributed to other states.

For convenience, we also define the transition from each  $s \in \mathcal{M}$  to the corresponding state  $s^{m \to n} \in \mathcal{N}$  for any  $\mathcal{M}_{f}^{\phi}$ , which is crucial in our analysis: for all  $h \in [H]$ ,

for 
$$m_f$$
:  $p_h^{\inf} \coloneqq P_h^{\inf, f, \phi}(n_f \mid m_f, \phi_h) = p - \sigma^+,$   
 $q_h^{\inf} \coloneqq P_h^{\inf, f, \phi}(n_f \mid m_f, 1 - \phi_h) = q - \sigma^+.$ 
(118)

1917 Then it is obvious that

1898 1899 1900

1905

1907

1918

1919

1922

1932

1933

1935

1938 1939

$$p_1^{\inf} = p_2^{\inf} = \cdots p_H^{\inf}, \quad q_1^{\inf} = q_2^{\inf} = \cdots q_H^{\inf},$$
 (119)

which motivates us to abbreviate them consistently as  $p^{\inf} \coloneqq p_1^{\inf}$  and  $q^{\inf} \coloneqq q_1^{\inf}$  later.

**Robust value functions and optimal policies.** We now define the robust value functions and identify the optimal policies for RMDP instances. For any RMDP  $\mathcal{M}_{f}^{\phi} \in \mathcal{M}(\mathcal{F}, \Phi)$ , let  $\tilde{\mu}^{\star,f,\phi} = \{\mu_{h}^{\star,f,\phi}\}_{h=1}^{H}$  represent the optimal policy, given that  $\nu$  is deterministic. At each step h, we use  $V_{h}^{\tilde{\mu},\sigma^+,f,\phi}$  and  $V_{h}^{\star,\sigma^+,f,\phi}$  to denote the robust value function of any policy  $\tilde{\mu}$  and the optimal policy  $\tilde{\mu}^{\star,f,\phi}$ , respectively, under uncertainty level  $\sigma^+$ . The following lemma highlights key properties of robust value functions and optimal policies; the proof is deferred to Appendix E.1.

1930 **Lemma 8** Consider any  $\mathcal{M}_{f}^{\phi} \in \mathcal{M}(\mathcal{F}, \Phi)$  and any policy  $\tilde{\mu}$ . Defining

$$m_h^{\widetilde{\mu},f,\phi} = p^{\inf}\widetilde{\mu}_h(\phi_h \mid m_f) + q^{\inf}\widetilde{\mu}_h(1 - \phi_h \mid m_f),$$
(120)

1934 it holds that

$$\forall h \in [H]: \quad V_h^{\tilde{\mu}, \sigma^+, f, \phi}(m_f) = m_h^{\tilde{\mu}, f, \phi} V_{h+1}^{\tilde{\mu}, \sigma^+, f, \phi}(n_f) + (1 - m_h^{\tilde{\mu}, f, \phi}) V_{h+1}^{\tilde{\mu}, \sigma^+, f, \phi}(m_f),$$
(121a)

$$\forall (s,h) \in \mathcal{N} \times [H]: \quad V_{h}^{\tilde{\mu},\sigma^{+},f,\phi}(s) = 1 + (1 - \sigma^{+})V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(s) + \sigma^{+}V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f}). \quad (121b)$$

In addition, for all  $h \in [H]$ , the optimal policy and the optimal value function obey 1941

1942 
$$\widetilde{\mu}_{h}^{\star,f,\phi}(\phi_{h} \mid m_{f}) = \widetilde{\mu}_{h}^{\star,f,\phi}(\phi_{h} \mid n_{f}) = 1,$$
 (122)

$$V_{h}^{\star,\sigma^{+},f,\phi}(m_{f}) = p^{\inf} V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(n_{f}) + (1-p^{\inf}) V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f}).$$
(123)

1944 1945 1946 1946 1946 1946 1947 Construction of the history/batch dataset. In the nominal environment  $\mathcal{M}_{f}^{\phi,n}$ , a batch dataset is generated with K independent sample trajectories, each of length H, according to (5) and based on the initial state distribution  $\varrho^{n}$  and behavior policy  $\tilde{\mu}^{n} = {\{\mu_{h}^{n}\}}_{h=1}^{H}$  satisfying

$$\varrho^{\mathsf{n}}(s) = \varrho(s) \quad \text{and} \quad \widetilde{\mu}_{h}^{\mathsf{n}}(a \,|\, s) = \frac{1}{2}, \qquad \forall (s, a, h) \in \mathcal{S}_{\mathsf{one}} \times \mathcal{A}_{\mathsf{one}} \times [H].$$
(124)

We define the nominal transition kernels for  $\mathcal{M}_{f}^{\phi,n}$ , where any state  $m_{i} \in \mathcal{M}$  transitions only to the corresponding  $n_{i} \in \mathcal{N}$  or remains at itself. For simplicity, for any  $s = m_{i} \in \mathcal{M}$ , we denote the corresponding state  $n_{i} \in \mathcal{N}$  as  $s^{m \to n}$ . The basic nominal transition kernel is defined as follows for all  $(h, s, a) \in [H] \times S_{one} \times \mathcal{A}_{one}$ :

$$P_{h}^{\star}(s' \mid s, a) = \begin{cases} (p + \Delta)\mathbb{1}(s' = s^{m \to n}) + (1 - p - \Delta)\mathbb{1}(s' = s) & \text{if } s \in \mathcal{M}, a = \phi_{h} \\ p\mathbb{1}(s' = s^{m \to n}) + (1 - p)\mathbb{1}(s' = s) & \text{if } s \in \mathcal{M}, a = 1 - \phi_{h} \\ \mathbb{1}(s' = s) & \text{if } s \in \mathcal{N}. \end{cases}$$
(125)

In words, the transition kernel of each  $\mathcal{M}_{f}^{\phi} \in \mathcal{M}(\mathcal{F}, \Phi)$  only differs slightly from the basic nominal transition kernel  $\mathcal{M}_{f}^{\phi,n}$  when  $s = m_{f}$ , which makes all the components within  $\mathcal{M}(\mathcal{F}, \Phi)$  close to each other.

1963 Specifically, p and q are set according to

$$0 \le p \le p + \Delta \le 1$$
 and  $0 \le q = p - \Delta$  (126)

for some p and  $\Delta > 0$ . Without loss of generality, let the uncertainty level be  $\sigma^+ \in (0, 1 - c_0]$  for some  $0 < c_0 < 1$ . Then taking  $c_2 \le \frac{1}{4}$  and  $c_1 \coloneqq \frac{c_0}{2} \le \frac{1}{4}$ , p and  $\Delta$  are set as

$$p = \begin{cases} \frac{c_2}{H}, & \text{if } \sigma^+ \le \frac{c_2}{2H} \\ \left(1 + \frac{c_1}{H}\right)\sigma^+ & \text{otherwise} \end{cases} \quad \text{and} \quad \Delta \le \begin{cases} \frac{c_2}{2H}, & \text{if } \sigma^+ \le \frac{c_2}{2H} \\ \frac{c_1}{H}\sigma^+ & \text{otherwise} \end{cases}$$
(127)

which establishes the fact that

$$+\Delta \ge p \ge q = p - \Delta \ge \max\left\{\frac{c_2}{2H}, \sigma^+\right\}.$$
(128)

1975 Combined with  $H \ge 2$ , it is easily verified that  $0 \le p + \Delta \le 1$  as follows:

p

when 
$$\sigma^{+} > \frac{c_{2}}{2H}$$
:  $\left(1 + \frac{c_{1}}{H}\right)\sigma^{+} + \frac{c_{1}}{H}\sigma^{+} \le 1 - c_{0} + \frac{2c_{1}}{H}\sigma^{+} \le 1 - \frac{c_{0}(H-1)}{H} < 1,$   
when  $\sigma^{+} \le \frac{c_{2}}{2H}$ :  $\frac{3c_{2}}{2H} \le 1.$  (129)

1979 1980 1981

1984

1988 1989 1990

1948 1949 1950

1965

1969 1970 1971

1972 1973

1974

1977 1978

In addition, let  $\overline{\varrho}(s)$  represents a state distribution supported on the state subset  $(m_f, n_f) \in \mathcal{M} \times \mathcal{N}$ :

$$\overline{\varrho}(s) = \frac{1}{CSA} \mathbb{1}(s = m_f) + \left(1 - \frac{1}{CSA}\right) \mathbb{1}(s = n_f),$$
(130)

where  $1(\cdot)$  is the indicator function, and C > 0 is some constant that will determine the concentrability coefficient  $C_r^*$  (as we shall detail momentarily) and obeys

$$\frac{1}{CSA} \le \frac{1}{4}.\tag{131}$$

As it turns out, for any MDP  $\mathcal{M}_{\phi}^{f}$ , the occupancy distributions of the above batch dataset are the same (due to symmetry) and admit the following simple characterization:

1993  
1994 
$$\forall (s,a) \in \mathcal{S}_{one} \times \mathcal{A}_{one}, \qquad \qquad d_1^{\mathsf{n}, P^{\phi, f}}(s,a) = \frac{1}{2}\overline{\varrho}(s), \qquad (132a)$$
1995

$$\begin{array}{l} \text{1996} \quad \forall (s,a,h) \in \mathcal{S}_{\text{one}} \times \mathcal{A}_{\text{one}} \times [H], \qquad \frac{\overline{\varrho}(s)}{2} \leq d_h^{\mathfrak{n},P^{\phi,f}}(s) \leq 2\overline{\varrho}(s), \quad \frac{\overline{\varrho}(s)}{4} \leq d_h^{\mathfrak{n},P^{\phi,f}}(s,a) \leq \overline{\varrho}(s). \end{array}$$

$$\begin{array}{l} \text{(132b)} \end{array}$$

In addition, we choose the following initial state distribution 

$$\varrho(s) = \begin{cases} \frac{1}{CSA}, & \text{if } s \in \mathcal{M} \\ 0, & \text{if } s \in \mathcal{N}. \end{cases}$$
(133)

With this choice of  $\rho$ , the single-policy clipped concentrability coefficient  $C_r^{\star}$  and the quantity C are intimately connected as follows: 

$$C \le C_{\rm r}^{\star} \le 2C. \tag{134}$$

The proof of the claim (132) and (134) are postponed to Appendix E.2. 

#### D.2 STEP 2: ESTABLISHING THE MINIMAX LOWER BOUND

(f

Recall our goal: for any policy estimator  $\tilde{\mu}$  computed based on the empirical dataset, we plan to control the quantity 

$$\max_{\phi)\in\mathcal{F}\times\Phi}\left\{V_{1}^{\star,\sigma^{+},f,\phi}(\varrho)-V_{1}^{\tilde{\mu},\sigma^{+},f,\phi}(\varrho)\right\}$$
(135)

with initial state distribution defined in (133).

**Step 1: converting the goal to estimate**  $(f, \phi)$ . Towards this, we make the following essential claim which shall be verified in Appendix E.3: letting 

$$\varepsilon \le \begin{cases} \frac{c_2}{H}, & \text{if } \sigma^+ \le \frac{c_2}{2H} \\ 1 & \text{otherwise} \end{cases}$$
(136)

and 

$$\Delta = c_5 \begin{cases} \frac{\varepsilon}{H^2}, & \text{if } \sigma^+ \le \frac{c_2}{2H} \\ \frac{\sigma^+ \varepsilon}{H} & \text{otherwise} \end{cases}$$
(137)

which satisfies (127), it leads to that for any policy  $\tilde{\mu}$  obeying 

$$\sum_{h=1}^{H} \left\| \tilde{\mu}_{h}(\cdot \,|\, m_{f}) - \tilde{\mu}_{h}^{\star, f, \phi}(\cdot \,|\, m_{f}) \right\|_{1} \ge \frac{H}{8},$$
(138)

one has 

$$V_1^{\star,\sigma^+,f,\phi}(m_f) - V_1^{\tilde{\mu},\sigma^+,f,\phi}(m_f) > \varepsilon.$$
 (139)

We are now ready to convert the task of estimating an optimal policy to estimating  $(f, \phi)$ . For this, let  $\mathbb{P}_{f,\phi}$  represent the probability distribution when the RMDP is  $\mathcal{M}_{f}^{\phi}$  for any  $(f,\phi) \in \mathcal{F} \times \Phi$ . Then, for any  $(f, \phi) \in \mathcal{F} \times \Phi$ , suppose that there exists a policy  $\tilde{\mu}$  achieving

$$\mathbb{P}_{f,\phi}\left\{V_1^{\star,\sigma^+,f,\phi}(m_f) - V_1^{\widetilde{\mu},\sigma^+,f,\phi}(m_f) \le \varepsilon\right\} \ge \frac{3}{4},\tag{140}$$

which in view of (139) indicates that we necessarily have

$$\mathbb{P}_{f,\phi}\left\{\sum_{h=1}^{H} \left\| \widetilde{\mu}_{h}(\cdot \mid m_{f}) - \widetilde{\mu}_{h}^{\star,f,\phi}(\cdot \mid m_{f}) \right\|_{1} < \frac{H}{8} \right\} \ge \frac{3}{4}.$$
(141)

Consequently, taking  $\widetilde{\phi} = \arg\min_{\phi \in \Phi} \sum_{h=1}^{H} \left\| \widetilde{\mu}_{h}(\cdot \mid m_{f}) - \widetilde{\mu}_{h}^{\star,f,\phi}(\cdot \mid m_{f}) \right\|_{1}$ , we are motivated to construct the estimate of  $\phi$  as  $\hat{\phi} = \tilde{\phi}$ . Namely, if  $\sum_{h=1}^{H} \left\| \widetilde{\mu}_h(\cdot \mid m_f) - \widetilde{\mu}_h^{\star,f,\phi}(\cdot \mid m_f) \right\|_1 < \frac{H}{8}$  holds for some  $\phi \in \Phi$ , then for any  $\phi' \in \Phi$  obeying  $\phi' \neq \phi$ , one has 

2045  
2046  
2047  
2048  
2049  
2050  
2051  

$$\sum_{h=1}^{H} \|\widetilde{\mu}_{h}(\cdot | m_{f}) - \widetilde{\mu}_{h}^{\star,f,\phi'}(\cdot | m_{f})\|_{1}$$

$$\sum_{h=1}^{H} \|\widetilde{\mu}_{h}(\cdot | m_{f}) - \widetilde{\mu}_{h}^{\star,f,\phi}(\cdot | m_{f})\|_{1}$$

 where the first inequality holds by the triangle inequality, and the last inequality follows from the assumption  $\sum_{h=1}^{H} \|\tilde{\mu}_h(\cdot | m_f) - \tilde{\mu}_h^{\star,f,\phi}(\cdot | m_f)\|_1 < \frac{H}{8}$  and the separation property of  $\phi \in \Phi$  (see (111)). Similarly, it shows that we have  $\hat{\phi} = \phi$  if

$$\sum_{h=1}^{H} \left\| \widetilde{\mu}_{h}(\cdot \mid m_{f}) - \widetilde{\mu}_{h}^{\star,f,\phi}(\cdot \mid m_{f}) \right\|_{1} < \frac{H}{8} < \sum_{h=1}^{H} \left\| \widetilde{\mu}_{h}(\cdot \mid m_{f}) - \widetilde{\mu}_{h}^{\star,f,\phi'}(\cdot \mid m_{f}) \right\|_{1}$$
(143)

holds for all  $\phi' \in \Phi$  that  $\phi' \neq \phi$ . It is clear that the above equation can be directly achieved when  $\sum_{h=1}^{H} \|\widetilde{\mu}_h(\cdot \mid m_f) - \widetilde{\mu}_h^{\star,f,\phi}(\cdot \mid m_f)\|_1 < \frac{H}{8}$ , which gives

$$\mathbb{P}_{f,\phi}\left[\widehat{\phi}=\phi\right] \ge \mathbb{P}_{f,\phi}\left\{\sum_{h=1}^{H}\left\|\widetilde{\mu}_{h}(\cdot\mid m_{f})-\widetilde{\mu}_{h}^{\star,f,\phi}(\cdot\mid m_{f})\right\|_{1}<\frac{H}{8}\right\} \ge \frac{3}{4}.$$
(144)

Step 2: developing the probability of error in testing multiple hypotheses. Next, we address the hypothesis testing problem over  $\phi \in \Phi$  and derive the information-theoretic lower bound for the probability of error. Specifically, we define the minimax probability of error as:

$$p_{\mathbf{e}} \coloneqq \inf_{(\widehat{f},\widehat{\phi})} \max_{(f,\phi)\in\mathcal{F}\times\Phi} \mathbb{P}_{f,\phi}\left(\widehat{\phi}\neq\phi\right)$$

where the infimum is taken over all possible tests  $\hat{\phi}$  constructed from the available batch dataset.

Given the dataset  $\mathcal{D}_0$  with K independent trajectories, let  $\varrho^{n,\phi}$  (and  $\varrho_h^{n,\phi}(s,a)$ ) represent the distribution vector (and distribution) of each sample tuple  $(s_h, a_h, s'_h)$  at time step h under the nominal transition kernel  $P^*$  for  $\mathcal{M}_f^{\phi,n}$ . Using this, along with Fano's inequality (Tsybakov, 2008, Theorem 2.2) and the additivity of KL divergence (Tsybakov, 2008, Page 85), we derive the following result:

$$p_{e} \geq 1 - K \frac{\max_{(\phi,\tilde{\phi})\in\Phi,\phi\neq\tilde{\phi}}\mathsf{KL}(\varrho^{\mathsf{n},\phi} \mid \varrho^{\mathsf{n},\phi}) + \log 2}{\log |\Phi|}$$

$$\stackrel{(i)}{\geq} 1 - \frac{8K}{H} \max_{(\phi,\tilde{\phi})\in\Phi,\phi\neq\tilde{\phi}}\mathsf{KL}(\varrho^{\mathsf{n},\phi} \mid \varrho^{\mathsf{n},\tilde{\phi}}) - \frac{8\log 2}{H}$$

(145)

2056 2057 2058

2059

2065

2066

2067 2068

2082 2083

2096 2097

2100

where (i) holds by  $|\Phi| \ge e^{H/8}$  and (ii) follows from  $H \ge 16 \log 2$ .

тт

Since the occupancy state distribution  $d_h^n$  is the same for any MDP  $\mathcal{M}_f^{\phi}$  for  $\phi \in \Phi$ , we apply the chain rule of KL divergence (Duchi, 2018, Lemma 5.2.8) and the Markov property of the independent sample trajectories to obtain:

 $\stackrel{(\mathrm{ii})}{\geq} \frac{1}{2} - \frac{8K}{H} \max_{(\phi, \widetilde{\phi}) \in \Phi, \phi \neq \widetilde{\phi}} \mathsf{KL}\big(\varrho^{\mathsf{n}, \phi} \,|\, \varrho^{\mathsf{n}, \widetilde{\phi}}\big),$ 

$$\mathsf{KL}(\varrho^{\mathsf{n},\phi} | \varrho^{\mathsf{n},\widetilde{\phi}}) = \sum_{h=1}^{n} \mathop{\mathbb{E}}_{s \sim d_{h}^{\mathsf{n}}(s)} \left[ \mathsf{KL}(P_{h}^{\star,\phi}(\cdot | s, a) \parallel P_{h}^{\star,\widetilde{\phi}}(\cdot | s, a)) \right]$$
$$\stackrel{(i)}{=} \frac{1}{2} \overline{\varrho}(m_{f}) \sum_{h=1}^{H} \sum_{a \in \{0,1\}} \left[ \mathsf{KL}(P_{h}^{\phi}(\cdot | m_{f}, a) \parallel P_{h}^{\widetilde{\phi}}(\cdot | m_{f}, a)) \right], \quad (146)$$

where (i) follows from applying (132) and obtaining the fact as

$$\mathbb{E}_{s \sim d_h^{\mathfrak{n}}(s)} \left[ \mathsf{KL} \left( P_h^{\star,\phi}(\cdot \mid s, a) \parallel P_h^{\star,\widetilde{\phi}}(\cdot \mid s, a) \right) \right] \\ - \sum_{s \sim d_h^{\mathfrak{n}}(s)} \int_{\sum \widetilde{u}_h^{\mathfrak{n}}(a \mid s) P_h^{\phi_h}(s' \mid s, a) \log \widetilde{\mu}_h^{\mathfrak{n}}(a \mid s) P_h^{\phi_h}}$$

$$= \sum_{s} d_h^{\mathbf{n}}(s) \left\{ \sum_{a,s'} \widetilde{\mu}_h^{\mathbf{n}}(a \,|\, s) P_h^{\phi_h}(s' \,|\, s, a) \log \frac{\widetilde{\mu}_h^{\mathbf{n}}(a \,|\, s) P_h^{\phi_h}(s' \,|\, s, a)}{\widetilde{\mu}_h^{\mathbf{n}}(a \,|\, s) P_h^{\widetilde{\phi}_h}(s' \,|\, s, a)} \right\}$$

$$= \frac{1}{2} \overline{\varrho}(m_f) \sum \sum P_h^{\phi_h}(s' \mid m_f, a) \log \frac{P_h^{\phi_h}(s' \mid m_f, a)}{P_h^{\phi_h}(s' \mid m_f, a)}$$

$$= \frac{1}{2}\overline{\varrho}(m_f) \sum_{a} \mathsf{KL}\left(P_h^{\phi_h}(\cdot \mid m_f, a) \parallel P_h^{\phi_h}(\cdot \mid m_f, a)\right).$$

2106 Consequently, combining (145) and (146) leads to

a

 $p_{\mathbf{e}} \geq \frac{1}{2} - \frac{4K}{H} \max_{(+,\widetilde{\mathbf{x}}) \in \mathbf{\Phi}, + <\widetilde{\mathbf{x}}} \left[ \overline{\varrho}(m_f) \sum_{h=1}^{H} \sum_{\mathbf{k} \in \mathbf{K}} \mathsf{KL} \left( P_h^{\phi_h}(\cdot \mid m_f, a) \parallel P_h^{\widetilde{\phi}_h}(\cdot \mid m_f, a) \right) \right].$ 

Thus, we turn to focus on terms in (147) now in different cases of the uncertainty level  $\sigma^+$ .

• For 
$$0 < \sigma^+ \leq \frac{c_2}{2H}$$
: If  $\phi_h = \widetilde{\phi}_h$ , it is obvious that

$$\sum_{a \in \{0,1\}} \mathsf{KL}\big(P_h^{\star,\phi}(\cdot \,|\, s,a) \parallel P_h^{\star,\widetilde{\phi}}(\cdot \,|\, s,a)\big) = 0.$$
(148)

(147)

Therefore, we consider the case of  $\phi_h \neq \tilde{\phi}_h$ . Without loss of generality, we suppose  $\phi_h = 0$  and  $\tilde{\phi}_h = 1$ , which indicates

$$\mathsf{KL}(P_{h}^{\star,\phi}(0 \mid m_{f}, 0) \parallel P_{h}^{\star,\widetilde{\phi}}(0 \mid m_{f}, 0)) \leq \frac{(p-q)^{2}}{q(1-q)} \stackrel{(i)}{=} \frac{\Delta^{2}}{q(1-q)}$$
$$\stackrel{(ii)}{=} \frac{(c_{5})^{2}\varepsilon^{2}}{H^{4}q(1-q)} \leq \frac{4(c_{5})^{2}\varepsilon^{2}}{c_{2}H^{3}}, \tag{149}$$

where the first inequality exists by applying Lemma 7, (i) follows from the definitions in (126), (ii) holds due to the definition in (137), and the last inequality arises from  $q = p - \Delta \ge \frac{c_2}{2H}$  (see (127)) and  $1 - q \ge 1 - p \ge 1 - \frac{c_2}{H} \ge \frac{1}{2}$ .

Similarly, we can establish the same bound for  $\mathsf{KL}(P_h^{\star,\phi}(0 \mid m_f, 1) \parallel P_h^{\star,\widetilde{\phi}}(0 \mid m_f, 1))$ . Summing up the results with the fact in (149), we arrive at

$$\sum_{e \in \{0,1\}} \mathsf{KL} \left( P_h^{\star,\phi}(\cdot \mid m_f, a) \parallel P_h^{\star,\widetilde{\phi}}(\cdot \mid m_f, a) \right) \le \frac{16(c_5)^2 \varepsilon^2}{c_2 H^3}.$$
 (150)

• For  $\frac{c_2}{2H} < \sigma^+ \le 1 - c_0$ : Following the same pipeline, it then boils down to control the main term as below:

$$\mathsf{KL}(P_{h}^{\star,\phi}(0 \mid m_{f}, 0) \parallel P_{h}^{\star,\widetilde{\phi}}(0 \mid m_{f}, 0)) \leq \frac{(p-q)^{2}}{q(1-q)} \stackrel{(\mathrm{i})}{=} \frac{\Delta^{2}}{q(1-q)}$$
$$\stackrel{(\mathrm{ii})}{=} \frac{(c_{5})^{2} \sigma^{+2} \varepsilon^{2}}{H^{2} q(1-q)} \leq \frac{2(c_{5})^{2} \sigma^{+} \varepsilon^{2}}{c_{0} H^{2}}, \qquad (151)$$

where (i) and (ii) follow from the definitions in (126) or (137). Here, the last inequality arises from

$$1 - q \ge 1 - p = 1 - (1 + \frac{c_1}{H})\sigma^+ \stackrel{(i)}{\ge} c_0 - \frac{c_1}{H} \stackrel{(ii)}{\ge} \frac{c_0}{2}$$
$$p \ge q = p - \Delta \stackrel{(iii)}{\ge} \sigma^+, \tag{152}$$

where (ii) holds by the definition of  $c_1 = \frac{c_0}{2}$ , and (iii) follows from (128). Consequently, we arrive at

$$\sum_{a \in \{0,1\}} \mathsf{KL}\big(P_h^{\star,\phi}(\cdot \mid s,a) \parallel P_h^{\star,\widetilde{\phi}}(\cdot \mid s,a)\big) \le \frac{8(c_5)^2 \sigma^+ \varepsilon^2}{c_0 H^2}.$$
(153)

Summing up (150) and (153), we achieve for any  $(\phi, \tilde{\phi}) \in \Phi$  with  $\phi \neq \tilde{\phi}$  and any time step  $h \in [H]$ 

$$\sum_{a \in \{0,1\}} \mathsf{KL}\big(P_h^{\star,\phi}(\cdot \mid m_f, a) \parallel P_h^{\star,\widetilde{\phi}}(\cdot \mid m_f, a)\big) \le \frac{16(c_5)^2 \varepsilon^2}{c_0 c_2 H^2} \max\{\sigma^+, 1/H\}.$$
(154)

Plugging (154) back to (147), under the definition in (133), we obtain

$$p_{\mathbf{e}} \geq \frac{1}{2} - \frac{4K}{H} \max_{(\phi, \tilde{\phi}) \in \Phi, \phi \neq \tilde{\phi}} \left[ \overline{\varrho}(m_f) \sum_{h=1}^{H} \sum_{a} \mathsf{KL} \left( P_h^{\phi_h}(\cdot \mid m_f, a) \parallel P_h^{\tilde{\phi}_h}(\cdot \mid m_f, a) \right) \right]$$

$$\geq \frac{1}{2} - \frac{4K}{H}\overline{\varrho}(m_f) \sum_{h=1}^{H} \frac{16(c_5)^2 \varepsilon^2}{c_0 c_2 H^2} \max\{\sigma^+, 1/H\}$$

$$\geq \frac{1}{2} - \frac{64K(c_5)^2 \varepsilon^2}{c_0 c_2 CSAH^2} \max\{\sigma^+, 1/H\} \geq \frac{1}{4},\tag{155}$$

2171 as long as the sample size T = KH of the dataset is selected as

$$T \le \frac{c_0 c_2 CSAH^3 \min\{1/\sigma^+, H\}}{256(c_5)^2 \varepsilon^2} \le \frac{c_0 c_2 C_{\rm r}^{\star} SAH^3 \min\{1/\sigma^+, H\}}{256(c_5)^2 \varepsilon^2}.$$
 (156)

**Step 3: summing up the results together.** We suppose that there exists an estimator  $\tilde{\mu}$  such that

$$\max_{(f,\phi\in\mathcal{F})\times\Phi} \mathbb{P}_{f,\phi}\left[\left\{V_1^{\star,\sigma^+,f,\phi}(\varrho) - V_1^{\widetilde{\mu},\sigma^+,f,\phi}(\varrho)\right\} \ge \varepsilon\right] < \frac{1}{4}.$$
(157)

2180 Then according to (135), we need

$$\forall w \in \mathcal{F}: \quad \max_{\phi \in \Phi} \mathbb{P}_{f,\phi} \left[ \left\{ V_1^{\star,\sigma^+,f,\phi}(m_f) - V_1^{\widetilde{\mu},\sigma^+,f,\phi}(m_f) \right\} \ge \varepsilon \right] < \frac{1}{4}.$$
(158)

To meet (158) for any  $w \in \mathcal{F}$ , we require

$$\forall \phi \in \Phi : \mathbb{P}_{f,\phi}\left\{V_1^{\star,\sigma^+,f,\phi}(m_f) - V_1^{\widetilde{\mu},\sigma^+,f,\phi}(m_f) < \varepsilon\right\} \ge \frac{3}{4},\tag{159}$$

which in view of (139) indicates that we necessarily have

$$\forall \phi \in \Phi : \quad \mathbb{P}_{f,\phi} \left\{ \sum_{h=1}^{H} \left\| \widetilde{\mu}_h(\cdot \mid m_f) - \widetilde{\mu}_h^{\star,f,\phi}(\cdot \mid m_f) \right\|_1 < \frac{H}{8} \right\} \ge \frac{3}{4}. \tag{160}$$

As a consequence, (144) indicates 2195

$$\forall \phi \in \Phi : \mathbb{P}_{f,\phi} \left[ \widehat{\phi} = \phi \right] \ge \frac{3}{4}.$$
(161)

To achieve (157), we here apply the fact in (161) to all  $w \in \mathcal{F}$ , which leads to the fact that one necessarily has

$$\forall (f,\phi) \in \mathcal{F} \times \Phi : \quad \mathbb{P}_{f,\phi} \left[ (\widehat{f},\widehat{\phi}) = (f,\phi) \right] \ge \frac{3}{4}.$$
(162)

However, this would contract with (155) as long as the sample size condition in (156) is satisfied. Thus, if the sample size obeys the condition (156), we can't achieve an estimate  $\tilde{\mu}$  that satisfies (157), which completes the proof.

## E AUXILIARY FACTS FOR THEOREM 2

2211 E.1 PROOF OF LEMMA 8

Since all RMDPs in  $\mathcal{M}(\mathcal{F}, \Phi)$  are constructed similarly for each  $w \in \mathcal{F}$  and  $\phi \in \Phi$ , we will focus on a specific RMDP  $\mathcal{M}_{f}^{\phi} \in \mathcal{M}(\mathcal{F}, \Phi)$ , with the results applicable to all other RMDPs in  $\mathcal{M}(\mathcal{F}, \Phi)$ .

 Part 1: ordering the robust value function over different states. Before proceeding, we introduce several facts and notations that will be useful throughout this section. First, for any  $\mathcal{M}_{f}^{\phi}$ and any policy  $\tilde{\mu}$ , we observe the following at the final step H + 1:

$$\forall s \in \mathcal{M} \cup \mathcal{N}: \quad V_{H+1}^{\widetilde{\mu}, \sigma^+, f, \phi}(s) = 0.$$
(163)

Then for step H, we can easily verify that

$$\forall s \in \mathcal{N}: \quad V_{H}^{\widetilde{\mu},\sigma^{+},f,\phi}(s) = \mathbb{E}_{a \sim \widetilde{\mu}_{H}(\cdot \mid s)} \left[ r_{H}(s,a) + \inf_{\substack{\mathcal{P} \in \mathcal{U}^{\sigma^{+}}(P_{H,s,a}^{f,\phi})}} \mathcal{P} V_{H+1}^{\widetilde{\mu},\sigma^{+},f,\phi} \right] = 1 \quad (164a)$$

$$\forall s \in \mathcal{M}: \quad V_{H}^{\widetilde{\mu},\sigma^{+},f,\phi}(s) = \mathbb{E}_{a \sim \widetilde{\mu}_{H}(\cdot \mid s)} \left[ r_{H}(s,a) + \inf_{\mathcal{P} \in \mathcal{U}^{\sigma^{+}}(P_{H,s,a}^{f,\phi})} \mathcal{P}V_{H+1}^{\widetilde{\mu},\sigma^{+},f,\phi} \right] = 0, \quad (164b)$$

which holds by (163) and the definition of the reward function (see (114)). The above fact directly indicates that

$$\forall (s,s') \in \mathcal{M} \times \mathcal{N} : \min_{\widetilde{s} \in \mathcal{S}} V_H^{\widetilde{\mu},\sigma^+,f,\phi}(\widetilde{s}) = V_H^{\widetilde{\mu},\sigma^+,f,\phi}(m_f) \le V_H^{\widetilde{\mu},\sigma^+,f,\phi}(s) < V_H^{\widetilde{\mu},\sigma^+,f,\phi}(s'),$$
(165a)

$$\forall (s,s') \in \mathcal{N} \times \mathcal{N} : \quad V_H^{\widetilde{\mu},\sigma^+,f,\phi}(s) = V_H^{\widetilde{\mu},\sigma^+,f,\phi}(s'). \tag{165b}$$

Then we introduce a claim which we will prove by induction in a moment as below:

$$\forall (h, s, s') \in [H] \times \mathcal{M} \times \mathcal{N} : \quad V_h^{\widetilde{\mu}, \sigma^+, f, \phi}(m_f) \le V_h^{\widetilde{\mu}, \sigma^+, f, \phi}(s) < V_h^{\widetilde{\mu}, \sigma^+, f, \phi}(s')$$
(166a)

$$\forall (s,s') \in \mathcal{N} \times \mathcal{N} : \quad V_h^{\widetilde{\mu},\sigma^+,f,\phi}(s) = V_h^{\widetilde{\mu},\sigma^+,f,\phi}(s').$$
(166b)

Note that the base case when the time step is H + 1 is verified in (165). Assume that the following fact at time step h + 1 holds

$$\forall (s,s') \in \mathcal{M} \times \mathcal{N} : \quad \min_{\widetilde{s} \in \mathcal{S}} V_{h+1}^{\widetilde{\mu},\sigma^+,f,\phi}(\widetilde{s}) = V_{h+1}^{\widetilde{\mu},\sigma^+,f,\phi}(m_f) \le V_{h+1}^{\widetilde{\mu},\sigma^+,f,\phi}(s) < V_{h+1}^{\widetilde{\mu},\sigma^+,f,\phi}(s'),$$
(167a)

$$\forall (s,s') \in \mathcal{N} \times \mathcal{N} : \quad V_{h+1}^{\widetilde{\mu},\sigma^+,f,\phi}(s) = V_{h+1}^{\widetilde{\mu},\sigma^+,f,\phi}(s'). \tag{167b}$$

Therefore, the rest of the proof focuses on proving the same property for time step h. For RMDP  $\mathcal{M}_{f}^{\phi} \in \mathcal{M}(\mathcal{F}, \Phi)$  and any policy  $\tilde{\mu}$ , we characterize the robust value function of different states separately:

• For state  $s \in \mathcal{N}$ : we observe that for any  $s \in \mathcal{N}$ ,

$$V_{h}^{\tilde{\mu},\sigma^{+},f,\phi}(s) = \mathbb{E}_{a\sim\tilde{\mu}_{h}(\cdot\mid s)} \left[ r_{h}(s,a) + \inf_{P\in\mathcal{U}^{\sigma^{+}}(P_{h,s,a}^{f,\phi})} PV_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi} \right]$$

$$\stackrel{(i)}{=} 1 + \mathbb{E}_{a\sim\tilde{\mu}_{h}(\cdot\mid s)} \left[ P_{h}^{\inf,f,\phi}(s\mid s,a)V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(s) \right] + \sigma^{+}V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f})$$

$$= 1 + (1 - \sigma^{+})V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(s) + \sigma^{+}V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f}), \qquad (168)$$

where (i) holds by  $r_h(s,a) = 1$  for all  $s \in \mathcal{N}$  (see (114)), the fact that  $\min_{\tilde{s}\in \mathcal{S}} V_{h+1}^{\tilde{\mu},\sigma^+,f,\phi}(\tilde{s}) = V_{h+1}^{\tilde{\mu},\sigma^+,f,\phi}(m_f)$  induced by the induction assumption (cf. (167)) and the definition of  $P_h^{\inf,f,\phi}(s \mid s, a)$  in (117), and the last equality follows from  $P^{f,\phi}(s \mid s, a) = 1$  for all  $(s, a) \in \mathcal{N} \times \mathcal{A}_{one}$ . Resorting to the induction assumption in (167), we have

$$\forall (s,s') \in \mathcal{N} \times \mathcal{N} : \quad V_h^{\widetilde{\mu},\sigma^+,f,\phi}(s) = V_h^{\widetilde{\mu},\sigma^+,f,\phi}(s'). \tag{169}$$

• For state  $m_f$ : first, the robust value function at state  $m_f$  obeys

 $V^{\widetilde{\mu},\sigma^+,f,\phi}(m_{\tau})$ 

$$= \mathbb{E}_{a \sim \widetilde{\mu}_{h}(\cdot \mid m_{f})} \left[ r_{h}(m_{f}, a) + \inf_{\substack{P \in \mathcal{U}^{\sigma^{+}}(P_{h, m_{f}, a}^{f, \phi})}} PV_{h+1}^{\widetilde{\mu}, \sigma^{+}, f, \phi} \right]$$
$$\stackrel{(i)}{=} 0 + \widetilde{\mu}_{h}(\phi_{h} \mid m_{f}) \inf_{\substack{P \in \mathcal{U}^{\sigma^{+}}(P_{h, m_{f}, \phi_{h}}^{f, \phi})}} PV_{h+1}^{\widetilde{\mu}, \sigma^{+}, f, \phi}$$

$$+ \widetilde{\mu}_{h}(1 - \phi_{h} \mid m_{f}) \inf_{P \in \mathcal{U}^{\sigma^{+}}(P_{h,m_{f},1-\phi_{h}}^{f,\phi})} PV_{h+1}^{\widetilde{\mu},\sigma^{+},f,\phi}$$

$$\stackrel{\text{(ii)}}{=} \widetilde{\mu}_{h}(\phi_{h} \mid m_{f}) \Big[ p^{\inf}V_{h+1}^{\widetilde{\mu},\sigma^{+},f,\phi}(n_{f}) + (1 - p^{\inf}) V_{h+1}^{\widetilde{\mu},\sigma^{+},f,\phi}(m_{f}) \Big]$$

$$+ \widetilde{\mu}_{h}(1 - \phi_{h} \mid m_{f}) \Big[ q^{\inf}V_{h+1}^{\widetilde{\mu},\sigma^{+},f,\phi}(n_{f}) + (1 - q^{\inf}) V_{h+1}^{\widetilde{\mu},\sigma^{+},f,\phi}(m_{f}) \Big]$$

$$\stackrel{\text{(iii)}}{=} m_h^{\tilde{\mu}, f, \phi} V_{h+1}^{\tilde{\mu}, \sigma^+, f, \phi}(n_f) + (1 - m_h^{\tilde{\mu}, f, \phi}) V_{h+1}^{\tilde{\mu}, \sigma^+, f, \phi}(m_f)$$
(170)

$$\leq (1 - \sigma^{+}) V_{h+1}^{\tilde{\mu}, \sigma^{+}, f, \phi}(n_{f}) + \sigma^{+} V_{h+1}^{\tilde{\mu}, \sigma^{+}, f, \phi}(m_{f}).$$
(171)

where (i) uses the definition of the robust value function and the reward function in (114), (ii) uses the induction assumption in (167) so that the minimum is attained by picking the choice specified in (118) to absorb probability mass to state  $m_f$ , and (iii) holds by plugging in the definition (120) of  $m_h^{\tilde{\mu},f,\phi}$ . Finally, the last inequality follows from the fact that function  $f(m) := mV_{h+1}^{\tilde{\mu},\sigma^+,f,\phi}(n_f) + (1-m)V_{h+1}^{\tilde{\mu},\sigma^+,f,\phi}(m_f)$  is monotonically increasing with m since  $V_{h+1}^{\tilde{\mu},\sigma^+,f,\phi}(n_f) > V_{h+1}^{\tilde{\mu},\sigma^+,f,\phi}(m_f)$  (see the induction assumption (167)), and the fact  $m_h^{\tilde{\mu},f,\phi} \leq 1 - \sigma^+$ .

Combining the above results with (169), we confirm the claim in (166).

**Part 2: deriving the optimal policy and optimal robust value function.** We shall characterize the optimal policy and corresponding optimal robust value function for different states separately:

• For states in M: Recall (170)

$$V_{h}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f}) = m_{h}^{\tilde{\mu},f,\phi} V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(n_{f}) + (1 - m_{h}^{\tilde{\mu},f,\phi}) V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f})$$
(172)

and the fact  $V_{h+1}^{\tilde{\mu},\sigma^+,f,\phi}(n_f) > V_{h+1}^{\tilde{\mu},\sigma^+,f,\phi}(m_f)$  in (166). We observe that (172) is monotonicity increasing with respect to  $m_h^{\tilde{\mu},f,\phi}$ , and  $m_h^{\tilde{\mu},f,\phi}$  is also increasing in  $\tilde{\mu}_h(\phi_h \mid m_f)$  (refer to the fact  $p^{\inf} \ge q^{\inf}$  since  $p \ge q$ ; see (126) and (118)). Consequently, the optimal policy and optimal robust value function in state  $m_f$  thus obey

$$\forall h \in [H]: \quad \tilde{\mu}_{h}^{\star,f,\phi}(\phi_{h} \mid m_{f}) = 1, \\ V_{h}^{\star,\sigma^{+},f,\phi}(m_{f}) = p^{\inf} V_{h+1}^{\star,\sigma^{+},f,\phi}(n_{f}) + (1 - p^{\inf}) V_{h+1}^{\star,\sigma^{+},f,\phi}(m_{f}). \quad (173)$$

For states s ∈ N: Recall the transitions in (125) and (113). Considering that the action does not influence the state transition for all states s ∈ N, without loss of generality, we choose the robust optimal policy obeying

$$\forall s \in \mathcal{N}: \quad \tilde{\mu}_h^{\star, f, \phi}(\phi_h \,|\, s) = 1. \tag{174}$$

2317 E.2 PROOF OF CLAIM (132) AND (134)

Proof of the claim (132). With the initial state distribution and behavior policy defined in (124), we have for any MDP  $\mathcal{M}_{\phi}^{f}$ ,

$$d_1^{\mathsf{n},P^{\phi,f}}(s) = \varrho^\mathsf{n}(s) = \overline{\varrho}(s),$$

which leads to

2336

2338 2339

2342 2343

2346 2347

2350

2354 2355

2357 2358 2359

2360

2361

$$\forall (m_f, a) \in \mathcal{M} \times \mathcal{A}_{one} : \quad d_1^{\mathsf{n}, P^{\phi, f}}(m_f, a) = \frac{1}{2}\overline{\varrho}(m_f).$$
(175)

Along with  $d_1^{\mathbf{n},P^{\phi,f}}(n_f,a) = \frac{1}{2}\overline{\varrho}(n_f) = 0$ , the claim (132a) is proved.

In view of (125), the state occupancy distribution at any step  $h = 2, 3, \dots, H$  obeys

$$d_{h}^{\mathbf{n},P^{\phi,f}}(m_{f}) \geq \mathbb{P}\left\{s_{h} = s' \mid s_{h-1} = m_{f}; \widetilde{\mu}^{\mathbf{n}}\right\}$$
  
$$\geq d_{h-1}^{\mathbf{n},P^{\phi,f}}(m_{f}) \left[\widetilde{\mu}_{h-1}^{\mathbf{n}}(\phi_{h-1} \mid m_{f})(1-p-\Delta) + \widetilde{\mu}_{h-1}^{\mathbf{n}}(1-\phi_{h-1} \mid m_{f})(1-p)\right]$$
  
$$\geq d_{h-1}^{\mathbf{n},P^{\phi,f}}(m_{f})(1-p-\Delta) \geq \cdots \geq d_{1}^{\mathbf{n},P^{\phi}}(m_{f}) \prod_{j=0}^{h-1} (1-p-\Delta)$$
  
$$\geq d_{1}^{\mathbf{n},P^{\phi}}(m_{f}) \left(1-p-\Delta\right)^{H} > \frac{\overline{\varrho}(m_{f})}{2}, \qquad (176)$$

where the last line makes use of the properties p and  $\Delta$  in (128) and

$$\left(1-p-\Delta\right)^{H} \ge \left(1-\frac{c_2}{2H}\right)^{H} \ge \left(1-\frac{1}{2H}\right)^{H} \ge \frac{1}{2},$$

provided that  $0 < c_2 < 1$ . In addition, as state  $n_f$  is an absorbing state and state  $m_f$  will only transfer to itself or state  $n_f$  at each time step, we directly achieve that

$$d_h^{\mathbf{n},P^{\phi,f}}(m_f) \le d_{h-1}^{\mathbf{n},P^{\phi,f}}(m_f) \le \dots \le d_1^{\mathbf{n},P^{\phi,f}}(m_f) \le \overline{\varrho}(m_f).$$
(177)

For state  $n_f$ , as it is absorbing, we directly have

$$d_{h}^{\mathbf{n},P^{\phi,f}}(n_{f}) = \mathbb{P}\left\{s_{h} = n_{f} \mid s_{h-1} = n_{f}; \widetilde{\mu}^{\mathbf{n}}\right\} \ge d_{h-1}^{\mathbf{n},P^{\phi,f}}(n_{f}) \ge \dots \ge d_{1}^{\mathbf{n},P^{\phi,f}}(n_{f}) = \overline{\varrho}(n_{f}).$$
(178)

According to the assumption in (131), it is easily verified that

 $d_h^{\mathbf{n},P^{\phi,f}}(n_f) \le 1 \le 2\overline{\varrho}(n_f).$ (179)

Finally, combining (176), (177, 178), (179), the definitions of  $P_h^{\star}(\cdot | s, a)$  in (125) and the Markov property, we arrive at for any  $(h, s) \in [H] \times S$ ,

$$\frac{\overline{\varrho}(s)}{2} \le d_h^{\mathsf{n},P^{\phi,f}}(s) \le 2\overline{\varrho}(s),\tag{180}$$

2356 which directly leads to

$$\frac{\overline{\varrho}(s)}{4} \le d_h^{\mathfrak{n}, P^{\phi, f}}(s, a) = \widetilde{\mu}_1^{\mathfrak{n}}(a \mid s) d_h^{\mathfrak{n}, P^{\phi, f}}(s) \le \overline{\varrho}(s).$$
(181)

**Proof of the claim (134).** Examining the definition of  $C_r^*$  in (22), we make the following observations.

$$\begin{array}{l} \textbf{2362} \qquad \textbf{. For } h = 1, \text{ we have} \\ \textbf{2363} \\ \textbf{2364} \\ \textbf{2365} \\ \textbf{2365} \\ \textbf{2365} \\ \textbf{2366} \\ \textbf{2366} \\ \textbf{2366} \\ \textbf{2367} \\ \textbf{2368} \\ \textbf{2368} \\ \textbf{2368} \\ \textbf{2369} \\ \textbf{2369} \\ \textbf{2370} \\ \textbf{2370} \\ \textbf{2371} \\ \textbf{2371} \\ \textbf{2372} \\ \textbf{2373} \\ \textbf{2373} \\ \textbf{1} \\ \textbf{1} \\ \textbf{1} \\ \textbf{2} \\ \textbf{1} \\ \textbf{2} \\ \textbf{2} \\ \textbf{2} \\ \textbf{2} \\ \textbf{2} \\ \textbf{3} \\ \textbf{3} \\ \textbf{3} \\ \textbf{4} \\ \textbf{3} \\ \textbf{4} \\ \textbf{1} \\ \textbf{1} \\ \textbf{4} \\ \textbf{5} \\ \textbf{4} \\ \textbf{1} \\ \textbf{1} \\ \textbf{5} \\ \textbf{6} \\ \textbf{6} \\ \textbf{6} \\ \textbf{7} \\ \textbf{6} \\ \textbf{6} \\ \textbf{6} \\ \textbf{7} \\ \textbf{7} \\ \textbf{6} \\ \textbf{6} \\ \textbf{6} \\ \textbf{7} \\ \textbf{7} \\ \textbf{6} \\ \textbf{6} \\ \textbf{6} \\ \textbf{7} \\ \textbf{7} \\ \textbf{6} \\ \textbf{6} \\ \textbf{7} \\ \textbf{7} \\ \textbf{6} \\ \textbf{6} \\ \textbf{7} \\ \textbf{7} \\ \textbf{6} \\ \textbf{7} \\ \textbf{7} \\ \textbf{6} \\ \textbf{7} \\ \textbf{7} \\ \textbf{7} \\ \textbf{6} \\ \textbf{7} \\ \textbf{$$

where (i) holds by  $d_1^{\star,P}(s) = \rho(s) = 0$  for all  $s \in \mathcal{N}$  (see (133)) and  $\tilde{\mu}_h^{\star,\phi}(\phi_h \mid s) = 1$  for all  $(s,h) \in \mathcal{M} \times [H]$  (see (122)), (ii) follows from the fact  $d_1^{\star,P}(s,\phi_1) = 1$  for all  $s \in \mathcal{M}$ , (iii) is verified in (132), and the last equality arises from the definition in (130). • Similarly, for  $h = 2, 3, \dots, H$ , we arrive at  $\max_{(s,a,P)\in\mathcal{S}_{\text{one}}\times\mathcal{A}_{\text{one}}\times\mathcal{U}^{\sigma}(P^{\phi})}\frac{\min\left\{d_{h}^{\star,P}(s,a),\frac{1}{4SA}\right\}}{d_{h}^{\mathsf{n},P^{\phi,f}}(s,a)}$  $\stackrel{\text{(i)}}{=} \max_{\substack{(s,P)\in\mathcal{S}\times\mathcal{U}^{\sigma}(P^{\phi})\\(s,P)\in\mathcal{M}\times\mathcal{U}^{\sigma}(P^{\phi})}} \frac{\min\left\{d_{h}^{\star,P}(s,\phi_{h}),\frac{1}{4SA}\right\}}{d_{h}^{\mathsf{n},P^{\phi,f}}(s,\phi_{h})}$   $\leq \max_{\substack{(s,P)\in\mathcal{M}\times\mathcal{U}^{\sigma}(P^{\phi})}} \frac{1}{4SAd_{h}^{\mathsf{n},P^{\phi,f}}(s,\phi_{h})}$  $\overset{(\mathrm{ii})}{\leq} \max_{s \in \mathcal{M}} \frac{1}{2SA\overline{\varrho}(s)} = 2C,$ (183)

where (i) holds by the optimal policy in (122) and the trivial fact that  $d_h^{\star,P}(s) = 0$  for all  $s \in \mathcal{N}$  (see (133) and (125)), (ii) arises from (132), and the last equality comes from (130).

Combining the above cases, we complete the proof by

$$\frac{C}{2} \leq C_{\mathbf{r}}^{\star} = \max_{(h,s,a,P) \in [H] \times \mathcal{S}_{\mathsf{one}} \times \mathcal{A}_{\mathsf{one}} \times \mathcal{U}^{\sigma}(P^{\phi})} \frac{\min\left\{d_{h}^{\star,P}(s,a), \frac{1}{4SA}\right\}}{d_{h}^{n,P^{\phi,f}}(s,a)} \leq C.$$

E.3 PROOF OF CLAIM (139) 

Recalling (121a) and (123), we first consider a more general form 

$$\begin{aligned} V_{h}^{\star,\sigma^{+},f,\phi}(m_{f}) &= V_{h}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f}) \\ = p^{\inf}V_{h+1}^{\star,\sigma^{+},f,\phi}(n_{f}) + (1-p^{\inf})V_{h+1}^{\star,\sigma^{+},f,\phi}(m_{f}) \\ &- \left(m_{h}^{\tilde{\mu},f,\phi}V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(n_{f}) + \left[1-m_{h}^{\tilde{\mu},f,\phi}\right]V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f})\right) \\ = \left(p^{\inf}-m_{h}^{\tilde{\mu},f,\phi}\right)V_{h+1}^{\star,\sigma^{+},f,\phi}(n_{f}) + m_{h}^{\tilde{\mu},f,\phi}\left(V_{h+1}^{\star,\sigma^{+},f,\phi}(n_{f}) - V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f})\right) \\ &+ (1-p^{\inf})\left(V_{h+1}^{\star,\sigma^{+},f,\phi}(m_{f}) - V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f})\right) - \left(p^{\inf}-m_{h}^{\tilde{\mu},f,\phi}\right)V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f})\right) \\ &= m_{h}^{\tilde{\mu},f,\phi}\left(V_{h+1}^{\star,\sigma^{+},f,\phi}(n_{f}) - V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f})\right) + (1-p^{\inf})\left(V_{h+1}^{\star,\sigma^{+},f,\phi}(m_{f}) - V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f})\right) \\ &+ \left(p^{\inf}-m_{h}^{\tilde{\mu},f,\phi}\right)\left(V_{h+1}^{\star,\sigma^{+},f,\phi}(m_{f}) - V_{h+1}^{\star,\sigma^{+},f,\phi}(m_{f})\right) \\ &\geq (1-p^{\inf})\left(V_{h+1}^{\star,\sigma^{+},f,\phi}(m_{f}) - V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f})\right) \\ &+ \left(p^{\inf}-m_{h}^{\tilde{\mu},f,\phi}\right)\left(V_{h+1}^{\star,\sigma^{+},f,\phi}(m_{f}) - V_{h+1}^{\star,\sigma^{+},f,\phi}(m_{f})\right) \\ &+ \frac{1}{2}(p-q)\|\tilde{\mu}_{h}^{\star,f,\phi}(\cdot|m_{f}) - \tilde{\mu}_{h}(\cdot|m_{f})\|_{1}\left(V_{h+1}^{\star,\sigma^{+},f,\phi}(n_{f}) - V_{h+1}^{\star,\sigma^{+},f,\phi}(m_{f})\right), \tag{184} \end{aligned}$$

where the last inequality holds since

2422  
2423 
$$p^{\inf} - m_h^{\tilde{\mu}, f, \phi} = (p^{\inf} - q^{\inf}) (1 - \tilde{\mu}_h(\phi_h \mid m_f))$$
2424 
$$= (p - q) (1 - \tilde{\mu}_h(\phi_h \mid m_f))$$
2425  
2426 
$$= \frac{1}{2} (p - q) (1 - \tilde{\mu}_h(\phi_h \mid m_f) + \tilde{\mu}_h(1 - \phi_h \mid m_f))$$
2427  
2428 
$$= \frac{1}{2} (p - q) \|\tilde{\mu}_h^{\star, f, \phi}(\cdot \mid m_f) - \tilde{\mu}_h(\cdot \mid m_f)\|_1, \quad (185)$$
2429

with the first equality holding by (120) and the second existing by (118).

To further control (184),  $V_h^{\star,\sigma^+,f,\phi}(n_f) - V_h^{\star,\sigma^+,f,\phi}(m_f)$  $\stackrel{\text{(i)}}{=} 1 + (1 - \sigma^+) V_{h+1}^{\star,\sigma^+,f,\phi}(n_f) + \sigma^+ V_{h+1}^{\star,\sigma^+,f,\phi}(m_f)$  $- \left( p^{\inf} V_{h+1}^{\star,\sigma^+,f,\phi}(n_f) + (1 - p^{\inf}) V_{h+1}^{\star,\sigma^+,f,\phi}(m_f) \right)$  $=1 + (1 - p^{\inf} - \sigma^+) \left( V_{h+1}^{\star, \sigma^+, f, \phi}(n_f) - V_{h+1}^{\star, \sigma^+, f, \phi}(m_f) \right)$  $\stackrel{\text{(ii)}}{=} 1 + (1-p) \left( V_{h+1}^{\star,\sigma^+,f,\phi}(n_f) - V_{h+1}^{\star,\sigma^+,f,\phi}(m_f) \right)$  $=\cdots = \sum_{j=0}^{H-h} (1-p)^j,$ (186)

where (i) follows from Lemma 8 and (ii) holds by (118). Then, we consider two cases w.r.t. the uncertainty level  $\sigma^+$  to control (186), respectively:

• When 
$$0 < \sigma^{+} \leq \frac{c_{2}}{2H}$$
: Recall  $p = \frac{c_{2}}{H}$  if  $\sigma^{+} \leq \frac{c_{2}}{2H}$ . In this case, applying (186), we have  
 $V_{h}^{\star,\sigma^{+},f,\phi}(n_{f}) - V_{h}^{\star,\sigma^{+},f,\phi}(m_{f})$   
 $= \sum_{j=0}^{H-h} (1-p)^{j} \geq \sum_{j=0}^{H-h} \left(1 - \frac{c_{2}}{H}\right)^{j} = \frac{1 - \left(1 - \frac{c_{2}}{H}\right)^{H-h+1}}{c_{2}/H} \geq \frac{2c_{2}(H-h+1)}{3}.$  (187)

Here, the final inequality holds by observing

$$\left(1 - \frac{c_2}{H}\right)^{H-h+1} \le \exp\left(-\frac{c_2(H-h+1)}{H}\right) \le 1 - \frac{2c_2(H-h+1)}{3H}, \quad (188)$$

where the first inequality holds by noticing  $c_2 < \frac{1}{2}$  and then  $1 - x \le \exp(-x)$ , and the last inequality holds by  $\exp(-x) \le 1 - \frac{2x}{3}$  for any  $0 \le x \le \frac{1}{2}$ . Plugging above fact in (187) back to (184), we arrive at

Plugging above fact in (187) back to (184), we arrive at  

$$V^{\star,\sigma^+}, f, \phi(m_{\sigma^+}) = V^{\mu}, \sigma^+, f, \phi(m_{\sigma^+})$$

$$\begin{aligned} &V_{h}^{\star,f}(m_{f}) - V_{h}^{\star,\sigma}(m_{f}) - V_{h}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f}) \\ &\geq (1 - p^{\inf}) \left( V_{h+1}^{\star,\sigma^{+},f,\phi}(m_{f}) - V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f}) \right) \\ &+ \frac{1}{2} (p - q) \left\| \widetilde{\mu}_{h}^{\star,f,\phi}(\cdot \mid m_{f}) - \widetilde{\mu}_{h}(\cdot \mid m_{f}) \right\|_{1} \frac{2c_{2}(H - h + 1)}{3}. \end{aligned}$$
(189)

Then invoking the assumption

$$\sum_{h=1}^{H} \left\| \widetilde{\mu}_{h}(\cdot \mid m_{f}) - \widetilde{\mu}_{h}^{\star, f, \phi}(\cdot \mid m_{f}) \right\|_{1} \ge \frac{H}{8}$$
(190)

in (138) and applying (189) recursively for  $h = 1, 2, \dots, H$  yields  $V^{\star,\sigma^+,f,\phi}(\dots, V^{\tilde{\mu},\sigma^+,f,\phi}(\dots, V))$ 

$$V_{1} = V_{1} = V_{1$$

where (i) follows from  $1 - p^{\inf} \ge 1 - p = 1 - \frac{c_2}{H}$ , and (ii) holds by

$$\forall h \in [H]: (1 - \frac{c_2}{H})^{h-1} \ge (1 - \frac{c_2}{H})^H \ge \frac{1}{2}b$$
 (192)

as long as  $c_2 \leq \frac{1}{2}$ . Here, (iii) arises from the definition of p, q in (126); (iv) can be verified by the fact that for any series  $0 \leq m_1, m_2, \dots, m_H \leq m_{\max}$  that obeys  $\sum_{h=1}^H m_h \geq y$ , one has

$$\sum_{h=1}^{H} m_h h \ge \sum_{h=1}^{\lfloor m_{\max}/n \rfloor} m_{\max} h, \tag{193}$$

and taking  $m_h = \|\widetilde{\mu}_{H-h+1}(\cdot | m_f) - \widetilde{\mu}_{H-h+1}^{\star,f,\phi}(\cdot | m_f)\|_1 \le 2 = m_{\max}$  and  $n = \frac{H}{8}$ . Consequently, observed from (191), the following inequality holds

$$V_1^{\star,\sigma^+,f,\phi}(m_f) - V_1^{\widetilde{\mu},\sigma^+,f,\phi}(m_f) \ge \frac{c_2\Delta}{6} \lfloor H/16 \rfloor (\lfloor H/16 \rfloor + 1) \ge c_3\Delta H^2 > \varepsilon \quad (194)$$

for some small enough constant  $c_3$  and letting  $\Delta = \frac{\varepsilon}{c_3 H^2}$ .

• When  $\frac{c_2}{2H} < \sigma^+ \le 1 - c_0$ : Similarly, recalling  $p = (1 + \frac{c_1}{H}) \sigma^+$  if  $\sigma^+ > \frac{c_2}{2H}$  and invoking (186) gives

$$V_{h}^{\star,\sigma^{+},f,\phi}(n_{f}) - V_{h}^{\star,\sigma^{+},f,\phi}(m_{f}) = \sum_{j=0}^{H-h} (1-p)^{j} = \sum_{j=0}^{H-h} \left(1 - \left(1 + \frac{c_{1}}{H}\right)\sigma^{+}\right)^{j}$$
$$\geq \frac{1 - \left(1 - \left(1 + \frac{c_{1}}{H}\right)\sigma^{+}\right)^{H-h+1}}{(1 + \frac{c_{1}}{H})\sigma^{+}}$$
$$\geq \frac{c_{2}(H-h+1)}{3\sigma^{+}H},$$
(195)

where the final inequality holds by observing

$$\left(1 - \left(1 + \frac{c_1}{H}\right)\sigma^+\right)^{H-h+1} \le \exp\left(-\left(1 + \frac{c_1}{H}\right)\sigma^+(H-h+1)\right)$$
$$\stackrel{(i)}{\le} \exp\left(-\frac{c_2}{2H}\left(1 + \frac{c_1}{H}\right)(H-h+1)\right)$$
$$\le 1 - \left(1 + \frac{c_1}{H}\right)\frac{c_2(H-h+1)}{3H}.$$
(196)

Here, (i) holds by observing  $\frac{c_2}{2H} < \sigma^+$ , and the last inequality holds by  $\left(1 + \frac{c_1}{H}\right) \leq 2$ ,  $c_2 \leq \frac{1}{2}$ , and the fact  $\exp(-x) \leq 1 - \frac{2x}{3}$  for any  $0 \leq x \leq \frac{1}{2}$ . Plugging above fact in (195) back to (184) gives

$$V_{h}^{\star,\sigma^{+},f,\phi}(m_{f}) - V_{h}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f})$$

$$\geq (1 - p^{\inf}) \left( V_{h+1}^{\star,\sigma^{+},f,\phi}(m_{f}) - V_{h+1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f}) \right)$$

$$+ \frac{1}{2} (p - q) \| \widetilde{\mu}_{h}^{\star,f,\phi}(\cdot \mid m_{f}) - \widetilde{\mu}_{h}(\cdot \mid m_{f}) \|_{1} \frac{c_{2}(H - h + 1)}{3\sigma^{+}H}.$$
(197)

Following the same routine to achieve (191), applying (197) recursively for  $h = 1, 2, \dots, H$  gives

$$V_{1}^{\star,\sigma^{+},f,\phi}(m_{f}) - V_{1}^{\tilde{\mu},\sigma^{+},f,\phi}(m_{f})$$

$$\geq \sum_{l=1}^{H} (1 - p^{\inf})^{h-1} (p - q) \frac{c_{2}(H - h + 1)}{6\sigma^{+}H} \left\| \tilde{\mu}_{h}^{\star,f,\phi}(\cdot \mid m_{f}) - \tilde{\mu}_{h}(\cdot \mid m_{f}) \right\|_{1}$$

2532  
2533  
2534  
2535  

$$\underbrace{(i)}_{h=1}^{2} \frac{c_2(p-q)}{6\sigma^+ H} \sum_{h=1}^{H} (1 - \frac{c_1}{H})^{h-1} (H-h+1) \| \widetilde{\mu}_h^{\star,f,\phi}(\cdot \mid m_f) - \widetilde{\mu}_h(\cdot \mid m_f) \|_1$$

2536  
2537
$$\stackrel{(ii)}{\geq} \frac{c_2 \Delta}{12\sigma^+ H} \lfloor H/16 \rfloor (\lfloor H/16 \rfloor + 1), \qquad (198)$$

where (i) follows from  $1 - p^{\inf} = 1 - (p - \sigma^+) = 1 - \frac{c_1}{H}\sigma^+$ , and (ii) holds by letting  $c_1 \leq \frac{1}{2}$  and following the same routine of (191). Consequently, (198) yields

$$V_1^{\star,\sigma^+,f,\phi}(m_f) - V_1^{\widetilde{\mu},\sigma^+,f,\phi}(m_f) \ge \frac{c_2\Delta}{12\sigma^+H} \lfloor H/16 \rfloor \left( \lfloor H/16 \rfloor + 1 \right) \ge \frac{c_4\Delta H}{\sigma^+} > \varepsilon,$$
(199)

which holds for some small enough constant  $c_4$  and letting  $\Delta = \frac{\sigma^+ \varepsilon}{c_4 H}$ .

#### 2548 F MULTIPLAYER GENERAL-SUM MARKOV GAMES 2549

In this section, we extend RTZ-VI-LCB to the setting of multi-player general-sum Markov games and present the corresponding theoretical guarantees.

# 2553 F.1 PROBLEM FORMULATION

2555 A robust general-sum Markov game is a tuple  $\mathcal{M}(\mathcal{S}, \{\mathcal{A}_i\}_{i=1}^m, H, \{\mathcal{U}_{\rho}^{\sigma_i}(P^0)\}_{i=1}^m, \{r_i\}_{i=1}^m)$  with m players, where S denotes the state space and H is the horizon length. We have m different action 2556 spaces, where  $\mathcal{A}_i$  is the action space for the *i*<sup>th</sup> player and  $|\mathcal{A}_i| = A_i$ . We let  $\mathcal{A} = \mathcal{A}_1 \times \cdots \times \mathcal{A}_m$ 2557 denote the joint action space, and let  $a := (a_1, \dots, a_m) \in \mathcal{A}$  denote the (tuple of) joint actions 2558 by all m players. A notable deviation from standard MGs is that: for  $1 \le i \le m$ , instead of 2559 assuming a fixed transition kernel, each  $i^{th}$  player anticipates that the transition kernel is allowed to 2560 be chosen arbitrarily from a prescribed uncertainty set  $\mathcal{U}_{\rho^{i}}^{\sigma_{i}}(P^{0})$ . Here, the uncertainty set  $\mathcal{U}_{\rho^{i}}^{\sigma_{i}}(P^{0})$ 2561 is constructed centered on  $P^0(\cdot|s, a)$ , with its size and shape defined by a certain distance metric  $\rho$  and a radius parameter  $\sigma_i > 0$ .  $r_i = \{r_{h,i}\}_{h \in [H]}$  is a collection of reward functions for the i<sup>th</sup> 2563 player, so that  $r_{h,i}(s, a)$  gives the reward received by the *i*<sup>th</sup> player if actions a are taken at state s 2564 at step h. 2565

The policy of the *i*<sup>th</sup> player is denoted as  $\pi_i := \{\pi_{h,i} : S \to \Delta_{\mathcal{A}_i}\}_{h \in [H]}$ . We denote the product policy of all players as  $\pi := \pi_1 \times \cdots \times \pi_M$ , and denote the policy of all players except the *i*<sup>th</sup> player as  $\pi_{-i}$ . We define  $V_{h,i}^{\pi}(s)$  as the expected cumulative reward that will be received by the *i*<sup>th</sup> player if starting at state *s* at step *h* and all players follow policy  $\pi$ . For any strategy  $\pi_{-i}$ , there also exists a *robust best response* of the *i*<sup>th</sup> player, which is a policy  $\mu^*(\pi_{-i})$  satisfying  $V_{h,i}^{\mu^*(\pi_{-i}),\pi_{-i},\sigma_i}(s) = \sup_{\pi_i} V_{h,i}^{\pi_i,\pi_{-i},\sigma_i}(s)$  for any  $(s,h) \in S \times [H]$ . For convenience, we denote  $V_{h,i}^{\star,\pi_{-i},\sigma_i} := V_{h,i}^{\mu^*(\pi_{-i}),\pi_{-i},\sigma_i}$ . The *Q*-functions of the robust best response can be defined similarly.

Similar to the definition of behavior policy  $(\mu^n, \nu^n)$ , we use the short-hand notation for the occupancy distribution w.r.t. the behavior policy  $\pi^n = (\pi_i^n, \pi_{-i}^n)$  as:  $\forall (h, s, a) \in [H] \times S \times A$ ,

$$d_h^{n,P^0}(s) = d_h^{\pi^n,P^0}(s) \coloneqq \mathbb{P}(s_h = s \,|\, s_1 \sim \varrho^n, \pi^n, P^0);$$
(200a)

 $d_h^{\mathsf{n},P^0}(s,\boldsymbol{a}) = d_h^{\pi^\mathsf{n},P^0}(s,\boldsymbol{a}) \coloneqq \mathbb{P}(s_h = s \,|\, s_1 \sim \varrho^\mathsf{n}, \pi^\mathsf{n}, P^0) \, \pi^\mathsf{n}(\boldsymbol{a} \,|\, s).$ 

Similarly, for any product policy 
$$\pi = (\pi_i, \pi_{-i})$$
, there is,  $\forall (h, s, a) \in [H] \times S \times A$ 

2580 2581 2582

2583 2584

2588

2589

2577 2578

2579

2538

2539

2540

2541 2542 2543

2544 2545

2546 2547

2550

2551

2552

$$d_h^{\pi_i,\pi_{-i},P}(s) \coloneqq \mathbb{P}(s_h = s \,|\, s_1 \sim \varrho, \pi, P); \tag{201a}$$

$$d_{h}^{\pi_{i},\pi_{-i},P}(s,\boldsymbol{a}) \coloneqq \mathbb{P}(s_{h}=s \,|\, s_{1} \sim \varrho, \pi, P) \,\pi_{i,h}(a_{i} \,|\, s) \,\pi_{-i,h}(\boldsymbol{a}_{-i} \,|\, s). \tag{201b}$$

Therefore, the robust variant of standard solution concepts—robust NE for Robust multi-player general-sum MGs is introcuded as follows: A product policy  $\pi$  is considered a *robust NE* if

$$\forall (s) \in \mathcal{S}, \quad V_1^{\pi,\sigma_i}(s) = V_h^{\star,\pi_{-i},\sigma^+}(s).$$
(202)

(200b)

A robust NE signifies that given the product policy  $(\pi)$  of the opponents, no player can enhance their outcome by deviating from their current policy unilaterally when each player accounts for the worst-case scenario within their uncertainty set  $\mathcal{U}_{a}^{\sigma_{i}}(P^{0})$  for all  $i = 1, 2, \dots, m$ . Since finding exact robust equilibria can be complex and may not always be feasible, practitioners often seek approximate equilibria. In this context, a product policy  $\pi \in \Delta(\mathcal{A})$  can be termed an  $\varepsilon$ -robust NE if

$$\operatorname{Gap}(\pi) \coloneqq \max\left\{\left\{V_{i,1}^{\star,\pi_{-i},\sigma_{i}}(\varrho) - V_{i,1}^{\pi,\sigma_{i}}(\varrho)\right\}_{i=1}^{m}\right\} \le \varepsilon,$$
(203)

2597 where

2595 2596

2598

2607

$$V_1^{\star,\pi_{-i},\sigma_i}(\varrho) = \mathbb{E}_{s\sim \varrho} V_1^{\star,\pi_{-i},\sigma_i}(s), \qquad \text{and} \qquad V_1^{\star,\sigma_i}(\varrho) = \mathbb{E}_{s\sim \varrho} V_1^{\star,\sigma_i}(s).$$

The existence of robust NE has been established for general divergence functions used in the uncertainty set by Blanchet et al. (2024).

**Learning objective** With a dataset collected from the nominal environment, our objective is to find a solution among the  $\varepsilon$ -robust NEs for the robust multi-player general-sum MG  $\mathcal{MG}_r$  with respect to a specified uncertainty set  $\mathcal{U}(P^0)$  around the nominal kernel, while minimizing the number of samples required under partial coverage of the state-action space.

# 2608 F.2 MULTI-RTZ-VI-LCB 2609

Here we present the Multi-RTZ-VI-LCB algorithm in Algorithm 4, which is an extension of Algorithm 2 for multi-player general-sum Markov games.

According to the empirical frequencies of state transitions, we can naturally construct an empirical estimate  $\hat{P}^0 = {\{\hat{P}_h^0\}}_{h=1}^H$  of  $P^0$ , where

$$\widehat{P}_{h}^{0}(s' \mid s, \boldsymbol{a}) = \begin{cases} \frac{1}{N_{h}(s, \boldsymbol{a})} \sum_{j=1}^{N} \mathbb{1}\left\{ \left(s_{j}, \boldsymbol{a}_{j}, s_{j}'\right) = (s, \boldsymbol{a}, s') \right\}, & \text{if } N_{h}(s, \boldsymbol{a}) > 0; \\ \frac{1}{S}, & \text{if } N_{h}(s, \boldsymbol{a}) = 0, \end{cases}$$
(204)

2617 2618 2619

2620

2623 2624

2615 2616

$$\widehat{r}_{i,h}\left(s,\boldsymbol{a}\right) = \begin{cases} r_{i,h}\left(s,\boldsymbol{a}\right), & \text{if } N_{h}\left(s,\boldsymbol{a}\right) > 0; \\ 0, & \text{if } N_{h}\left(s,\boldsymbol{a}\right) = 0, \end{cases}$$
(205)

for any  $(i, h, s, a, s') \in [m] \times [H] \times S \times A \times B \times S$ . Besides,  $N_h(s, a)$  represents the total number of sample transitions from (s, a) at step h, and

$$N_h(s, \boldsymbol{a}) \coloneqq \sum_{j=1}^N \mathbb{1}\left\{ (s_j, \boldsymbol{a}_j) = (s, \boldsymbol{a}) \right\}.$$
(206)

2625 2626

Before the details of Multi-RTZ-VI-LCB, we extend Algorithm 1 as Algorithm 3, which reduces statistical dependencies and produces a distributionally equivalent dataset  $\mathcal{D}_0$  with independent samples. Similar to Lemma 1, we present the following lemma concerning the dataset  $\mathcal{D}_0$ , whose proof is similar to the context in Appendix C.1.

-	-	-
- 3		- 73
1		. 7
_	~	~

2632	Algorithm 3: Two-stage subsampling technique for Multi-RTZ-VI-LCB.	
------	--	--

1 **Input:** Dataset  $\mathcal{D}$ , probability  $\delta$ .

2634 2 Step 1: Data Partitioning. Split  $\mathcal{D}$  into two equal-sized subsets,  $\mathcal{D}^{\mathsf{m}}$  and  $\mathcal{D}^{\mathsf{a}}$ , each containing K/2 trajectories.

**Step 2: Defining Transition Bounds.** For step h and state s, denote the number of transitions from  $\mathcal{D}^{\mathsf{m}}$  (resp.  $\mathcal{D}^{\mathsf{a}}$ ) as  $N_h^{\mathsf{m}}(s)$  (resp.  $N_h^{\mathsf{a}}(s)$ ). Construct the trimmed count as:

$$N_h^{\mathsf{t}}(s) \coloneqq \max\left\{N_h^{\mathsf{a}}(s) - 10\sqrt{N_h^{\mathsf{a}}(s)\log\frac{HS}{\delta}}, 0\right\}.$$

2640 2641

2638

4 Step 3: Generating Subsampled Dataset. Randomly sample transitions (quadruples of the  
form 
$$(s, a, h, s')$$
) from  $\mathcal{D}^m$  uniformly. For each  $(s, h) \in \mathcal{S} \times [H]$ , include  
 $\min\{N_h^t(s), N_h^m(s)\}$  transitions in the new dataset  $\mathcal{D}^t$ .

<sup>2645</sup> <sup>5</sup> Output: Set  $\mathcal{D}_0 = \mathcal{D}^t$ .

**Lemma 9** The dataset produced by the two-stage subsampling method is distributionally identical to  $\mathcal{D}_0$  with probability at least  $1 - 8\delta$ , where  $\{N_h(s, \boldsymbol{a})\}$  are independent of the sample transitions in  $\mathcal{D}^{0}$  and obey:  $\forall (h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$ , 

$$N_h(s, \boldsymbol{a}) \ge \frac{K d_h^{\mathsf{n}}(s, \boldsymbol{a})}{8} - 5\sqrt{K d_h^{\mathsf{n}}(s, \boldsymbol{a}) \log \frac{KH}{\delta}}.$$
(207)

Algorithm 4: Multi-RTZ-VI-LCB.

1 Initialization: Set uncertainty levels  $\sigma_i$  for  $i = 1, 2, \dots, m$ ; set  $\widehat{V}_{i h}^{\sigma_i}(s) = H$  and  $\widehat{Q}_{ih}^{\sigma_i}(s, \boldsymbol{a}) = H \text{ for all } (i, s, \boldsymbol{a}, h) \in [m] \times \mathcal{S} \times \mathcal{A} \times [H+1].$ <sup>2</sup> Compute the empirical reward function  $\hat{r}$  using (13) and the empirical transition kernel  $\hat{P}_0$ using (12). **3** for  $h = H, H - 1, \dots, 1$  do for *player* i = 1, 2, ..., m do Update the robust Q-value estimate as  $\widehat{Q}_{i,h}^{\sigma_{i}}\left(s,\boldsymbol{a}\right) = \min\left\{\widehat{r}_{i,h}\left(s,\boldsymbol{a}\right) + \inf_{P \in \mathcal{U}^{\sigma_{i}}\left(\widehat{P}_{h,s,\boldsymbol{a}}^{0}\right)} P\widehat{V}_{i,h+1}^{\sigma_{i}} + \beta_{i,h}\left(s,\boldsymbol{a},\widehat{V}_{i,h+1}^{\sigma_{i}}\right), H\right\},$ with  $\beta_{i,h}\left(s, \boldsymbol{a}, V\right) = \min\left\{\max\left\{\sqrt{\frac{C_{n}\log\frac{KH}{\delta}}{N_{h}(s, \boldsymbol{a})}}\mathsf{Var}_{\hat{P}_{h, s, \boldsymbol{a}}^{0}}(V), \frac{2C_{n}H\log\frac{KH}{\delta}}{N_{h}(s, \boldsymbol{a})}\right\}, H\right\}.$ **Compute** Nash policy for each  $s \in S$  as  $\pi_{h}\left(s\right) = \left(\pi_{i,h}\left(s\right), \pi_{-i,h}\left(s\right)\right) = \mathsf{ComputNash}\left(\widehat{Q}_{i,h}^{\sigma_{i}}\left(s,\cdot\right)\right),$ **Update** the robust value estimate for each  $s \in S$  as  $\widehat{V}_{i,h}^{\sigma_{i}}(s) = \mathbb{E}_{\boldsymbol{a} \sim \pi_{h}(s)} \left[ \widehat{Q}_{i,h}^{\sigma_{i}}(s, \boldsymbol{a}) \right].$ **Output**: The product policy  $\hat{\pi}(s) = {\pi_h(s)}_{h=1}^H$  with  $\pi_h(s) = \prod_{i=1}^m \pi_{i,h}(s)$ . Based on Algorithm 4, we propose a model-based approach for solving robust multi-player general-sum MGs using an approximate  $\widehat{P}^0$  for  $P^0$ , as summarized in Algorithm 4. Similar to (18), we can tackle the multi-player general-sum MGs problem as:  $\inf_{P \in \mathcal{U}^{\sigma_i}(\widehat{P}_{h,s,a}^0)} P\widehat{V}_{i,h+1}^{\sigma_i} = \max_{\alpha \in [\min_s \widehat{V}_{i,h+1}^{\sigma_i}, \max_s \widehat{V}_{i,h+1}^{\sigma_i}]} \Big\{ \widehat{P}_{h,s,a}^0 \Big[ \widehat{V}_{i,h+1}^{\sigma_i} \Big]_{\alpha} - \sigma_i \left( \alpha - \min_{s'} \Big[ \widehat{V}_{i,h+1}^{\sigma_i} \Big]_{\alpha}(s') \right) \Big\}.$ (208)where  $\left[\widehat{V}_{i,h+1}^{\sigma_i}\right]_{\alpha}$  respectively denote the clipped versions of  $\widehat{V}_{i,h+1}^{\sigma_i} \in \mathbb{R}^S$  based on some level  $\alpha \geq 0$ , as follows. 

$$\left[\widehat{V}_{i,h+1}^{\sigma_i}\right]_{\alpha}(s) := \begin{cases} \widehat{V}_{i,h+1}^{\sigma_i}(s), & \text{if } \widehat{V}_{i,h+1}^{\sigma_i}(s) > \alpha; \\ \alpha. & \text{otherwise;} \end{cases}$$
(209)

ANALYSIS OF MULTI-ME-NASH-QL F.3 

In this subsection, we prove Theorem 3, which can separated into three steps as the proof of Theorem 1. 

First of all, similar to Assumption 1, we measure the distributional discrepancy between the historical data and the target data to assess the effectiveness of the historical dataset for achieving the desired goal. We propose a novel assumption for robust multi-agent general-sum MGs as:

Assumption 2 (Robust multiple clipped concentrability) The behavior policies of the historical dataset D satisfies

2706

2719 2720 2721

2723 2724 2725

2732

2746 2747

**Step 1: decoupling statistical dependency** Before bounding  $\operatorname{Gap}(\widehat{\pi})$ , we introduce an important lemma whose proof is similar to Lemma 3 in Appendix C.3, quantifying the difference between  $\widehat{P}$ and P when projected in the direction of the value function.

 $\max\left\{\left\{\sup_{\substack{(\pi_{-i},s,\boldsymbol{a},h,P)\in\Delta(\mathcal{A}_{-i})\times\mathcal{S}\times\mathcal{A}\times[H]\times\mathcal{U}^{\sigma_{i}}(P^{0})}}\frac{\min\left\{d_{h}^{\pi_{i}^{\star},\pi_{-i},P}(s,\boldsymbol{a}),\frac{1}{S\sum_{i=1}^{m}A_{i}}\right\}}{d_{h}^{\mathfrak{n},P^{0}}(s,\boldsymbol{a})}\right\}_{i=1}^{m}\right\}\leq C_{\mathrm{mr}}^{\star}$ 

**2711 Lemma 10** Instate the assumptions in Theorem 3. Consider any vector  $V \in \mathbb{R}^S$  with  $||V||_{\infty} \leq H$  **2712** for all  $(i, h, s, a) \in [m] \times [H] \times S \times A$  satisfying  $N_h(s, a) > 0$ . With probability at least  $1 - \delta$ , **2713** one has

$$\left|\inf_{P \in \mathcal{U}^{\sigma_{i}}(\widehat{P}^{0}_{h,s,\boldsymbol{a}})} PV - \inf_{P \in \mathcal{U}^{\sigma_{i}}(P^{0}_{h,s,\boldsymbol{a}})} PV\right| \leq C_{4} \sqrt{\frac{1}{N_{h}\left(s,\boldsymbol{a}\right)}} \mathsf{Var}_{\widehat{P}^{0}_{h,s,\boldsymbol{a}}}\left(V\right) \log \frac{KH}{\delta} + C_{4} \frac{H \log \frac{KH}{\delta}}{N_{h}\left(s,\boldsymbol{a}\right)}$$
(211)

for some sufficiently large constant  $C_4 > 0$ , and

$$\operatorname{Var}_{\widehat{P}_{h,s,a}^{0}}\left(V\right) \leq 2\operatorname{Var}_{P_{h,s,a}^{0}}\left(V\right) + O\left(\frac{H^{2}}{N_{h}\left(s,a\right)}\log\frac{KH}{\delta}\right).$$

$$(212)$$

2722 With Lemma 10, we can now have

$$\inf_{P \in \mathcal{U}^{\sigma_i}(\widehat{P}^0_{h,s,\boldsymbol{a}})} PV - \inf_{P \in \mathcal{U}^{\sigma_i}(P^0_{h,s,\boldsymbol{a}})} PV \le \beta_h(s,\boldsymbol{a},V)$$
(213)

(210)

for any  $(i, h, s, a) \in [m] \times [H] \times S \times A$  satisfying  $N_h(s, a) \ge 1$ .

Therefore, we conclude that  $\widehat{Q}_{i,h}^{\sigma_i}(s, a)$  is an optimistic estimation of  $\widehat{Q}_{i,h}^{\pi,\sigma_i}(s, a)$  for any  $i = 1, 2, \dots, m$ , which is summarized below, whose proof is similar to Lemma 4 in Appendix C.4.

**2730** Lemma 11 With probability exceeding  $1 - \delta$ , it holds that

$$\widehat{Q}_{i,h}^{\sigma_i}(s, \boldsymbol{a}) \ge Q_{i,h}^{\star, \widehat{\pi}_{-i}, \sigma_i}(s, \boldsymbol{a}) \qquad and \qquad \widehat{V}_{i,h}^{\sigma_i}(s) \ge V_{i,h}^{\star, \widehat{\pi}_{-i}, \sigma_i}(s).$$
(214)

2733 Besides, we introduce another key lemma highlighting the difference between robust multi-player general-sum MGs and standard multi-player general-sum MGs from the same idea of Lemma 5, as shown below.
2736

**2737 Lemma 12** Consider any multi-player general-sum MGs  $\mathcal{MG}_r = \{S, \{\mathcal{A}_i\}_{i=1}^m, H, \{\mathcal{U}_{\rho}^{\sigma_i}(P^0)\}_{i=1}^m, \{r_i\}_{i=1}^m\}$  and the uncertainty set  $\{\mathcal{U}_{\rho}^{\sigma_i}(P^0)\}_{i=1}^m(\cdot)$  with TV **2739** distance. The optimistic robust value function estimate  $\widehat{V}_{i,h}^{\sigma_i}$ :

$$\forall (i,h) \in [m] \times [H]: \quad \max_{s \in \mathcal{S}} \widehat{V}_{i,h}^{\sigma_i} - \min_{s \in \mathcal{S}} \widehat{V}_{i,h}^{\sigma_i} \leq \min\left\{\frac{(H+1)\left(1 - (1 - \sigma_i)^{H-h}\right)}{\sigma_i}, H\right\}.$$

2744 Step 2: decomposing the error  $Gap(\hat{\pi})$  The goal of our algorithm is to output an  $\varepsilon$ -robust NE policy  $(\hat{\pi})$  satisfying  $Gap(\hat{\pi})$  in (203), i.e.,

$$\operatorname{Gap}(\widehat{\pi}) \coloneqq \max\left\{\left\{V_{i,1}^{\star,\pi_{-i},\sigma_i}(\varrho) - V_{i,1}^{\pi,\sigma_i}(\varrho)\right\}_{i=1}^m\right\} \le \varepsilon.$$

According to the relationship in Lemma 11, under the definition of  $\mathcal{A}_{-i} \coloneqq \mathcal{A}_1 \times \cdots \times \mathcal{A}_{i-1} \times \mathcal{A}_{i+1} \times \cdots \times \mathcal{A}_m$ , we obtain

$$V_{h}^{\star,\widehat{\pi}_{-i,h},\sigma^{+}}(s) \leq \widehat{V}_{i,h}^{\sigma_{i}}(s) = \max_{\pi_{i}\in\Delta(\mathcal{A}_{i})} \min_{\pi_{-i}\in\Delta(\mathcal{A}_{-i})} \mathbb{E}_{\boldsymbol{a}\sim\pi} \left[\widehat{Q}_{i,h}^{\sigma_{i}}(s,\boldsymbol{a})\right]$$

$$\leq \min_{\max_{\pi_{-i}\in\Delta(\mathcal{A}_{-i})}} \mathbb{E}_{\boldsymbol{a}\sim(\pi_{i}^{\star}(s),\pi_{-i}(s))} \left[Q_{i,h}^{\sigma_{i}}(s,\boldsymbol{a})\right], \quad (215)$$

where the first equality comes from line 8 in Algorithm 4. Therefore, there exists a deterministic policy  $\pi_{-i}^{\mathsf{d}}: \mathcal{S} \leftarrow \Delta(\mathcal{A}_{-i})$  satisfying that for any  $s \in \mathcal{S}$ 

$$\pi_{-i}^{\mathsf{d}}(s) \coloneqq \arg\min_{\pi_{-i} \in \Delta(\mathcal{A}_i)} \mathbb{E}_{\boldsymbol{a} \sim (\pi_i^{\star}(s), \pi_{-i}(s))} \left[ Q_{i,h}^{\sigma_i}(s, \boldsymbol{a}) \right].$$
(216)

Before starting, we introduce several useful notations: 

> • The state-action space covered by the behavior policy  $\pi^n$  in the nominal transition kernel  $P^0$  is denoted as

$$\mathcal{C}^{\mathsf{n}} = \{(h, s, \boldsymbol{a}) : d_{h}^{\mathsf{n}}(s, \boldsymbol{a}) > 0\}.$$
(217)

-

• The set of potential state occupancy distributions w.r.t. the policy  $(\pi_i^{\star}(s), \pi_{-i}^{\mathsf{b}}(s))$  in a model within the uncertainty set  $P \in \mathcal{U}^{\sigma_i}(P^0)$  for any  $(i,h) \in [m] \times [H]$  is denoted as

$$\mathcal{D}_{i,h}^{\mathsf{pi}} \coloneqq \left\{ \left[ d_h^{\pi_i^*(s), \pi_{-i}^{\mathsf{b}}(s), P}(s) \right]_{s \in \mathcal{S}} : P \in \mathcal{U}^{\sigma_i} \left( P^0 \right) \right\};$$
(218)

$$\mathcal{D}_{i,h}^{\mathsf{pai}} \coloneqq \left\{ \left[ d_h^{\pi_i^*(s), \pi_{-i}^{\mathsf{b}}(s), P}(s, \boldsymbol{a}) \right]_{(s, \boldsymbol{a}) \in \mathcal{S} \times \mathcal{A}} : P \in \mathcal{U}^{\sigma_i} \left( P^0 \right) \right\}.$$
(219)

• For convenience and without ambiguity, we introduce an additional notation for  $(i, h) \in$  $[m] \times [H]$  as

$$\beta_{i,h}^{\pi_i^\star,\pi_{-i}^{\mathsf{b}}}(s) = \mathbb{E}_{\boldsymbol{a}\sim(\pi_i^\star(s),\pi_{-i}^{\mathsf{b}}(s))}\beta_{i,h}\left(s,\boldsymbol{a},\widehat{V}_{i,h+1}^{\sigma_i}\right).$$

In particular, the vector  $\beta_{i,h}^{\pi_i^*,\pi_{-i}^b} \in \mathbb{R}^S$  is defined with its *s*-th item given by  $\beta_{i,h}^{\pi_i^*,\pi_{-i}^b}(s)$ .

• Similarly, we can define the notation related to rewards for  $(i, h) \in [m] \times [H]$  as

$$\widehat{r}_{i,h}^{\pi_{i}^{\star},\pi_{-i}^{\mathsf{b}}}(s) = \mathbb{E}_{\boldsymbol{a} \sim (\pi_{i}^{\star}(s),\pi_{-i}^{\mathsf{b}}(s))} \widehat{r}_{i,h}(s,\boldsymbol{a}).$$

According to the update rule in line 4 in Algorithm 4 and robust Bellman equality similar to (31), we derive

$$\leq \widehat{V}_{i,h}^{\sigma_{i}}(s) - V_{h}^{\pi_{i}^{*},\pi_{-i}^{b},\sigma^{+}}(s)$$

$$\leq \mathbb{E}_{\boldsymbol{a} \sim (\pi_{i}^{*}(s),\pi_{-i}^{b}(s))} \inf_{P \in \mathcal{U}^{\sigma_{i}}\left(\widehat{P}_{h,s,a}^{0}\right)} P\widehat{V}_{i,h+1}^{\sigma_{i}} + \beta_{i,h}^{\pi_{i}^{*},\pi_{-i}^{b}}(s)$$

$$- \mathbb{E}_{\boldsymbol{a} \sim (\pi_{i}^{*}(s),\pi_{-i}^{b}(s))} \inf_{P \in \mathcal{U}^{\sigma_{i}}\left(P_{h,s,a}^{0}\right)} PV_{i,h+1}^{\pi_{i}^{d},\pi_{-i}^{*},\sigma^{+}}$$

$$\stackrel{(i)}{\leq} \mathbb{E}_{\boldsymbol{a} \sim (\pi_{i}^{\star}(s), \pi_{-i}^{\mathsf{b}}(s))} \left[ \inf_{P \in \mathcal{U}^{\sigma_{i}}\left(P_{h, s, a}^{0}\right)} P\widehat{V}_{i, h+1}^{\sigma_{i}} - \inf_{P \in \mathcal{U}^{\sigma_{i}}\left(P_{h, s, a}^{0}\right)} PV_{i, h+1}^{\pi_{i}^{\mathsf{d}}, \pi_{-i}^{\star}, \sigma_{i}} \right] + 2\beta_{i, h}^{\pi_{i}^{\star}, \pi_{-i}^{\mathsf{b}}}(s)$$

$$\stackrel{(ii)}{\leq} \mathbb{E}_{\boldsymbol{a} \sim (\pi_{i}^{\star}(s), \pi_{-i}^{\mathsf{b}}(s))} \left[ P_{i, h, s, a}^{\inf, V}\left(\widehat{V}_{i, h+1}^{\sigma_{i}} - V_{i, h+1}^{\pi_{i}^{\mathsf{d}}, \pi_{-i}^{\star}, \sigma_{i}}\right) \right] + 2\beta_{i, h}^{\pi_{i}^{\star}, \pi_{-i}^{\mathsf{b}}}(s).$$

$$(220)$$

Here, (ii) is valid under the notation

 $V_{\sigma}^{\star,\widehat{\pi}_{-i},\sigma^+}(s) - V_{\sigma}^{\widehat{\pi},\sigma^+}(s)$ 

 $P_{i,h,s,\boldsymbol{a}}^{\mathrm{inf},V}\coloneqq \mathrm{argmin}_{P\in\mathcal{U}^{\sigma^+}\left(P_{h,s,\boldsymbol{a}}^0\right)}PV_{i,h+1}^{\pi_i^{\sharp},\pi_{-i}^{\star},\sigma^+}$ 

(221)

and consequently,  $\inf_{P \in \mathcal{U}^{\sigma_i}(P_{h,s,a}^0)} PV_{i,h+1}^{\pi_i^i, \pi_{-i}^*, \sigma^+} = P_{i,h,s,a}^{\inf V_i, \pi_{-i}^*, \sigma^+}, \text{ and } \inf_{P \in \mathcal{U}^{\sigma_i}(P_{h,s,a}^0)} P\widehat{V}_{i,h+1}^{\sigma_i} \le P_{i,h,s,a}^{\inf V_i, h+1}.$ 

Besides, (i) in (220) exists due to (213) in Lemma 10 for  $N_h(s, a) > 0$  and 

$$\left|\inf_{P\in\mathcal{U}^{\sigma_i}\left(P^0_{h,s,a}\right)}P\widehat{V}^{\sigma_i}_{i,h+1} - \inf_{P\in\mathcal{U}^{\sigma_i}\left(\widehat{P}^0_{h,s,a}\right)}P\widehat{V}^{\sigma_i}_{i,h+1}\right| \le H = \beta^{\pi^\star_i,\pi^\flat_{-i}}_{i,h}(s)$$
(222)

for  $N_h(s, \boldsymbol{a}) = 0$ . 

For ease of proof, we introduce a notation as  $\check{P}_{i,h,s}^{\inf,V} := \mathbb{E}_{\boldsymbol{a} \sim (\pi_i^*(s), \pi_{-i}^{\mathsf{b}}(s))} P_{i,h,s,\boldsymbol{a}}^{\inf,V}$ . Furthermore, we define a sequence of matrices  $\check{P}_{i,h}^{\inf,V} \in \mathbb{R}^{S \times S}$ . We can utilizing (220) recursively over the time steps  $h, h + 1, \dots, H$  and derive 

$$\begin{aligned}
\mathbf{V}_{i,h}^{\star,\widehat{\pi}_{-i},\sigma_{i}}(s) - \mathbf{V}_{i,h}^{\star,\sigma_{i}}(s) &\leq \widehat{V}_{i,h}^{\sigma_{i}}(s) - \mathbf{V}_{i,h}^{\pi_{i}^{d},\pi_{-i}^{\star},\sigma_{i}}(s) \\
&\leq \check{P}_{i,h}^{\inf,V}\left(\widehat{V}_{i,h+1}^{\sigma_{i}} - \mathbf{V}_{i,h+1}^{\pi_{i}^{d},\pi_{-i}^{\star},\sigma_{i}}\right) + 2\beta_{i,h}^{\pi_{i}^{\star},\pi_{-i}^{b}}(s) \\
&\leq \check{P}_{i,h}^{\inf,V}\check{P}_{i,h+1}^{\inf,V}\left(\widehat{V}_{i,h+2}^{\sigma_{i}} - \mathbf{V}_{i,h+2}^{\pi_{i}^{d},\pi_{-i}^{\star},\sigma^{+}}\right) + 2\check{P}_{i,h}^{\inf,V}\beta_{i,h+1}^{\pi_{i}^{d},\pi_{-i}^{\star}} + 2\beta_{i,h}^{\pi_{i}^{\star},\pi_{-i}^{b}}(s) \\
&\leq \cdots \leq 2\sum_{i'=h}^{H} \left(\prod_{j=h}^{i'-1}\check{P}_{i,j}^{\inf,V}\right)\beta_{i,i'}^{\pi_{i}^{\star},\pi_{-i}^{b}}, \\
&\qquad (223)
\end{aligned}$$

where we define  $\left(\prod_{j=h}^{i'-1} \check{P}_{i,j}^{\inf,V}\right) = I$  for convenience. 

For any  $d_h^{\pi_h^*,\pi_{-i}^b} \in \mathcal{D}_h^{\mathsf{pi}}$  (cf. (41)), taking inner product with (46) yields 

$$\left\langle d_{h}^{\pi_{i}^{\star},\pi_{-i}^{\mathsf{b}}}, V_{i,h}^{\star,\widehat{\pi}_{-i},\sigma_{i}}(s) - V_{i,h}^{\star,\sigma_{i}}(s) \right\rangle \leq \left\langle d_{h}^{\pi_{i}^{\star},\pi_{-i}^{\mathsf{b}}}, 2\sum_{i'=h}^{H} \left( \prod_{j=h}^{i'-1} \check{P}_{i,j}^{\inf,V} \right) \beta_{i,i'}^{\pi_{i}^{\star},\pi_{-i}^{\mathsf{b}}} \right\rangle$$
$$= 2\sum_{i'=h}^{H} \left\langle d_{i'}^{\mathsf{p},\pi_{i}^{\star},\pi_{-i}^{\mathsf{b}}}, \beta_{i,i'}^{\pi_{i}^{\star},\pi_{-i}^{\mathsf{b}}} \right\rangle,$$
(224)

where

$$d_{i'}^{\mathbf{p},\pi_i^{\mathsf{d}},\pi_{-i}^{\star}} \coloneqq \left[ \left( d_h^{\pi_i^{\star},\pi_{-i}^{\mathsf{b}}} \right)^{\top} \left( \prod_{j=h}^{i'-1} \check{P}_{i,j}^{\inf,V} \right) \right]^{\top} \in \mathcal{D}_{i'}^{\mathsf{pi}}$$
(225)

by the definition of  $\mathcal{D}_{i'}^{pi}$  (cf. (218)) for all  $i' = h + 1, \dots, H$ .

Next, we control  $\langle d_{i'}^{\mathbf{p},\pi_i^{\star},\pi_{-i}^{\mathbf{b}}}, \beta_{i,i'}^{\pi_i^{\star},\pi_{-i}^{\mathbf{b}}} \rangle$  utilizing concentrability. First of all, according to the definition of penalty, we demonstrate that the pessimistic penalty satisfies

$$\begin{split} \beta_{i,i'}(s, \boldsymbol{a}, \hat{V}) &\leq \max\left\{ \sqrt{\frac{C_{\mathsf{n}} \log \frac{KH}{\delta}}{N_{i}\left(s, \boldsymbol{a}\right)}} \mathsf{Var}_{\hat{P}_{i,s,a}^{0}}(\hat{V}), \frac{2C_{\mathsf{n}}H \log \frac{KH}{\delta}}{N_{i}\left(s, \boldsymbol{a}\right)} \right\} \\ &\leq \sqrt{\frac{C_{\mathsf{n}} \log \frac{KH}{\delta}}{N_{i}\left(s, \boldsymbol{a}\right)}} \mathsf{Var}_{\hat{P}_{i,s,a}^{0}}(\hat{V}) + \frac{2C_{\mathsf{n}}H \log \frac{KH}{\delta}}{N_{i}\left(s, \boldsymbol{a}\right)}}{\frac{(i)}{\delta} \sqrt{\frac{C_{\mathsf{n}} \log \frac{KH}{\delta}}{N_{i}\left(s, \boldsymbol{a}\right)}} \left(2\mathsf{Var}_{P_{i,s,a}^{0}}(\hat{V}) + \frac{C_{\mathsf{0}}H^{2}}{N_{i}\left(s, \boldsymbol{a}\right)} \log \frac{KH}{\delta}\right)} + \frac{2C_{\mathsf{n}}H \log \frac{KH}{\delta}}{N_{i}\left(s, \boldsymbol{a}\right)} \end{split}$$

$$\sum_{k=1}^{2860} \sum_{i=1}^{(ii)} \sqrt{\frac{2C_{\mathsf{n}}\log\frac{KH}{\delta}}{N_{i}\left(s,\boldsymbol{a}\right)}} \operatorname{Var}_{P_{i,s,\boldsymbol{a}}^{0}}\left(\widehat{V}\right) + \frac{\left(2C_{\mathsf{n}} + \sqrt{C_{\mathsf{n}}C_{0}}\right)H\log\frac{KH}{\delta}}{N_{i}\left(s,\boldsymbol{a}\right)}$$
(226)

where (i) holds by applying (212) for some sufficiently large  $C_0$  and (ii) exists follows from the Cauchy-Schwarz inequality. Therefore, combining the definition of  $\beta_{i,i'}^{\pi_i^*,\pi_{-i}^b}(s)$ , we obtain

$$\langle d_{i'}^{\mathbf{p},\pi_{i}^{\star},\pi_{-i}^{\mathbf{b}}},\beta_{i,i'}^{\pi_{i}^{\star},\pi_{-i}^{\mathbf{b}}}\rangle = \sum_{s\in\mathcal{S}} d_{i'}^{\mathbf{p},\pi_{i}^{\star},\pi_{-i}^{\mathbf{b}}}(s)\beta_{i,i'}^{\pi_{i}^{\star},\pi_{-i}^{\mathbf{b}}}(s) = \sum_{s\in\mathcal{S}} d_{i'}^{\mathbf{p},\pi_{i}^{\star},\pi_{-i}^{\mathbf{b}}}(s)\mathbb{E}_{\boldsymbol{a}\sim(\pi_{i}^{\star}(s),\pi_{-i}^{\mathbf{b}}(s))}\beta_{i,i'}(s,\boldsymbol{a},\hat{V}) = \sum_{(s,\boldsymbol{a})\in\mathcal{S}\times\mathcal{A}\times\mathcal{B}} d_{i'}^{\mathbf{p},\pi_{i}^{\star},\pi_{-i}^{\mathbf{b}}}(s)\mathbb{1}\{a_{i}=\pi_{i}^{\star}(s)\}\pi_{-i}^{\mathbf{d}}(\mathbf{a}_{-i}|s)\beta_{i,i'}(s,\boldsymbol{a},\hat{V}) = \sum_{(s,a_{i})\in\mathcal{S}\times\mathcal{A}} d_{i'}^{\mathbf{p},\pi_{i}^{\star},\pi_{-i}^{\mathbf{b}}}(s,a_{i},\pi_{-i}^{\mathbf{b}}(s))\beta_{i,i'}(s,\pi_{i}^{\mathbf{d}}(s),\mathbf{a}_{-i},\hat{V}),$$
(227)

where the last equation holds due to the definition in (201b). Then, we observe  $d_h^{\mathbf{p},\pi_i^\star,\pi_{-i}^\flat}(s, \mathbf{a}) \in \mathcal{D}_h^{\mathsf{pai}}$  (cf. (219)). Thereafter, we divide the bound (227) into two cases.

For the first case, i.e.,  $s \in S$  where  $\max_{P \in \mathcal{U}^{\sigma_i}(P^0)} d_{i'}^{\pi_i^{\star}, \pi_{-i}^{\mathsf{b}}, P}(s, a_i, \pi_{-i}^{\mathsf{b}}(s)) = 0$ , it follows from the definition (cf. (218)) that for any  $d_{i'}^{\mathsf{p}, \pi_i^{\star}, \pi_{-i}^{\mathsf{b}}}(s, a_i, \pi_{-i}^{\mathsf{b}}(s)) \in \mathcal{D}_i^{\mathsf{pai}}$ , it satisfies that

$$d_{i'}^{\mathbf{p},\pi_i^{\star},\pi_{-i}^{\mathsf{b}}}(s,a_i,\pi_{-i}^{\mathsf{b}}(s)) = 0.$$
(228)

For the second case, i.e.,  $s \in S$  where  $\max_{P \in \mathcal{U}^{\sigma^+}(P^0)} d_{i'}^{\pi_i^*, \pi_{-i}^{\mathsf{b}}, P}(s, a_i, \pi_{-i}^{\mathsf{b}}(s)) > 0$ , by the assumption in (210)

$$\max_{P \in \mathcal{U}^{\sigma_{i}}(P^{0})} \frac{\min\left\{d_{i'}^{\pi_{i}^{\star}, \pi_{-i}^{\mathsf{b}}, P}\left(s, a_{i}, \pi_{-i}^{\mathsf{b}}(s)\right), \frac{1}{S\sum_{i=1}A_{i}}\right\}}{d_{i'}^{\mathsf{n}}\left(s, a_{i}, \pi_{-i}^{\mathsf{b}}(s)\right)} \leq C_{\mathsf{r}}^{\star} < \infty$$

It implies that

$$d_{i'}^{\mathsf{n}}(s, a_i, \pi_{-i}^{\mathsf{b}}(s)) > 0 \quad \text{and} \quad (i', s, a_i, \pi_{-i}^{\mathsf{b}}(s)) \in \mathcal{C}^{\mathsf{n}}.$$
(229)

**2893** Lemma 9 tells that with probability at least  $1 - 8\delta$ ,

$$N_{i'}(s, a_{i}, \pi_{-i}^{\mathsf{b}}(s)) \geq \frac{Kd_{i'}^{\mathsf{n}}(s, a_{i}, \pi_{-i}^{\mathsf{b}}(s))}{8} - 5\sqrt{Kd_{i'}^{\mathsf{n}}(s, a_{i}, \pi_{-i}^{\mathsf{b}}(s)) \log \frac{KH}{\delta}}$$

$$\stackrel{(i)}{\geq} \frac{Kd_{i'}^{\mathsf{n}}(s, a_{i}, \pi_{-i}^{\mathsf{b}}(s))}{16}$$

$$\stackrel{(ii)}{\geq} \frac{K \max_{P \in \mathcal{U}^{\sigma_{i}}(P^{0})} \min \left\{ d_{i'}^{\pi_{i}^{*}, \pi_{-i}^{\mathsf{b}}, P}(s, a_{i}, \pi_{-i}^{\mathsf{b}}(s)), \frac{1}{S\sum_{i=1}A_{i}} \right\}}{16C_{\mathsf{r}}^{*}}$$

$$\geq \frac{K \min \left\{ d_{i'}^{\mathsf{p}, \pi_{i}^{*}, \pi_{-i}^{\mathsf{b}}}(s, a_{i}, \pi_{-i}^{\mathsf{b}}(s)), \frac{1}{S\sum_{i=1}A_{i}} \right\}}{16C_{\mathsf{r}}^{*}}, \qquad (230)$$

2906 where (ii) comes from Assumption 2 and (i) holds due to

$$Kd_{i'}^{\mathsf{n}}\left(s, a_{i}, \pi_{-i}^{\mathsf{b}}(s)\right) \geq c_{0} \frac{HS\sum_{i=1}A_{i}}{d_{\mathsf{m}}^{\mathsf{n}}} \log \frac{KH}{\delta} f\left(\{\sigma_{i}\}_{i=1}^{m}, H\right) d_{i'}^{\mathsf{n}}\left(s, a_{i}, \pi_{-i}^{\mathsf{b}}(s)\right)$$
$$\geq c_{0}HS\sum_{i=1}A_{i} \log \frac{KH}{\delta} f\left(\{\sigma_{i}\}_{i=1}^{m}, H\right) \geq 1600 \log \frac{KH}{\delta}, \qquad (231)$$

where  $f({\sigma_i}_{i=1}^m, H) = \min\left\{\left\{\frac{(H\sigma_i - 1 + (1 - \sigma_i)^H)}{(\sigma_i)^2}\right\}_{i=1}^m, H\right\}$ , the first inequality follows from condition (29), and the second inequality follows from

$$d_{\mathsf{m}}^{\mathsf{n}} = \min_{h,s,a_{i},\pi_{-i}^{\mathsf{b}}(s)} \left\{ d_{h}^{\mathsf{n}}(s,\pi_{i}^{\mathsf{d}}(s),\mathbf{a}_{-i}) : d_{h}^{\mathsf{n}}(s,\pi_{i}^{\mathsf{d}}(s),\mathbf{a}_{-i}) > 0 \right\} \le d_{i'}^{\mathsf{n}}(s,a_{i},\pi_{-i}^{\mathsf{b}}(s)).$$
(232)

 $=\sum_{(s,a_i)\in\mathcal{S}\times\mathcal{A}_i}d_{i'}^{\mathbf{p},\pi_i^\star,\pi_{-i}^\mathsf{b}}(s,a_i,\pi_{-i}^\mathsf{b}(s))\beta_{i,i'}(s,a_i,\pi_{-i}^\mathsf{b}(s),\hat{V})$ 

Combining the results in (49) and (50), we arrive at 

 $\langle d_{i'}^{\mathbf{p},\pi_i^\star,\pi_{-i}^\mathsf{b}}, \beta_{i,i'}^{\pi_i^\star,\pi_{-i}^\mathsf{b}} \rangle$ 

 $\leq \sum_{(s,a_i)\in\mathcal{S}\times\mathcal{A}_i} d_{i'}^{\mathbf{p},\pi_i^\star,\pi_{-i}^{\mathbf{b}}}(s,a_i,\pi_{-i}^{\mathbf{b}}(s)) \Biggl( \frac{16C_{\mathbf{r}}^\star\left(2C_{\mathbf{n}}+\sqrt{C_{\mathbf{n}}C_0}\right)H\log\frac{KH}{\delta}}{K\min\left\{d_{i'}^{\mathbf{p},\pi_i^\star,\pi_{-i}^{\mathbf{b}}}(s,a_i,\pi_{-i}^{\mathbf{b}}(s)),\frac{1}{S\sum_{i=1}A_i}\right\}}$  $+ \sqrt{\frac{32C_{\mathsf{r}}^{\star}C_{\mathsf{n}}\log\frac{KH}{\delta}}{K\min\left\{d_{i'}^{\mathsf{p},\pi_{i}^{\star},\pi_{-i}^{\mathsf{b}}}(s,a_{i},\pi_{-i}^{\mathsf{b}}(s)),\frac{1}{S\sum_{i=1}A_{i}}\right\}}}\mathsf{Var}_{P_{i,s,a_{i},\pi_{-i}^{\mathsf{b}}}^{0}}(\widehat{V})\bigg).$ (233)

Similar to the proof in Appendix B.2, we are ready to bound  $V_{i,1}^{\star,\sigma_i}(\varrho) - V_{i,1}^{\hat{\pi}_i,\star,\sigma_i}(\varrho)$ . There exists some sufficiently large constants  $C_1, C_2, C_3 > 0$ , and

 $\leq \sum_{(s,a_i)\in\mathcal{S}\times\mathcal{A}_i} d_{i'}^{\mathbf{p},\pi_i^{\star},\pi_{-i}^{\mathsf{b}}}(s,a_i,\pi_{-i}^{\mathsf{b}}(s)) \sqrt{\frac{2C_{\mathsf{n}}\log\frac{KH}{\delta}}{N_i\left(s,a_i,\pi_{-i}^{\mathsf{b}}(s)\right)}} \mathsf{Var}_{P_{i,s,a_i,\pi_{-i}^{\mathsf{b}}(s)}^0}\left(\widehat{V}\right)}$ 

 $+\sum_{(s,a_i)\in\mathcal{S}\times\mathcal{A}_i} d_{i'}^{\mathsf{p},\pi_i^\star,\pi_{-i}^\mathsf{b}}(s,a_i,\pi_{-i}^\mathsf{b}(s)) \frac{\left(2C_\mathsf{n}+\sqrt{C_\mathsf{n}C_0}\right)H\log\frac{KH}{\delta}}{N_i\left(s,a_i,\pi_{-i}^\mathsf{b}(s)\right)}$ 

$$V_{i,1}^{\star,\hat{\pi}_{-i},\sigma_{i}}(\varrho) - V_{i,1}^{\star,\sigma_{i}}(\varrho) \leq \sqrt{\frac{C_{\mathsf{r}}^{\star}C_{1}H^{3}S\sum_{i=1}A_{i}\log\frac{KH}{\delta}}{K}} \min\left\{\frac{2(H\sigma_{i}-1+(1-\sigma_{i})^{H})}{(\sigma_{i})^{2}},H\right\}}{+\frac{C_{\mathsf{r}}^{\star}C_{2}H^{2}S\sum_{i=1}A_{i}\log\frac{KH}{\delta}}{K}}{\sum_{i=1}A_{i}\log\frac{KH}{\delta}}\min\left\{\frac{2(H\sigma_{i}-1+(1-\sigma_{i})^{H})}{(\sigma_{i})^{2}},H\right\}}{\leq \sqrt{\frac{C_{\mathsf{r}}^{\star}C_{3}H^{3}S\sum_{i=1}A_{i}\log\frac{KH}{\delta}}{K}}\min\left\{\frac{2(H\sigma_{i}-1+(1-\sigma_{i})^{H})}{(\sigma_{i})^{2}},H\right\}},$$
(234)

where the last inequality follows from condition (29). 

**Step 3: summing up the results** Consequently, we obtain the upper bound of  $V_{i,1}^{\star,\hat{\pi}_{-i},\sigma_i}(\varrho)$  –  $V_{i\,1}^{\widehat{\pi},\sigma_i}(\varrho)$  in (234). which directly leads to

$$\operatorname{Gap}(\widehat{\pi}) \le c_1 \sqrt{\frac{C_{\mathsf{r}}^{\star} H^2 S \sum_{i=1}^m A_i \log \frac{KH}{\delta}}{K}} \min\left\{\left\{\frac{2(H\sigma_i - 1 + (1 - \sigma_i)^H)}{(\sigma_i)^2}\right\}_{i=1}^m, H\right\}, \quad (235)$$

for some sufficiently large  $c_1$  and

$$K \ge HS \sum_{i=1} A_i \log \frac{KH}{\delta} \min\left\{\left\{\frac{2(H\sigma_i - 1 + (1 - \sigma_i)^H)}{(\sigma_i)^2}\right\}_{i=1}^m, H\right\}.$$