

# SAMPLE-EFFICIENT DISTRIBUTIONALLY ROBUST MULTI-AGENT REINFORCEMENT LEARNING VIA ONLINE INTERACTION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Well-trained multi-agent systems can fail when deployed in real-world environments due to model mismatches between the training and deployment environments, caused by environment uncertainties including noise or adversarial attacks. Distributionally Robust Markov Games (DRMGs) enhance system resilience by optimizing for worst-case performance over a defined set of environmental uncertainties. However, current methods are limited by their dependence on simulators or large offline datasets, which are often unavailable. This paper pioneers the study of online learning in DRMGs, where agents learn directly from environmental interactions without prior data. We introduce the *Multiplayer Optimistic Robust Nash Value Iteration (MORNAVI)* algorithm and provide the first provable guarantees for this setting. Our theoretical analysis demonstrates that the algorithm achieves low regret and efficiently finds the optimal robust policy for uncertainty sets measured by Total Variation divergence and Kullback-Leibler divergence. These results establish a new, practical path toward developing truly robust multi-agent systems.

## 1 INTRODUCTION

Multi-agent reinforcement learning (MARL), along with its stochastic game-based mathematical formulation (Shapley, 1953; Littman, 1994), has emerged as a cornerstone paradigm for intelligent multi-agent systems capable of complex, coordinated behavior. It provides the theoretical and algorithmic foundation for enabling multiple agents to learn, adapt, and make sequential decisions in shared, dynamic environments. Its practical impacts span from strategic gaming, where MARL agents have achieved superhuman mastery (Silver et al., 2016; Vinyals et al., 2019); autonomous transportation, where it is used to coordinate fleets of vehicles to navigate complex traffic scenarios (Shalev-Shwartz et al., 2016; Hua et al., 2024); and distributed robotics, where teams of robots learn to execute tasks (Lowe et al., 2017; Matignon et al., 2012).

Despite the remarkable progress in MARL, a fundamental and pervasive challenge severely restricts its reliable deployment in the physical world: the *Sim-to-Real* gap (Zhao et al., 2020; Peng et al., 2018). A standard pipeline of RL involves training extensively within a high-fidelity simulator and then deploying in practice, as training directly in the real world can be prohibitively expensive, time-consuming, or dangerously unsafe. However, any simulator inevitably fails to capture the full richness and complexity of the real world, omitting subtle physical effects, unpredictable sensor noise, unmodeled system dynamics, or latent environmental factors (Padakandla et al., 2020; Rajeswaran et al., 2016). Consequently, a policy that appears optimal within the clean confines of a simulation can prove to be brittle and perform poorly—or even fail catastrophically—when deployed into the noisy, unpredictable environment it was designed for.

This vulnerability to model mismatch is magnified exponentially in the multi-agent context: this uncertainty is amplified through a cascading feedback loop of agent interactions. A minor, unmodeled perturbation that affects one agent can cause it to deviate from its expected behavior. This deviation alters the environment for its peers, who in turn must adapt their policies. Their adaptations further change the dynamics for all other agents, including the one first affected. This can trigger a chain of unpredictable responses, destabilizing the collective strategy and leading to a highly non-stationary

learning environment far more volatile than that caused by strategic adaptation alone (Papoudakis et al., 2019; Canese et al., 2021; Wong et al., 2023). The entire multi-agent system becomes fragile, as the intricate inter-agent dependencies act as amplifiers for even the smallest model inaccuracies.

To inoculate MARL agents against such environmental uncertainty, the framework of Distributionally Robust Markov Games (DRMGs) offers a principled and powerful solution (Zhang et al., 2020; Kardeş et al., 2011). Rather than trusting a single, nominal model of the environment (the simulator), the DRMG approach embraces a principle of pessimism. It defines an uncertainty set of plausible environment models centered around the nominal one. The agents’ goal is to maximize the worst-case expected returns across the entire uncertainty set. This robust optimization strategy yields two profound benefits. First, it provides a formal performance guarantee: if the true environment lies within the uncertainty set, the policy’s performance is guaranteed to be no worse than the optimized worst-case value. Second, it acts as a powerful regularizer, forcing agents to discover simpler and more generalizable policies that are inherently less sensitive to minor perturbations, thereby enhancing generalization even to environments outside the specified set (Vinitsky et al., 2020; Abdullah et al., 2019; Liu et al., 2025).

However, despite its theoretical appeal, the current body of research on DRMGs is built upon assumptions that create a critical disconnect from the realities of many high-stakes applications. The prevailing algorithmic frameworks fall into two main categories: those that assume access to a generative model (Shi et al., 2024b; Jiao & Li, 2024), which is tantamount to having a perfect, queryable oracle or simulator, and those designed for the offline setting (Li et al., 2025; Blanchet et al., 2023), which presuppose the existence of a large, static, and sufficiently comprehensive dataset collected beforehand. These assumptions are untenable in precisely the domains where robustness is most crucial. Consider applications in autonomous systems (Demontis et al., 2022) or personalized healthcare (Alaa Eldin, 2023; Lu et al., 2021). In these settings, creating a high-fidelity simulator is often impossible, and pre-collecting a dataset that covers all critical scenarios is infeasible. Agents have no choice but to learn online, through direct, sequential interaction with the complex and unknown real world. In this online paradigm, data is not a free commodity to be sampled at will; it is earned through experience, where every action has a real cost and naive exploration can lead to severe or irreversible outcomes. This necessitates a new class of algorithms that can navigate the exploration-exploitation tradeoff under the additional burden of worst-case environmental uncertainty.

We aim for robustness that survives contact with reality: agents must cope with misspecification while learning purely from experience. Without simulators or sizable offline datasets, existing approaches struggle to bridge theory and practice. This shortfall clarifies the gap we address and motivates our central question of our work: *How to design a provably effective online algorithms for distributionally robust Markov games?*

In this paper, we answer the above question by designing a model-based online algorithm for DRMGs and providing corresponding theoretical guarantees. Our contributions are summarized as follows.

**Hardness in Online DRMGs:** We first revealed the inherent hardness of online learning in DRMGs. Specifically, we showed that the online learning can suffer from the support shifting issue, where the support of the worst-case kernel is not fully covered by the support of the nominal environment, by constructing a hard instance that achieve an  $\Omega(K \min\{H, \prod_i A_i\})$ -regret for any algorithm. Moreover, we use another example to show that even without the support shifting issue, the regret can still have a minimax lower bound of  $\Omega(\sqrt{K \prod_i A_i})$ . Here,  $K$  is the number of iteration episodes,  $H$  is the DRMG horizon, and  $\prod_i A_i$  is the size of the joint action space. These results directly imply the hardness of online learning, comparing to other well-posed learning schemes including generative model (Shi et al., 2024a; Jiao & Li, 2024) or offline learning (Li et al., 2025).

**A Framework for Online Robust MARL:** We introduce  $f$ -MORNAVI, a novel model-based meta-algorithm designed specifically for online learning in DRMGs. Our framework pioneers a dual approach that synergizes the *pessimism* required for robust optimization with the *optimism* essential for provably efficient online exploration. At its core,  $f$ -MORNAVI learns the nominal environment model from online interactions and then incorporates a carefully constructed, data-driven bonus term,  $\beta$ . This bonus term is uniquely tailored to the geometry of the chosen uncertainty set, guiding exploration while guaranteeing that the learned policy is robust to worst-case model perturbations. We further present two concrete instantiations of our framework for uncertainty sets defined by Total Variation (TV) distance and Kullback-Leibler (KL) divergence.

**Near-Optimal Regret Bounds for Online DRMGs:** We establish the first known theoretical guarantees for online learning in general-sum DRMGs by providing rigorous, high-probability regret bounds for our algorithms. The regret measures the performance gap between our algorithm and an optimal robust policy, thus formally characterizing the sample complexity needed to solve the DRMG. We further prove that our algorithms converge to an  $\epsilon$ -optimal robust policy with high sample efficiency (see Corollary 1). Our results are significant as they are the first to demonstrate that finding a robust equilibrium in a general-sum DRMG is achievable in a sample-efficient manner through online interaction, without requiring a simulator or a pre-collected dataset.

## 2 PROBLEM FORMULATION

### 2.1 DISTRIBUTIONALLY ROBUST MARKOV GAMES

A *Distributionally Robust Markov Game* (DRMG) can be specified as  $\mathcal{MG}_{\text{rob}} = \{\mathcal{M}, \mathcal{S}, \mathcal{A}, H, \{\mathcal{P}_i\}_{i \in \mathcal{M}}, r\}$ , where  $\mathcal{M} = \{1, \dots, m\}$  is the set of  $m$  agents,  $\mathcal{S} = \{1, 2, \dots, S\}$  denotes the finite state space,  $\mathcal{A}$  denotes the joint action space for all agents as  $\mathcal{A} = \mathcal{A}_1 \times \dots \times \mathcal{A}_m$ , where  $\mathcal{A}_i = \{1, 2, \dots, A_i\}$  being the action space of agent  $i$ ,  $H$  denotes the horizon length. We consider non-stationary DRMGs, i.e.,  $r$  is the reward function:  $r = \{r_{i,h}\}_{1 \leq i \leq m, 1 \leq h \leq H}$  with  $r_{i,h} : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ . Specifically, for any  $(i, h, s, \mathbf{a}) \in \mathcal{M} \times [H] \times \mathcal{S} \times \mathcal{A}$ ,  $r_{i,h}(s, \mathbf{a})$  is the immediate (deterministic) reward received by the  $i$ -th agent in state  $s$  when the joint action profile is  $\mathbf{a}$ . The major difference between a DRMG and a standard Markov game is the transition kernel. Instead of having a fixed transition kernel, agents in a DRMG maintain their own uncertainty sets of transition kernels  $\mathcal{P}_i$ , to capture the potential environment uncertainties in their perspective. At each step, the environment does not transit following a fixed transition kernel, instead, it transits following an arbitrary kernel from the uncertainty set.

In this work, we mainly consider uncertainty sets specified by  $f$ -divergence (Sason & Verdú, 2016). Drawing inspiration from the rectangularity condition in robust single-agent RL (Iyengar, 2005; Wiesemann et al., 2013; Zhou et al., 2021; Shi et al., 2023), and following standard DRMG studies (Shi et al., 2024b;a; Zhang et al., 2020), we consider the *agent-wise*  $(s, \mathbf{a})$ -*rectangular uncertainty set*, due to its computational tractability. Namely, for each agent  $i$ , the DRMG specify an uncertainty set  $\mathcal{P}_i$ , which is independently defined over all horizons, states, and joint actions:

$$\mathcal{P}_i = \bigotimes_{(h,s,\mathbf{a}) \in [H] \times \mathcal{S} \times \mathcal{A}} \mathcal{P}_{i,h,f}^{\rho_i}(s, \mathbf{a}), \quad (1)$$

where  $\bigotimes$  denotes the Cartesian product. At step  $h$ , if all agents take a joint action  $\mathbf{a}_h$  at the state  $s_h$ , each agent anticipates that the transition kernel is allowed to be chosen arbitrarily from the prescribed uncertainty set  $\mathcal{P}_{i,h,f}^{\rho_i}(s_h, \mathbf{a}_h)$ . Here, the uncertainty set  $\mathcal{P}_{i,h,f}^{\rho_i}(s, \mathbf{a})$  is constructed centered on a *nominal kernel*  $P^* : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ :

**Definition 1** ( $f$ -Divergence Uncertainty Set). The  $f$ -divergence uncertainty set is defined as:

$$\mathcal{P}_{i,h,f}^{\rho_i}(s, \mathbf{a}) = \left\{ P_h \in \Delta(\mathcal{S}) : f\left(P_h, P_h^*(\cdot|s, \mathbf{a})\right) \leq \rho_i \right\},$$

where the  $f$ -divergence is defined as  $f(P_h, P_h^*(\cdot|s, \mathbf{a})) = \sum_{s' \in \mathcal{S}} f\left(\frac{P_h(s')}{P_h^*(s'|s, \mathbf{a})}\right) P_h^*(s'|s, \mathbf{a})$ .

The  $f$ -divergence uncertainty sets with different  $f$  have been extensively studied in distributionally robust RL (Clavier et al., 2023; Shi et al., 2023; Panaganti et al., 2022; Yang et al., 2022; Wang et al., 2024e; Zhang et al., 2025). In this work, we focus on TV and KL-divergence.

**Robust Value Functions.** For a DRMG, each agent aims to maximize its own worst-case performance over all possible transition kernels in its own (possibly different) prescribed uncertainty set. The strategy of agent  $i$  taking actions is captured by a policy  $\pi_i = \{\pi_{i,h} : \mathcal{S} \rightarrow \Delta(\mathcal{A}_i)\}_{h=1}^H$ . Since the immediate rewards and transition kernels are determined by the joint actions, the worst-case performance of the  $i$ -th agent over its own uncertainty set  $\mathcal{P}_i$  is determined by a joint policy  $\pi = \{\pi_h : \mathcal{S} \rightarrow \Delta(\mathcal{A})\}_{h=1}^H$ , which we refer to as the robust value function  $V_{i,h}^{\pi, \rho_i}$  and the robust

$Q$ -function  $Q_{i,h}^{\pi,\rho_i}$ , for an initial state  $s$  and initial action  $\mathbf{a}$ :

$$Q_{i,h}^{\pi,\rho_i}(s, \mathbf{a}) \triangleq \inf_{\tilde{P} \in \mathcal{P}_i} \mathbb{E}_{\pi, \tilde{P}} \left[ \sum_{t=h}^H r_{i,t}(s_t, \mathbf{a}_t) \mid s_h = s, \mathbf{a}_h = \mathbf{a} \right], V_{i,h}^{\pi,\rho_i}(s) \triangleq \sum_{\mathbf{a}} \pi(\mathbf{a}|s) Q_{i,h}^{\pi,\rho_i}(s, \mathbf{a}),$$

where the expectation is taken over the randomness of the joint policy  $\pi$  and the kernel  $\tilde{P}$ .

**Solutions to DRMGs.** Due to different objectives, the goal of a DRMG is to achieve some notions of equilibrium (Fudenberg & Tirole, 1991). We begin by formalizing the best-response policy.

For any given joint policy  $\pi$ , we use  $\pi_{-i}$  to represent the policies of all agents excluding the  $i$ -th agent. The agent  $i$ 's best response policy to  $\pi_{-i}$ ,  $\pi_i^{\dagger, \rho_i}(\pi_{-i})$ , is the policy that maximizes its own robust value function, at the give step  $h$  and state  $s$ :

$$\pi_i^{\dagger, \rho_i}(\pi_{-i}) \triangleq \arg \max_{\pi'_i \in \Delta(\mathcal{A}_i)} V_{i,h}^{(\pi_{-i} \times \pi'_i), \rho_i}(s). \quad (2)$$

The corresponding robust value function is denoted as

$$V_{i,h}^{\dagger, \pi_{-i}, \rho_i}(s) \triangleq \max_{\pi'_i \in \Delta(\mathcal{A}_i)} V_{i,h}^{\pi'_i \times \pi_{-i}, \rho_i}(s). \quad (3)$$

As noted, the objective in a DRMG is to compute an equilibrium policy (Fudenberg & Tirole, 1991): each agent's policy is a best response to the others, so no single agent can improve its robust value by deviating while the rest remain fixed. Standard notions of equilibrium include *robust Nash Equilibrium (NE)*, *robust Coarse Correlated Equilibrium (CCE)*, and *robust Correlated Equilibrium (CE)* (all of them exist (Blanchet et al., 2023)). A DRMG aims to find some approximated equilibrium:

**Robust  $\varepsilon$ -NE.** A product policy  $\pi \in \Delta(\mathcal{A}_1) \times \dots \times \Delta(\mathcal{A}_m)$  is an *robust- $\varepsilon$  NE* if for any  $s \in \mathcal{S}$ :

$$\text{gap}_{\text{NE}}(\pi, s) \triangleq \max_{i \in \mathcal{M}} \left\{ V_{i,1}^{\dagger, \pi_{-i}, \rho_i}(s) - V_{i,1}^{\pi, \rho_i}(s) \right\} \leq \varepsilon.$$

Robust NE ensures that, the agent  $i$ 's policy induced by the NE is a best response policy to the remaining agents' joint policy (up to  $\varepsilon$ ), thus no agent can improve its worst-case performance—evaluated over its own uncertainty set  $\mathcal{P}_i$ —by unilaterally deviating from the NE.

**Robust  $\varepsilon$ -CCE.** A (possibly correlated) joint policy  $\pi \in \Delta(\mathcal{A})$  is an *robust- $\varepsilon$  CCE* if for any  $s \in \mathcal{S}$ :

$$\text{gap}_{\text{CCE}}(\pi, s) \triangleq \max_{i \in \mathcal{M}} \left\{ V_{i,1}^{\dagger, \pi_{-i}, \rho_i}(s) - V_{i,1}^{\pi, \rho_i}(s) \right\} \leq \varepsilon.$$

Robust CCE relaxes the notion of NE by allowing for potentially correlated policies, while still ensuring that no agent has an incentive to unilaterally deviate from it.

**Robust  $\varepsilon$ -CE.** A joint policy  $\pi \in \Delta(\mathcal{A})$  is an *robust- $\varepsilon$  CE* if for any  $s \in \mathcal{S}$ :

$$\text{gap}_{\text{CE}}(\pi, s) \triangleq \max_{i \in \mathcal{M}} \left\{ \max_{\phi \in \Phi_i} V_{i,1}^{\phi \diamond \pi, \rho_i}(s) - V_{i,1}^{\pi, \rho_i}(s) \right\} \leq \varepsilon.$$

Here, a strategy modification  $\phi \triangleq \{\phi_{h,s}\}_{(h,s) \in [H] \times \mathcal{S}}$  for player  $i$  is a set of  $[H] \times \mathcal{S}$  functions from  $\mathcal{A}_i$  to itself. Let  $\Phi_i$  denote the set of all possible strategy modifications for player  $i$ . Given a joint policy  $\pi$ , applying a modification  $\phi$  yields a new joint policy  $\phi \diamond \pi$ , which matches  $\pi$  everywhere except that at each state  $s$  and timestep  $h$ , player  $i$ 's action  $a_i$  is replaced by  $\phi_{h,s}(a_i)$ .

**Online Learning in DRMGs.** We consider online learning in DRMGs, aiming to compute equilibria  $\{\text{NASH}, \text{CCE}, \text{CE}\}$  via interaction with the nominal environment  $P^*$  over  $K \in \mathbb{N}$  episodes. Each episode starts from  $s_1^k$ , proceeds with a policy  $\pi^k$  chosen from experience, and ends with an update for the next round. We use *robust regret* as our performance metric, which compares the learned outcome to the target equilibrium in the presence of model error.

**Definition 2 (Robust Regret).** Let  $\pi^k$  be the execution policy in the  $k^{\text{th}}$  episode. After a total of  $K$  episodes, the corresponding robust regret is defined as

$$\text{Regret}_{\{\text{NASH}, \text{CCE}, \text{CE}\}}(K) = \sum_{k=1}^K \text{gap}_{\{\text{NASH}, \text{CCE}, \text{CE}\}}(\pi^k, s_1^k).$$

Notably, if an algorithm has a sub-linear regret, it achieves a robust equilibrium as  $K \rightarrow \infty$ .

### 3 OPTIMISTIC ROBUST NASH VALUE ITERATION

We then present Multiplayer Optimistic Robust Nash Value Iteration for  $f$ -Divergence Uncertainty Set ( $f$ -MORNAVI), a meta-algorithm for episodic, finite-horizon DRMGs with interactive data collection.  $f$ -MORNAVI handles general  $f$ -divergences, with emphasis on KL and TV.

---

#### Algorithm 1: $f$ -MORNAVI

---

```

1: Input: Uncertainty level  $\rho_i > 0$  for all  $i \in \mathcal{M}$ .
2: Initialize: Dataset  $\mathbb{D} = \emptyset$ 
3: for episode  $k = 1, \dots, K$  do
4:   * Nominal Transition Estimation *
5:   Compute the transition kernel estimator  $\hat{P}_h^k(s, \mathbf{a}, s')$  as given in eq. 4.
6:   * Optimistic Robust Planning *
7:   Set  $\bar{V}_{H+1}^{k, \rho_i}(\cdot) = \underline{V}_{H+1}^{k, \rho_i}(\cdot) = 0$  for all  $i \in \mathcal{M}$ .
8:   for step  $h = H, \dots, 1$  do
9:     For all  $(s, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$  and  $i \in \mathcal{M}$ , update  $\bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a})$  [eq. 5] and  $\underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a})$  [eq. 6].
10:    For all  $s \in \mathcal{S}$ , update  $\pi_h^k(\cdot|s)$  by eq. 7.
11:    For all  $s \in \mathcal{S}$  and  $i \in \mathcal{M}$ , update  $\bar{V}_{i,h}^{k, \rho_i}(s)$  and  $\underline{V}_{i,h}^{k, \rho_i}(\cdot)$  by eq. 8.
12:  end for
13:  * Execution of policy and data collection *
14:  Receive initial State  $s_1^k \in \mathcal{S}$ 
15:  for step  $h = 1, \dots, H$  do
16:    Take action  $\mathbf{a}_h^k \sim \pi_h^k(\cdot | s_h^k)$ , observe reward  $r_h(s_h^k, \mathbf{a}_h^k)$  and next state  $s_{h+1}^k$ .
17:  end for
18:  Set  $\mathbb{D} = \mathbb{D} \cup \{(s_h^k, \mathbf{a}_h^k, s_{h+1}^k)\}_{h=1}^H$ .
19: end for
20: Output: Return policy  $\pi^{\text{out}} = \{\pi^k\}_{k=1}^K$ .

```

---

#### 3.1 ALGORITHM DESIGN

Our algorithm has the following three stages.

**Stage 1: Nominal Transition Estimation (Line 4).** At the start of each episode  $k \in [K]$ , we maintain an estimate of the nominal kernel  $P^*$  using the historical data  $\mathbb{D} = \{(s_h^\tau, \mathbf{a}_h^\tau, s_{h+1}^\tau)\}_{\tau=1, h=1}^{k-1, H}$  collected from past interactions with the training environment. Specifically,  $f$ -MORNAVI updates the empirical transition kernel for each tuple  $(h, s, \mathbf{a}, s') \in [H] \times \mathcal{S} \times \mathcal{A} \times \mathcal{S}$  as follows:

$$\hat{P}_h^k(s'|s, \mathbf{a}) = \frac{N_h^k(s, \mathbf{a}, s')}{N_h^k(s, \mathbf{a})} (\text{if } N_h^k(s, \mathbf{a}) > 0), \text{ and } \hat{P}_h^k(s'|s, \mathbf{a}) = \frac{1}{|\mathcal{S}|} (\text{if } N_h^k(s, \mathbf{a}) = 0), \quad (4)$$

where  $N_h^k(s, \mathbf{a}, s')$  and  $N_h^k(s, \mathbf{a})$ , are calculated on the current dataset  $\mathbb{D}$  by  $N_h^k(s, \mathbf{a}, s') = \sum_{\tau=1}^{k-1} \mathbf{1}\{(s_h^\tau, \mathbf{a}_h^\tau, s_{h+1}^\tau) = (s, \mathbf{a}, s')\}$ , and  $N_h^k(s, \mathbf{a}) = \sum_{s' \in \mathcal{S}} N_h^k(s, \mathbf{a}, s')$ . Note that we adopt a model-based approach that estimates transition kernels. Although this leads to higher memory consumption, model-free DRMGs are inherently challenging due to the non-linearity of worst-case expectation w.r.t. nominal kernels, which makes model-free estimators biased or sample-inefficient (Liu et al., 2022; Wang et al., 2023c; 2024d; Zhang et al., 2025).

**Stage 2: Optimistic Robust Planning (Lines 5–10).** The  $f$ -MORNAVI constructs the episode policy  $\pi^k$  via optimistic robust planning based on the empirical model  $\hat{P}^k$ . This involves estimating an upper bound on the robust value function, following the principle of Upper-Confidence-Bound (UCB) methods, which are well-established in online vanilla RL (Auer & Ortner, 2010; Azar et al., 2017; Zanette & Brunskill, 2019; Zhang et al., 2021b; Ménard et al., 2021; Zhang et al., 2024), and this optimism encourages exploration of less-visited state–action pairs.

To this end,  $f$ -MORNAVI maintains a bonus term at each episode  $k$ , capturing the gap between the robust value function under  $\hat{P}^k$  and that under the true model. This bonus is added to the robust

Bellman estimate to ensure its optimism. Specifically, for each  $(h, s, \mathbf{a}) \in [H] \times \mathcal{S} \times \mathcal{A}$ , we set

$$\bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) = \min \{r_{i,h}(s, \mathbf{a}) + \sigma_{\hat{\mathcal{P}}_{i,h,f}^{\rho_i}(s, \mathbf{a})}[\bar{V}_{i,h+1}^{k,\rho_i}] + \beta_{i,h,f}^k(s, \mathbf{a}), H\}. \quad (5)$$

$$\underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) = \max \{r_{i,h}(s, \mathbf{a}) + \sigma_{\hat{\mathcal{P}}_{i,h,f}^{\rho_i}(s, \mathbf{a})}[\underline{V}_{i,h+1}^{k,\rho_i}] - \beta_{i,h,f}^k(s, \mathbf{a}), 0\}, \quad (6)$$

here,  $\sigma_{\mathcal{P}}[V] = \inf_{P \in \mathcal{P}} \mathbb{E}_P[V]$  is the support function of  $V$  over the uncertainty set  $\mathcal{P}$ , and can be calculated through its dual representation (see Lemma 1);  $\hat{\mathcal{P}}_{i,h,f}^{\rho_i}$  is the uncertainty set centered at  $\hat{P}^k$  from eq. 4:  $\hat{\mathcal{P}}_{i,h,f}^{\rho_i}(s, \mathbf{a}) = \{P_h \in \Delta(\mathcal{S}) : f(P_h, \hat{P}_h(\cdot|s, \mathbf{a})) \leq \rho_i\}$ .

Each of these estimates in eq. 5 and eq. 6 are based on estimated robust Bellman operators (see Appendix C for details) and a bonus term  $\beta_{i,h,f}^k(s, \mathbf{a}) \geq 0$ . The bonus term is constructed (we will discuss the construction later) to ensure the estimation becomes a confidence interval of the true robust value function, i.e.,  $Q_{i,h}^{\dagger, \pi^{-i}, \rho_i}(s, \mathbf{a}) \in [\underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}), \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})]$ , with high probability.

**EQUILIBRIUM subroutine (Line 8).** Given robust  $Q$ -function estimates  $\bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})$  and  $\underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})$  for  $i \in \mathcal{M}$  at step  $h$ , the sub-routine EQUILIBRIUM  $\in \{\text{NASH}, \text{CCE}, \text{CE}\}$  finds a corresponding equilibrium  $\pi_h^k(\cdot|s)$  for the matrix-form game with pay-off matrices  $\{\bar{Q}_{i,h}^{k,\rho_i}(s, \cdot)\}_{i \in \mathcal{M}}$ :

$$\pi_h^k(\cdot|s) \leftarrow \text{EQUILIBRIUM}\left(\left\{\bar{Q}_{i,h}^{k,\rho_i}(s, \cdot)\right\}_{i \in \mathcal{M}}\right). \quad (7)$$

Note that finding a NE can be PPAD-hard (Daskalakis et al., 2009), but computing CE or CCE remains tractable in polynomial time (Liu et al., 2021).

We then update the estimation of  $V_h^{\dagger, \pi^{-i}, \rho}$  as

$$\bar{V}_{i,h}^{k,\rho_i}(s) = \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [\bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})] \quad \text{and} \quad \underline{V}_{i,h}^{k,\rho_i}(s) = \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [\underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})]. \quad (8)$$

Note that while the lower estimate in eq. 6 does not influence policy execution directly, it plays a crucial role in constructing valid exploration bonuses and ensuring strong theoretical guarantees. By leveraging both upper and lower bounds, the algorithm performs optimistic robust planning, enabling structured, uncertainty-aware exploration that balances exploration, exploitation, and robustness.

**Stage 3: Execution of Policy and Data Collection (Lines 11–17).** After evaluating the policy  $\{\pi_h^k\}_{h=1}^H$  for episode  $k$ , the learner takes action based on  $\pi_h^k$  and observes the reward  $r_h(s_h^k, \mathbf{a}_h^k)$  and next state  $s_{h+1}^k$ , which get appended to the historical dataset collected till episode  $k-1$ .

## 4 HARDNESS OF ONLINE LEARNING

In this section, we aim to discuss the inherent hardness of online learning in DRMGs from two aspects: (1) When there is the support shift issue, no MARL algorithm can obtain a sub-linear regret on a certainty DRMG; (2) Even if there is no support shift issue, there exists a DRMG such that any online algorithm suffers from the curse of multi-agency. This is a separation between DRMGs with interactive data collection and generative model/offline data, and also between DRMGs with non-robust MGs, showing the inherent challenges of online DRMGs.

### 4.1 HARDNESS WITH SUPPORT SHIFT

Support shift (Lu et al., 2024) refers to the case that the support of the worst-case transition kernel is not covered by the support of the nominal kernel. It can happen when, for instance, the uncertainty set is defined through TV. It will result in a challenge that, for those states that is not covered by the nominal kernel, there is no data available, so that the agent can never learn the optimal robust policy efficiently. Specifically, we derive the following result to illustrate the hardness.

**Theorem 1.** There exists a TV-DRMG, such that any online learning algorithm suffers the following regret lower bound:

$$\inf_{\mathcal{ALG}} \mathbb{E}[\text{Regret}_{\text{NASH}}(K)] \geq \Omega\left(\rho K \cdot \min\{H, \prod_{i \in \mathcal{M}} A_i\}\right).$$

Our construction is deferred to Example 1 in Appendix. This regret bound is linear in the number of episodes  $K$ , creating a combinatorial explosion that makes the problem information-theoretically intractable. Moreover, our result shows that when the game horizon  $H$  is large enough, the minimax lower bound depends on the joint action space, showing the hardness of online learning compared to generative models and offline settings.

## 4.2 HARDNESS WITHOUT SUPPORT SHIFT

We then illustrate the hardness of online DRMGs when there is no support shift. Note that when the uncertainty set is defined through, e.g., KL divergence, the worst-case support will be covered by the nominal one, so there will not be any support shift. However, we construct another example to show that, even without the support shift, the online learning can still be challenging and inefficient.

**Theorem 2** (Lower Bound for Robust Learning without Support Shift). There exists a DRMG, such that any learning algorithm suffers the following cumulative regret lower bound over  $K$  episodes:

$$\inf_{\mathcal{ALG}} \mathbb{E}[\text{Regret}_{\text{NASH}}(K)] \geq \Omega\left(\sqrt{K \prod_{i \in \mathcal{M}} A_i}\right).$$

Our construction is in Example 2 in Appendix. This result illustrates that, even without any support shift, some hard instance can require at least  $\Omega\left(\sqrt{K \prod_{i \in \mathcal{M}} A_i}\right)$  regret. Our result hence suggests that the dependence on the joint action space may be inevitable in online robust learning, which suffer from the curse of multi-agency.

## 5 THEORETICAL GUARANTEES

We then develop the theoretical results of our algorithm under both TV and KL sets.

### 5.1 REGRET BOUND FOR TOTAL VARIATION

As discussed in Section 4, no efficient algorithm can be expected due to the support shifting issue. We hence adopt a standard fail-state assumption (Lu et al., 2024; Liu et al., 2024) to ensure the worst-case kernel support will be covered by the nominal one, bypassing the issue.

**Assumption 1** (Failure States). For any agent  $i$ , there exists an (agent-specified) set of failure states  $\mathcal{S}_f^i \subseteq \mathcal{S}$ , such that  $r_i(s, \mathbf{a}) = 0$ , and  $P_h^*(s'|s, \mathbf{a}) = 1$ ,  $\forall \mathbf{a} \in \mathcal{A}, \forall s \in \mathcal{S}_f^i, \forall s' \in \mathcal{S}_f^i$ .

This assumption is only for TV case. Assumption 1 is a standard assumption in single-agent robust RL studies (Panaganti et al., 2022; Lu et al., 2024), and we adapt it to multi-agent cases.

We then present our theoretical guarantees.

**Theorem 1** (Upper bound of TV-MORNAVI). Denote  $\rho_{\min} := \min_{i \in \mathcal{M}} \rho_i$ . For any  $\delta \in (0, 1)$ ,

we set  $\beta_{i,h,f}^k(s, \mathbf{a})$  as  $\sqrt{\frac{c_1 \iota \text{Var}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \frac{V_{i,h+1}^{k, \rho_i} + V_{i,h+1}^{k, \rho_i}}{2} \right]}{N_h^k(s, \mathbf{a}) \vee 1}} + \frac{c_2 H^2 S \iota}{\sqrt{N_h^k(s, \mathbf{a}) \vee 1}} + \frac{2 \mathbb{E}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} [V_{i,h+1}^{k, \rho_i} - V_{i,h+1}^{k, \rho_i}]}{H} + \frac{1}{\sqrt{K}}$ , where  $\iota = \log \left( S^2 (\prod_{i=1}^m A_i) H^2 K^{3/2} / \delta \right)$  and  $c_1, c_2$  are absolute constants. Then under Assumption 1, for EQUILIBRIUM being one of  $\{\text{NASH}, \text{CE}, \text{CCE}\}$ , with probability at least  $1 - \delta$ , the regret of our TV-MORNAVI algorithm can be bounded as:

$$\text{Regret}_{\{\text{NASH}, \text{CCE}, \text{CE}\}}(K) = \tilde{O} \left( \sqrt{\min \{ \rho_{\min}^{-1}, H \} H^2 S K \left( \prod_{i \in \mathcal{M}} A_i \right)} \right),$$

where  $f(K) = \tilde{O}(g(K))$  means  $f(K) \leq \text{Poly}(\log(K)) \cdot g(K)$  for sufficiently large  $K$  and some polynomial of  $\log(K)$ .

### 5.2 REGRET BOUND FOR KL-DIVERGENCE

We then study the regret bound of KL-divergence set. As discussed, KL set is free from supporting issue hence no additional assumption is required. Our regret bound result is as follows.

**Theorem 2.** For any  $\delta$ , set  $\beta_{i,h,f}^k(s, \mathbf{a})$  in KL-DRMG as  $\frac{2c_f H}{\rho_i} \sqrt{\frac{\iota}{(N_h^k(s, \mathbf{a}) \vee 1) \hat{P}_{\min, h}^k(s, \mathbf{a})}} + \sqrt{\frac{1}{K}}$ , where  $\hat{P}_{\min, h}^k(s, \mathbf{a}) = \min_{s' \in \mathcal{S}} \{\hat{P}_h^k(s'|s, \mathbf{a}) : \hat{P}_h^k(s'|s, \mathbf{a}) > 0\}$ ,  $\iota = \log \left( S^2 \left( \prod_{i=1}^m A_i \right) H^2 K^{3/2} / \delta \right)$ , and  $c_f$  is an absolute constant. Then for EQUILIBRIUM being one of  $\{\text{NASH}, \text{CE}, \text{CCE}\}$ , with probability at least  $1 - \delta$ , it holds that

$$\text{Regret}_{\{\text{NASH}, \text{CCE}, \text{CE}\}}(K) = \tilde{O} \left( \sqrt{H^4 \exp(2H^2) K S \left( \prod_{i \in \mathcal{M}} A_i \right) \left( \rho_{\min}^2 P_{\min}^* \right)^{-1}} \right), \quad (9)$$

here,  $P_{\min}^* \triangleq \min_{(s, \mathbf{a}, s', h): P_h(s'|s, \mathbf{a}) > 0} P(s'|s, \mathbf{a})$  is the smallest positive entry of the nominal kernel.

We note that  $\exp(H)$  term in KL results can be replaced by  $P_{\min}^{-1}$  (Panaganti & Kalathil, 2022; Blanchet et al., 2023), and both of these terms are inevitable.

### 5.3 SAMPLE COMPLEXITY

As a direct corollary, we derive the sample complexity to learn an  $\varepsilon$ -equilibrium. Using a standard online-to-batch conversion (Cesa-Bianchi et al., 2001), we have the following results.

**Corollary 1** (Sample Complexity). *With probability at least  $1 - \delta$ , and under the settings of Theorem 1 and Theorem 2, the number of samples required to find an  $\epsilon$ -approximate equilibrium is bounded as:*

$$KH = \begin{cases} \tilde{O} \left( \epsilon^{-2} \min \{ \rho_{\min}^{-1}, H \} H^3 S \left( \prod_{i \in \mathcal{M}} A_i \right) \right), & \text{for TV-DRMG} \\ \tilde{O} \left( \epsilon^{-2} H^5 \exp(2H^2) S \left( \prod_{i \in \mathcal{M}} A_i \right) \left( \rho_{\min}^2 P_{\min}^* \right)^{-1} \right), & \text{for KL-DRMG} \end{cases}.$$

Our results hence implies that, despite the inherent hardness of online learning in DRMGs, our algorithm is able to learn an equilibrium with efficient sample complexity. As we shall discussed in the next section, our complexity bounds are near-optimal (expect the term  $\prod_{i \in \mathcal{M}} A_i$ ), which hence implies the efficiency of our method.

## 6 COMPARISON WITH PRIOR WORKS AND DISCUSSION

We then compare our results with prior works (the detailed Comparisons are shown in Table 1).

A substantial body of research on DRMGs has focused on two primary settings: (i) generative model setting, where the agents can freely sample from all state-action pairs (Shi et al., 2024a;b; Jiao & Li, 2024); (ii) offline setting, which relies on a comprehensive, pre-collected dataset (Blanchet et al., 2023; Li et al., 2025). As we discuss in Section 4, both of these avoid exploration and are therefore easier than the online regime we consider. Despite this added difficulty, our algorithm attains complexities comparable to those reported for the generative and offline settings.

For both uncertainty sets, our results match or improve upon previous results and the minimax lower bound in all parameters except for the action-product term,  $\prod_i A_i$ , under the generative model setting. In the offline setting, if the dataset is generated uniformly, the convergence coefficients  $C_{u/p}^*$  from (Li et al., 2025; Blanchet et al., 2023) introduce an additional  $\prod_i A_i$  term into the sample complexity. Consequently, our results also match or surpass the offline complexity in all parameter dependence. This raises an important open question:

**Can any online DRMG learning algorithm (or even under generative model settings) overcome the curse of multi-agency and eliminate the dependence on  $\prod_i A_i$ ?**

While some works (Shi et al., 2024a; Jiao & Li, 2024; Li et al., 2025; Ma et al., 2023) have achieved independence from  $\prod_i A_i$ , it remains unclear whether these improvements are applicable to general DRMGs. Specifically, the results in (Shi et al., 2024a) and (Jiao & Li, 2024) are developed for special uncertainty sets with desirable properties. For instance, the fictitious TV uncertainty set in (Shi et al., 2024a) allows the global transition kernel to be estimated from a single agent’s local information; And robust RL under contamination models is known to be equivalent to a non-robust problem with a



Table 1: Comparison with prior results.  $C_{u/p}^*$  are coverage coefficients for offline learning. In (Li et al., 2025),  $f(H, \rho) = (H\rho - 1 + (1 - \rho)^H)/\rho^2$ . The  $\exp(H)$  term in KL results can be replaced by  $P_{\min}^{-1}$  directly (Panaganti & Kalathil, 2022; Blanchet et al., 2023).

Setting & Algorithm	Uncertainty Set	Sample Complexity
<b>Generative</b> (Shi et al., 2024b)	TV	$\tilde{O}(\epsilon^{-2} H^3 S(\prod_{i \in \mathcal{M}} A_i) \min\{\rho_{\min}^{-1}, H\})$
<b>Generative</b> (Jiao & Li, 2024)	Contamination	$\tilde{O}(\epsilon^{-2} H^3 S(\sum_{i \in \mathcal{M}} A_i) \min\{\rho_{\min}^{-1}, H\})$
<b>Generative</b> (Shi et al., 2024a)	TV (fictitious)	$\tilde{O}(\epsilon^{-4} H^6 S(\sum_{i \in \mathcal{M}} A_i) \min\{\rho_{\min}^{-1}, H\})$
<b>Offline</b> (Blanchet et al., 2023)	KL	$\tilde{O}(\epsilon^{-2} \rho_{\min}^{-2} C_u^* H^4 \exp(H) S^2(\prod_{i \in \mathcal{M}} A_i))$
	TV	$\tilde{O}(\epsilon^{-2} C_u^* H^4 S^2(\prod_{i \in \mathcal{M}} A_i))$
<b>Offline</b> (Li et al., 2025)	TV	$\tilde{O}(\epsilon^{-2} C_p^* H^4 S(\sum_{i=1}^m A_i) \min\{f(H, \rho), H\})$
<b>Online</b> (Ma et al., 2023)	KL	$\tilde{O}(\epsilon^{-2} H^5 S(\max_i\{A_i\})^2)$ (with an oracle)
<b>Online</b> (Our work)	TV	$\tilde{O}(\epsilon^{-2} H^3 S(\prod_{i \in \mathcal{M}} A_i) \min\{\rho_{\min}^{-1}, H\})$
	KL	$\tilde{O}(\epsilon^{-2} \rho_{\min}^{-2} (P_{\min}^*)^{-1} H^5 \exp(2H^2) S(\prod_{i \in \mathcal{M}} A_i))$
<b>Generative</b> <i>Lower bound</i> (Shi et al., 2024b)	TV	$\Omega(\epsilon^{-2} H^3 S(\max_{i \in \mathcal{M}} A_i) \min\{\rho_{\min}^{-1}, H\})$

specific discount factor (Wang et al., 2023a). And the improvement in the offline setting is attributed to the benefits of the coverage coefficient.

The only online method (which also breaks the curse of multi-agency) is presented in (Ma et al., 2023). However, their algorithm relies on additional assumptions about uncertainty sets and a powerful oracle. This oracle is required to provide an  $\epsilon$ -accurate estimation of the worst-case performance,  $\sigma_{\mathcal{P}_i}[V]$  (see Theorem 12 of (Ma et al., 2023)), without any need for exploration. A central challenge in the analysis of robust learning algorithms is precisely quantifying this estimation error, as demonstrated in works like (Shi et al., 2023; Xu et al., 2023; Panaganti & Kalathil, 2022; Liu & Xu, 2024). By assuming the existence of such an oracle, they bypass this core challenge, which significantly reduces their sample complexity. Moreover, their results need additional assumptions on the radius  $\rho$ . For instance, it is assumed that  $\rho \leq \frac{P_{\min}^*}{H}$ , whereas ours do not require any of them.

Therefore, it is still uncertain whether the complexity reduction in these papers is a blessing of their specific uncertainty set structures, the properties of offline coverage coefficients, or the use of an estimation oracle. Furthermore, based on our discussion in Section 4, it is not clear whether the minimax lower bound for online DRMGs is independent of the size of the joint action space. We, therefore, leave the exploration of this direction, including whether practical relaxations and techniques can avoid it, for future work.

## 7 CONCLUSION

In this paper, we introduced the Multiplayer Optimistic Robust Nash Value Iteration (MORNAVI) algorithm, pioneering the study of online learning in DRMGs. Our work provides the first provable guarantees for this challenging setting, demonstrating that MORNAVI achieves low regret and efficiently identifies optimal robust policies for TV-divergence and KL-divergence uncertainty sets. This research establishes a practical path toward developing truly robust multi-agent systems that learn directly from environmental interactions. Despite the inherent hardness of online DRMGs, our algorithm achieves complexity results comparable to generative model and offline settings. This work also highlights a critical open question: whether online DRMG learning algorithms can overcome the curse of multi-agency and eliminate the dependence on the joint action space size. Future work will explore this fundamental challenge to advance the scalability of robust MARL.

## REFERENCES

- Mohammed Amin Abdullah, Hang Ren, Haitham Bou Ammar, Vladimir Milenkovic, Rui Luo, Mingtian Zhang, and Jun Wang. Wasserstein robust reinforcement learning. *arXiv preprint arXiv:1907.13196*, 2019.
- Baraa Alaa Eldin. Why applying deep reinforcement learning in healthcare is hard. <https://medium.com/@baraa.alaa.eldin/why-applying-deep-reinforcement-learning-in-healthcare-is-hard-ffc6e05ab7ca>, 2023. Accessed: 2025-07-28.
- Peter Auer and Ronald Ortner. UCB revisited: Improved regret bounds for the stochastic multi-armed bandit problem. *Periodica Mathematica Hungarica*, 61(1-2):55–65, 2010.
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax Regret Bounds for Reinforcement Learning. In *International conference on machine learning*, pp. 263–272. PMLR, 2017.
- Kishan Panaganti Badrinath and Dileep Kalathil. Robust reinforcement learning using least squares policy iteration with provable performance guarantees. In *Proc. International Conference on Machine Learning (ICML)*, 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.
- 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*, 36:66845–66859, 2023.
- Alexander Bukharin, Yan Li, Yue Yu, Qingru Zhang, Zhehui Chen, Simiao Zuo, Chao Zhang, Songan Zhang, and Tuo Zhao. Robust multi-agent reinforcement learning via adversarial regularization: Theoretical foundation and stable algorithms. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, volume 36, pp. 68121–68133, 2023.
- Lucian Busoniu, Robert Babuska, and Bart De Schutter. A comprehensive survey of multiagent reinforcement learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 38(2):156–172, 2008.
- Lorenzo Canese, Gian Carlo Cardarilli, Luca Di Nunzio, Rocco Fazzolari, Daniele Giardino, Marco Re, and Sergio Spanò. Multi-agent reinforcement learning: A review of challenges and applications. *Applied Sciences*, 11(11):4948, 2021.
- Nicoló Cesa-Bianchi, Alex Conconi, and Claudio Gentile. On the Generalization Ability of On-Line Learning Algorithms. *Advances in neural information processing systems*, 14, 2001.
- Zixiang Chen, Dongruo Zhou, and Quanquan Gu. Almost optimal algorithms for two-player zero-sum linear mixture markov games. In Sanjoy Dasgupta and Nika Haghtalab (eds.), *Proceedings of The 33rd International Conference on Algorithmic Learning Theory*, volume 167 of *Proceedings of Machine Learning Research*, pp. 227–261. PMLR, 29 Mar–01 Apr 2022. URL <https://proceedings.mlr.press/v167/chen22d.html>.
- 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.
- Constantinos Daskalakis. On the complexity of approximating a nash equilibrium. *ACM Transactions on Algorithms (TALG)*, 9(3):1–35, 2013.

- Constantinos Daskalakis, Paul W Goldberg, and Christos H Papadimitriou. The complexity of computing a nash equilibrium. *Communications of the ACM*, 52(2):89–97, 2009.
- Ambra Demontis, Maura Pintor, Luca Demetrio, Kathrin Grosse, Hsiao-Ying Lin, Chengfang Fang, Battista Biggio, and Fabio Roli. A survey on reinforcement learning security with application to autonomous driving, 2022. URL <https://arxiv.org/abs/2212.06123>.
- Xiaotie Deng, Ningyuan Li, David Mguni, Jun Wang, and Yaodong Yang. On the complexity of computing markov perfect equilibrium in general-sum stochastic games. *National Science Review*, 10(1):nwac256, 2023.
- Jing Dong, Jingwei Li, Baoxiang Wang, and Jingzhao Zhang. Online policy optimization for robust mdp. *arXiv preprint arXiv:2209.13841*, 2022.
- Devdatt P Dubhashi and Alessandro Panconesi. *Concentration of Measure for the Analysis of Randomized Algorithms*. Cambridge University Press, 2009.
- Songtao Feng, Ming Yin, Yu-Xiang Wang, Jing Yang, and Yingbin Liang. Improving sample efficiency of model-free algorithms for zero-sum markov games. *arXiv preprint arXiv:2308.08858*, 2023.
- Arlington M Fink. Equilibrium in a stochastic  $n$ -person game. *Journal of science of the hiroshima university, series ai (mathematics)*, 28(1):89–93, 1964.
- Drew Fudenberg and Jean Tirole. *Game theory*. MIT press, 1991.
- Debamita Ghosh, George K. Atia, and Yue Wang. Provably near-optimal distributionally robust reinforcement learning in online settings, 2025. URL <https://arxiv.org/abs/2508.03768>.
- Amy Greenwald, Keith Hall, Roberto Serrano, et al. Correlated q-learning. In *ICML*, volume 3, pp. 242–249, 2003.
- Songyang Han, Sanbao Su, Sihong He, Shuo Han, Haizhao Yang, Shaofeng Zou, and Fei Miao. What is the solution for state-adversarial multi-agent reinforcement learning? *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=HyqSwNhM3x>.
- Thomas Dueholm Hansen, Peter Bro Miltersen, and Uri Zwick. Strategy iteration is strongly polynomial for 2-player turn-based stochastic games with a constant discount factor. *Journal of the ACM (JACM)*, 60(1):1–16, 2013.
- Sihong He, Songyang Han, Sanbao Su, Shuo Han, Shaofeng Zou, and Fei Miao. Robust multi-agent reinforcement learning with state uncertainty, 2023. URL <https://arxiv.org/abs/2307.16212>.
- Junling Hu and Michael P Wellman. Nash q-learning for general-sum stochastic games. *Journal of machine learning research*, 4(Nov):1039–1069, 2003.
- Min Hua, Dong Chen, Xinda Qi, Kun Jiang, Zemin Eitan Liu, Quan Zhou, and Hongming Xu. Multi-agent reinforcement learning for connected and automated vehicles control: Recent advancements and future prospects, 2024. URL <https://arxiv.org/abs/2312.11084>.
- Garud N Iyengar. Robust Dynamic Programming. *Mathematics of Operations Research*, 30(2): 257–280, 2005.
- Yuchen Jiao and Gen Li. Minimax-optimal multi-agent robust reinforcement learning. *arXiv preprint arXiv:2412.19873*, 2024.
- Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is Q-learning provably efficient? In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, pp. 4868–4878, 2018.
- Chi Jin, Qinghua Liu, Yuanhao Wang, and Tiancheng Yu. V-learning—a simple, efficient, decentralized algorithm for multiagent rl. *arXiv preprint arXiv:2110.14555*, 2021.

- Yujia Jin, Vidya Muthukumar, and Aaron Sidford. The complexity of infinite-horizon general-sum stochastic games. *arXiv preprint arXiv:2204.04186*, 2022.
- Erim Kardeş, Fernando Ordóñez, and Randolph W Hall. Discounted robust stochastic games and an application to queueing control. *Operations research*, 59(2):365–382, 2011.
- Na Li, Yuchen Jiao, Hangguan Shan, and Shefeng Yan. Provable memory efficient self-play algorithm for model-free reinforcement learning. In *The Twelfth International Conference on Learning Representations*, 2024.
- Na Li, Zewu Zheng, Wei Ni, Hangguan Shan, Wenjie Zhang, and Xinyu Li. Sample efficient robust offline self-play for model-based reinforcement learning. Manuscript, OpenReview preprint, 2025. URL <https://openreview.net/forum?id=3lXZjsir0e>.
- Shihui Li, Yi Wu, Xinyue Cui, Honghua Dong, Fei Fang, and Stuart Russell. Robust multi-agent reinforcement learning via minimax deep deterministic policy gradient. In *Proc. Conference on Artificial Intelligence (AAAI)*, volume 33, pp. 4213–4220, 2019.
- Yan Li, Guanghui Lan, and Tuo Zhao. First-order policy optimization for robust markov decision process. *arXiv preprint arXiv:2209.10579*, 2022.
- Zhipeng Liang, Xiaoteng Ma, Jose Blanchet, Jiheng Zhang, and Zhengyuan Zhou. Single-trajectory distributionally robust reinforcement learning. *arXiv preprint arXiv:2301.11721*, 2023.
- Jieyu Lin, Kristina Dzevaroska, Sai Qian Zhang, Alberto Leon-Garcia, and Nicolas Papernot. On the robustness of cooperative multi-agent reinforcement learning, 2020. URL <https://arxiv.org/abs/2003.03722>.
- Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In *Machine learning proceedings 1994*, pp. 157–163. Elsevier, 1994.
- Michael L Littman and Csaba Szepesvári. A generalized reinforcement-learning model: Convergence and applications. In *ICML*, volume 96, pp. 310–318, 1996.
- Michael L Littman et al. Friend-or-foe q-learning in general-sum games. In *ICML*, volume 1, pp. 322–328, 2001.
- Guangyi Liu, Suzan Iloglu, Michael Caldara, Joseph W Durham, and Michael M. Zavlanos. Distributionally robust multi-agent reinforcement learning for dynamic chute mapping. In *Proc. International Conference on Machine Learning (ICML)*, 2025.
- Qinghua Liu, Tiancheng Yu, Yu Bai, and Chi Jin. A sharp analysis of model-based reinforcement learning with self-play. In *Proc. International Conference on Machine Learning (ICML)*, pp. 7001–7010. PMLR, 2021.
- Zhishuai Liu and Pan Xu. Minimax Optimal and Computationally Efficient Algorithms for Distributionally Robust Offline Reinforcement Learning. *Advances in Neural Information Processing Systems*, 37:86602–86654, 2024.
- Zhishuai Liu, Weixin Wang, and Pan Xu. Upper and lower bounds for distributionally robust off-dynamics reinforcement learning. *arXiv preprint arXiv:2409.20521*, 2024.
- Zijian Liu, Qinxun Bai, Jose Blanchet, Perry Dong, Wei Xu, Zhengqing Zhou, and Zhengyuan Zhou. Distributionally robust Q-learning. In *Proc. International Conference on Machine Learning (ICML)*, pp. 13623–13643. PMLR, 2022.
- Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, pp. 6379–6390, 2017.
- Miao Lu, Han Zhong, Tong Zhang, and Jose Blanchet. Distributionally robust reinforcement learning with interactive data collection: Fundamental hardness and near-optimal algorithm. *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

- MingYu Lu, Zachary Shahn, Daby Sow, Finale Doshi-Velez, and Li-wei H Lehman. Is deep reinforcement learning ready for practical applications in healthcare? A sensitivity analysis of duel-DDQN for hemodynamic management in sepsis patients. In *AMIA annual symposium proceedings*, volume 2020, pp. 773, 2021.
- Shaocong Ma, Ziyi Chen, Shaofeng Zou, and Yi Zhou. Decentralized robust v-learning for solving markov games with model uncertainty. *Journal of Machine Learning Research*, 24(371):1–40, 2023.
- Xiaoteng Ma, Zhipeng Liang, Jose Blanchet, Mingwen Liu, Li Xia, Jiheng Zhang, Qianchuan Zhao, and Zhengyuan Zhou. Distributionally Robust Offline Reinforcement Learning with Linear Function Approximation. *arXiv preprint arXiv:2209.06620*, 2022.
- Weichao Mao and Tamer Başar. Provably efficient reinforcement learning in decentralized general-sum markov games. *Dynamic Games and Applications*, 13(1):165–186, 2023.
- Laetitia Matignon, Guillaume J Laurent, and Nadine Le Fort-Piat. Independent reinforcement learners in cooperative markov games: a survey regarding coordination problems. *The Knowledge Engineering Review*, 27(1):1–31, 2012.
- Andreas Maurer and Massimiliano Pontil. Empirical Bernstein Bounds and Sample Variance Penalization. *arXiv preprint arXiv:0907.3740*, 2009.
- Pierre Ménard, Omar Darwiche Domingues, Xuedong Shang, and Michal Valko. UCB Momentum Q-learning: Correcting the bias without forgetting. In *International Conference on Machine Learning*, pp. 7609–7618. PMLR, 2021.
- Arnab Nilim and Laurent El Ghaoui. Robust Control of Markov Decision Processes with Uncertain Transition Matrices. *Operations Research*, 53(5):780–798, 2005.
- Afshin Oroojlooy and Davood Hajinezhad. A review of cooperative multi-agent deep reinforcement learning. *Applied Intelligence*, 53(11):13677–13722, 2023.
- Sindhu Padakandla, Prabuchandran KJ, and Shalabh Bhatnagar. Reinforcement learning algorithm for non-stationary environments. *Applied Intelligence*, 50(11):3590–3606, 2020.
- Kishan Panaganti and Dileep Kalathil. Sample Complexity of Robust Reinforcement Learning with a Generative Model. In *International Conference on Artificial Intelligence and Statistics*, pp. 9582–9602. PMLR, 2022.
- Kishan Panaganti, Zaiyan Xu, Dileep Kalathil, and Mohammad Ghavamzadeh. Robust reinforcement learning using offline data. *arXiv preprint arXiv:2208.05129*, 2022.
- Georgios Papoudakis, Filippos Christianos, Arrasy Rahman, and Stefano V Albrecht. Dealing with non-stationarity in multi-agent deep reinforcement learning. *arXiv preprint arXiv:1906.04737*, 2019.
- Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE international conference on robotics and automation (ICRA)*, pp. 3803–3810. IEEE, 2018.
- Aravind Rajeswaran, Sarvjeet Ghotra, Balaraman Ravindran, and Sergey Levine. Epopt: Learning robust neural network policies using model ensembles. *arXiv preprint arXiv:1610.01283*, 2016.
- Shyam Sundhar Ramesh, Pier Giuseppe Sessa, Yifan Hu, Andreas Krause, and Ilija Bogunovic. Distributionally robust model-based reinforcement learning with large state spaces, 2023. URL <https://arxiv.org/abs/2309.02236>.
- Igal Sason and Sergio Verdú.  $f$ -divergence inequalities. *IEEE Transactions on Information Theory*, 62(11):5973–6006, 2016.
- Shai Shalev-Shwartz et al. Safe, multi-agent, reinforcement learning for autonomous driving. *arXiv preprint arXiv:1610.03295*, 2016.

- Lloyd S Shapley. Stochastic games. *Proceedings of the national academy of sciences*, 39(10): 1095–1100, 1953.
- Laixi Shi and Yuejie Chi. Distributionally Robust Model-Based Offline Reinforcement Learning with Near-Optimal Sample Complexity. *Journal of Machine Learning Research*, 25(200):1–91, 2024.
- Laixi Shi, Gen Li, Yuting Wei, Yuxin Chen, Matthieu Geist, and Yuejie Chi. The Curious Price of Distributional Robustness in Reinforcement Learning with a Generative Model. *Advances in Neural Information Processing Systems*, 36:79903–79917, 2023.
- Laixi Shi, Jingchu Gai, Eric Mazumdar, Yuejie Chi, and Adam Wierman. Breaking the curse of multiagency in robust multi-agent reinforcement learning. *arXiv preprint arXiv:2409.20067*, 2024a.
- Laixi Shi, Eric Mazumdar, Yuejie Chi, and Adam Wierman. Sample-Efficient Robust Multi-Agent Reinforcement Learning in the Face of Environmental Uncertainty. *arXiv preprint arXiv:2404.18909*, 2024b.
- David Silver, Aja Huang, Christopher Maddison, Arthur Guez, Laurent Sifre, George Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of go with deep neural networks and tree search. *Nature*, 529:484–489, 01 2016. doi: 10.1038/nature16961.
- Ziang Song, Song Mei, and Yu Bai. When can we learn general-sum markov games with a large number of players sample-efficiently? *arXiv preprint arXiv:2110.04184*, 2021.
- Eugene Vinitsky, Yuqing Du, Kanaad Parvate, Kathy Jang, Pieter Abbeel, and Alexandre Bayen. Robust reinforcement learning using adversarial populations. *arXiv preprint arXiv:2008.01825*, 2020.
- Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, L. Sifre, Trevor Cai, John P. Agapiou, Max Jaderberg, Alexander Sasha Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden, Yury Sulsky, James Molloy, Tom Le Paine, Caglar Gulcehre, Ziyun Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul, Timothy P. Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575:350–354, 2019. URL <https://api.semanticscholar.org/CorpusID:204972004>.
- He Wang, Laixi Shi, and Yuejie Chi. Sample complexity of offline distributionally robust linear markov decision processes. *arXiv preprint arXiv:2403.12946*, 2024a.
- Qiu hao Wang, Chin Pang Ho, and Marek Petrik. Policy gradient in robust mdps with global convergence guarantee. In *Proc. International Conference on Machine Learning (ICML)*, pp. 35763–35797. PMLR, 2023a.
- Shengbo Wang, Nian Si, Jose Blanchet, and Zhengyuan Zhou. A Finite Sample Complexity Bound for Distributionally Robust Q-learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 3370–3398. PMLR, 2023b.
- Shengbo Wang, Nian Si, Jose Blanchet, and Zhengyuan Zhou. Sample Complexity of Variance-Reduced Distributionally Robust Q-Learning. *Journal of Machine Learning Research*, 25(341):1–77, 2024b.
- Shengbo Wang, Nian Si, Jose Blanchet, and Zhengyuan Zhou. On the foundation of distributionally robust reinforcement learning, 2024c. URL <https://arxiv.org/abs/2311.09018>.
- Yudan Wang, Yue Wang, Yi Zhou, Alvaro Velasquez, and Shaofeng Zou. Data-driven robust multi-agent reinforcement learning. In *2022 IEEE 32nd International Workshop on Machine Learning for Signal Processing (MLSP)*, pp. 1–6. IEEE, 2022.

- Yudan Wang, Shaofeng Zou, and Yue Wang. Model-free robust reinforcement learning with sample complexity analysis. In *Proc. International Conference on Uncertainty in Artificial Intelligence (UAI)*, 2024d.
- Yue Wang and Shaofeng Zou. Online Robust Reinforcement Learning with Model Uncertainty. *Advances in Neural Information Processing Systems*, 34:7193–7206, 2021.
- Yue Wang, Alvaro Velasquez, George K Atia, Ashley Prater-Bennette, and Shaofeng Zou. Model-free robust average-reward reinforcement learning. In *Proc. International Conference on Machine Learning (ICML)*, pp. 36431–36469. PMLR, 2023c.
- Yue Wang, Zhongchang Sun, and Shaofeng Zou. A Unified Principle of Pessimism for Offline Reinforcement Learning under Model Mismatch. *Advances in Neural Information Processing Systems*, 37:9281–9328, 2024e.
- Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem. Robust Markov Decision Processes. *Mathematics of Operations Research*, 38(1):153–183, 2013.
- Annie Wong, Thomas Bäck, Anna V Kononova, and Aske Plaat. Deep multiagent reinforcement learning: Challenges and directions. *Artificial Intelligence Review*, 56(6):5023–5056, 2023.
- Qiaomin Xie, Yudong Chen, Zhaoran Wang, and Zhuoran Yang. Learning zero-sum simultaneous-move markov games using function approximation and correlated equilibrium. In *Proc. Annual Conference on Learning Theory (CoLT)*, pp. 3674–3682. PMLR, 2020.
- Zaiyan Xu, Kishan Panaganti, and Dileep Kalathil. Improved Sample Complexity Bounds for Distributionally Robust Reinforcement Learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 9728–9754. PMLR, 2023.
- Wenhao Yang, Liangyu Zhang, and Zhihua Zhang. Toward Theoretical Understandings of Robust Markov Decision Processes: Sample Complexity and Asymptotics. *The Annals of Statistics*, 50(6): 3223–3248, 2022.
- Wenhao Yang, Han Wang, Tadashi Kozuno, Scott M Jordan, and Zhihua Zhang. Robust markov decision processes without model estimation. *arXiv preprint arXiv:2302.01248*, 2023.
- Andrea Zanette and Emma Brunskill. Tighter Problem-Dependent Regret Bounds in Reinforcement Learning without Domain Knowledge using Value Function Bounds. In *International Conference on Machine Learning*, pp. 7304–7312. PMLR, 2019.
- Chi Zhang, Zain Ulabedeen Farhat, George K. Atia, and Yue Wang. Model-free offline reinforcement learning with enhanced robustness. In *Proc. International Conference on Learning Representations (ICLR)*, 2025.
- Kaiqing Zhang, Tao Sun, Yunzhe Tao, Sahika Genc, Sunil Mallya, and Tamer Basar. Robust multi-agent reinforcement learning with model uncertainty. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, 2020.
- Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of reinforcement learning and control*, pp. 321–384, 2021a.
- Runyu Zhang, Yang Hu, and Na Li. Soft Robust MDPs and Risk-Sensitive MDPs: Equivalence, Policy Gradient, and Sample Complexity. *arXiv preprint arXiv:2306.11626*, 2023.
- Zihan Zhang, Xiangyang Ji, and Simon Du. Is Reinforcement Learning More Difficult Than Bandits? A Near-optimal Algorithm Escaping the Curse of Horizon. In *Conference on Learning Theory*, pp. 4528–4531. PMLR, 2021b.
- Zihan Zhang, Yuxin Chen, Jason D Lee, and Simon S Du. Settling the sample complexity of online reinforcement learning. In *Proc. Annual Conference on Learning Theory (CoLT)*, pp. 5213–5219. PMLR, 2024.

- Wenshuai Zhao, Jorge Peña Queralta, and Tomi Westerlund. Sim-to-real transfer in deep reinforcement learning for robotics: a survey. In *2020 IEEE symposium series on computational intelligence (SSCI)*, pp. 737–744. IEEE, 2020.
- Zhengqing Zhou, Zhengyuan Zhou, Qinxun Bai, Linhai Qiu, Jose Blanchet, and Peter Glynn. Finite-Sample Regret Bound for Distributionally Robust Offline Tabular Reinforcement Learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 3331–3339. PMLR, 2021.
- Ziyuan Zhou, Guanjun Liu, and Mengchu Zhou. A robust mean-field actor-critic reinforcement learning against adversarial perturbations on agent states. *IEEE Transactions on Neural Networks and Learning Systems*, 35(10):14370–14381, October 2024. ISSN 2162-2388. doi: 10.1109/tnnls.2023.3278715. URL <http://dx.doi.org/10.1109/TNNLS.2023.3278715>.



## A USE OF LARGE LANGUAGE MODELS

We used ChatGPT only as a general-purpose assistant for language editing and typesetting. Its role was limited to (i) improving grammar, style, and readability, and (ii) LaTeX support—adjusting algorithm placement, tidying BibTeX entries and citation styles, and resolving compile issues (e.g., Type-3 font warnings and package conflicts). All ideas, derivations, and final claims were conceived, checked, and validated by the authors, who bear full responsibility for the paper’s content.

## B RELATED WORKS

In this section we discuss other related works.

**Single-Agent Robust RL.** Robust RL for single-agent settings has been extensively studied across a wide range of formulations. In particular, a substantial body of work has examined the generative-model setting (Clavier et al., 2023; Liu et al., 2022; Panaganti & Kalathil, 2022; Ramesh et al., 2023; Shi et al., 2023; Wang et al., 2023b; 2024c;b; Xu et al., 2023; Yang et al., 2022; 2023), where the agent is assumed to have access to a simulator. These studies develop distributionally robust RL algorithms under various uncertainty sets, including TV, KL,  $\chi^2$ , and Wasserstein divergences. Another, and arguably more challenging, line of research focuses on the offline setting (Blanchet et al., 2023; Ma et al., 2022; Panaganti et al., 2022; Shi & Chi, 2024; Zhang et al., 2023; Liu & Xu, 2024; Wang et al., 2024e; Blanchet et al., 2023; Wang et al., 2024a). In this setting, the agent must learn exclusively from a fixed offline dataset, without the ability to collect additional online samples. Finally, we consider the online setting (Badrinath & Kalathil, 2021; Dong et al., 2022; Li et al., 2022; Liang et al., 2023; Wang & Zou, 2021), where the agent learns exclusively through direct interaction with the environment. Prior work spans model-based, model-free, and policy-gradient approaches, with some methods, such as the policy optimization algorithm of (Dong et al., 2022), achieving sublinear regret guarantees.

**Robust MARL.** Besides the distributionally robust Markov games we considered in our paper, there are also other works investigate robustness in MARL for cooperative tasks, where all agents share a unified objective. (Bukharin et al., 2023) enhance robustness through adversarial regularization, perturbing the environment to encourage Lipschitz-continuous policies. (Lin et al., 2020) explore adversarial attacks on MARL agents as a means of improving resilience, while (Li et al., 2019) extend this approach to continuous action spaces by modifying the MADDPG algorithm (Lowe et al., 2017) to focus on worst-case actions—a narrower interpretation of worst-case optimization in robust RL. (Wang et al., 2022) studied robust MARL with network agents.

Another line of research focuses on the robustness in MARL under observation uncertainty, under the formulation of partially observable MDPs. The framework of observation-robust games is proposed in (He et al., 2023; Han et al., 2024). Observation-robust cooperative MARL is studied in (Zhou et al., 2024).

**Non-Robust Markov Games.** Markov games (MGs), or stochastic games, introduced by (Shapley, 1953), form the standard foundation for multi-agent reinforcement learning (MARL), particularly in equilibrium learning. Comprehensive surveys such as (Busoniu et al., 2008; Oroojlooy & Hajinezhad, 2023; Zhang et al., 2021a) offer thorough coverage of the field’s evolution. Early work in MARL focused on asymptotic convergence guarantees (Littman et al., 2001; Littman & Szepesvári, 1996), whereas recent research emphasizes finite-sample analyses to establish non-asymptotic guarantees, especially for learning Nash equilibria (NE)—a central solution concept. The existence of NE in general-sum MGs was shown by (Fink, 1964), and the algorithmic foundation was laid by the seminal work of (Littman, 1994). Classical algorithms such as Nash-Q (Hu & Wellman, 2003), FF-Q (Littman et al., 2001), and correlated-Q learning (Greenwald et al., 2003) were proposed to compute NE and its variants. However, computing NE in general-sum multi-player settings remains PPAD-complete (Daskalakis, 2013), and no polynomial-time algorithms exist for this case (Jin et al., 2022; Deng et al., 2023). In contrast, the two-player zero-sum setting admits tractable solutions, with the first polynomial-time algorithm developed by (Hansen et al., 2013). To address the computational intractability in general-sum MGs, attention has shifted to weaker notions like CE and CCE, with polynomial-time algorithms such as V-learning (Jin et al., 2021; Mao & Başar, 2023; Song et al., 2021) and Nash value iteration (Liu et al., 2021) enabling efficient computation. Furthermore, significant progress in finite-sample analysis—spanning both model-based

and model-free algorithms—has been achieved in the two-player zero-sum setting, as evidenced by (Bai & Jin, 2020; Xie et al., 2020; Cui et al., 2023; Chen et al., 2022; Liu et al., 2021; Feng et al., 2023; Li et al., 2024), advancing the theoretical understanding of equilibrium learning in standard MARL without robustness considerations.

## C DRMGM WITH $f$ -DIVERGENCE UNCERTAINTY SET

We review the formulation of DRMGM with  $f$ -divergence uncertainty sets. This framework operates under the  $\mathcal{S} \times \mathcal{A}$ -rectangularity assumption, where the nominal transition probability  $P^*$  and the agent-specific radius  $\rho_i$  for  $i \in \mathcal{M}$  define the robust problem as per Definition 1.

**Lemma 1** (Strong duality for  $f$ -divergence). *Let  $\mathcal{P}_f^{\rho_i}(s, \mathbf{a})$  be an  $f$ -divergence uncertainty set as defined in Definition 1. For any value function  $V_i : \mathcal{S} \rightarrow \mathbb{R}_+$  and a nominal transition kernel  $P^* : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ , the worst-case expected value,  $\sigma_{\mathcal{P}_f^{\rho_i}(s, \mathbf{a})}[V_i] := \inf_{P \in \mathcal{P}_f^{\rho_i}(s, \mathbf{a})} [\mathbb{P}V_i](s, \mathbf{a})$ , admits a dual representation given by:*

$$\sigma_{\mathcal{P}_{i,h,f}^{\rho_i}(s, \mathbf{a})}[V] = \sup_{\lambda \geq 0, \eta \in \mathbb{R}} \left\{ -\lambda \sum_{s \in \mathcal{S}} P^*(s) f^* \left( \frac{\eta - V(s)}{\lambda} \right) - \lambda \rho_i + \eta \right\},$$

where  $f^*$  is the convex conjugate of  $f$ .

The detailed proof is given in Lemma B.1 of (Yang et al., 2022).

**Corollary 2** (Dual representation for TV and KL-divergence). *Under the assumption of  $\mathcal{S} \times \mathcal{A}$ -rectangularity, the dual representation from Lemma 1 simplifies to the following for two specific cases of  $f$ -divergence. For any value function  $V : \mathcal{S} \rightarrow [0, H]$  and a nominal distribution  $P_h^*$  over the next states:*

**TV-Divergence.** *For an uncertainty set defined by TV-divergence, where  $f(t) = \frac{1}{2}|t - 1|$ , the robust expectation  $\sigma_{\mathcal{P}_{i,h,TV}^{\rho_i}(s, \mathbf{a})}[V_i]$  is expressed as:*

$$\sigma_{\mathcal{P}_{i,h,TV}^{\rho_i}(s, \mathbf{a})}[V_i] = \sup_{\eta \in [0, H]} \left\{ -\mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} [\max(0, \eta - V_i)] - \frac{\rho}{2} \max(0, \eta - \min_{s' \in \mathcal{S}} V_i(s')) + \eta \right\}. \quad (10)$$

**KL-Divergence.** *For an uncertainty set defined by KL-divergence, with  $f(t) = t \log(t)$ , the robust expectation  $\sigma_{\mathcal{P}_{i,h,KL}^{\rho_i}(s, \mathbf{a})}[V_i]$  is expressed as:*

$$\sigma_{\mathcal{P}_{i,h,KL}^{\rho_i}(s, \mathbf{a})}[V_i] = \sup_{\eta \in [\underline{\eta}, H/\rho_i]} \left\{ -\eta \log \left( \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{V_i}{\eta} \right\} \right] \right) - \eta \rho_i \right\}. \quad (11)$$

## ROBUST BELLMAN EQUATIONS.

Analogous to standard MGs, the following proposition provides the robust Bellman equation for DRMGMs. In particular, the robust value functions  $V_{i,h}^{\pi, \rho_i}(s)$  associated with any joint policy  $\pi$  for all  $(i, h, s) \in \mathcal{M} \times [H] \times \mathcal{S}$  obeys the following proposition given below:

**Proposition 1** (Robust Bellman Equation). *Under the  $\mathcal{S} \times \mathcal{A}$ -rectangularity assumption, for any nominal transition kernel  $P^*$  and joint policy  $\pi$ , the robust Bellman equation holds for any  $(i, h, s, \mathbf{a})$ :*

$$Q_{i,h}^{\pi, \rho_i}(s, \mathbf{a}) = r_{i,h}(s, \mathbf{a}) + \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\pi, \rho_i}] \quad (12)$$

$$V_{i,h}^{\pi, \rho_i}(s) = \mathbb{E}_{\mathbf{a} \sim \pi_h(\cdot|s)} [Q_{i,h}^{\pi, \rho_i}(s, \mathbf{a})] \quad (13)$$

The detailed proof of Proposition 1 for finite-horizon RMDP is given in (Blanchet et al., 2023, Proposition 2.3). We emphasize that the robust Bellman equation in 13 is fundamentally grounded in the agent-wise  $(s, \mathbf{a})$ -rectangularity condition imposed on the uncertainty set. This condition decouples the dependencies of uncertainty across agents, state-action pairs, and time steps, thereby enabling the recursive structure of the Bellman equation.

## D HARDNESS OF MULTI-AGENT ONLINE LEARNING

### D.1 HARDNESS WITH SUPPORT SHIFT

**Example 1** (The “Initial Shock” Game). Consider a class of  $N$ -agent DRMGs,  $\{M_{\mathbf{a}^*}\}_{\mathbf{a}^* \in \mathcal{A}}$ , parameterized by a “secret escape route”  $\mathbf{a}^* \in \mathcal{A}$ .

- **Action Spaces:**  $A_i = M$  for each agent. The joint action space has size  $|\mathcal{A}| = \prod_{i \in [N]} A_i = M^N$ .
- **States, Horizon, Rewards:**  $\mathcal{S} = \{s_{\text{good}}, s_{\text{bad}}\}$ , horizon  $H$ , initial state  $s_1 = s_{\text{good}}$ , and rewards are defined as

$$r_i(s, \mathbf{a}) = \begin{cases} 1, & \text{if } s = s_{\text{good}} \text{ or if } (s = s_{\text{bad}} \text{ and } \mathbf{a} = \mathbf{a}^*) \\ 0, & \text{if } s = s_{\text{bad}} \text{ and } \mathbf{a} \neq \mathbf{a}^* \end{cases}.$$

- **Dynamics:** The system dynamics create the trap.
  - From  $\mathbf{s}_{\text{good}}$ : Nominally, the system stays in  $s_{\text{good}}$ . An adversary can force a transition to  $s_{\text{bad}}$  with probability  $\rho$ .
  - From  $\mathbf{s}_{\text{bad}}$ : This is the trap. The only way to escape is to play the secret joint action:

$$\text{Next State} = \begin{cases} s_{\text{good}}, & \text{if } \mathbf{a} = \mathbf{a}^* \\ s_{\text{bad}}, & \text{if } \mathbf{a} \neq \mathbf{a}^* \end{cases}.$$

- **Uncertainty Set:** The uncertainty is non-zero only at the first step.
  - At  $h = 1$  and  $s_1 = s_{\text{good}}$ : The uncertainty set is a TV-ball with radius  $\rho$ .
  - For all  $h > 1$  or  $s \neq s_{\text{good}}$ : There is no uncertainty ( $\rho = 0$ ). The transition is the nominal one.

**Theorem 3.** For the “Initial Shock” DRMG, any decentralized online learning algorithm suffers the following best-response regret lower bound:

$$\inf_{\mathcal{ALG}} \sup_{\mathbf{a}^* \in \mathcal{A}} \mathbb{E}[\text{Regret}_i(K)] \geq \Omega \left( \rho K \cdot \min \left\{ H, \prod_{i \in [N]} A_i \right\} \right).$$

*Proof. Step 1: Decomposing the Per-Episode Regret.* The best-response regret for Agent 1 in an episode is  $\text{Regret}_1^k = V_{1,1}^{\dagger, \pi^{-i}, \rho} - V_{1,1}^{\pi, \rho}$ . We expand this using the robust Bellman equation at  $s_1 = s_{\text{good}}$ , where uncertainty exists.

$$\begin{aligned} \text{Regret}_1^k &= \left( 1 + (1 - \rho) V_{1,2}^{\dagger, \pi^{-i}, \rho}(s_{\text{good}}) + \rho V_{1,2}^{\dagger, \pi^{-i}, \rho}(s_{\text{bad}}) \right) \\ &\quad - \left( 1 + (1 - \rho) V_{1,2}^{\pi, \rho}(s_{\text{good}}) + \rho V_{1,2}^{\pi, \rho}(s_{\text{bad}}) \right) \\ &= (1 - \rho) \left( V_{1,2}^{\dagger, \pi^{-i}, \rho}(s_{\text{good}}) - V_{1,2}^{\pi, \rho}(s_{\text{good}}) \right) + \rho \left( V_{1,2}^{\dagger, \pi^{-i}, \rho}(s_{\text{bad}}) - V_{1,2}^{\pi, \rho}(s_{\text{bad}}) \right). \end{aligned}$$

Since there is no uncertainty for  $h > 1$ , the transition from  $s_{\text{good}}$  at  $h = 2$  is deterministically to  $s_{\text{good}}$  at  $h = 3$ . Thus,  $V_{1,2}(s_{\text{good}})$  is a constant independent of the policy in the trap state, which means  $V_{1,2}^{\dagger, \pi^{-i}, \rho}(s_{\text{good}}) = V_{1,2}^{\pi, \rho}(s_{\text{good}})$ . The first term is exactly zero, and thus we have that

$$\text{Regret}_1^k = \rho \left( V_{1,2}^{\dagger, \pi^{-i}, \rho}(s_{\text{bad}}) - V_{1,2}^{\pi, \rho}(s_{\text{bad}}) \right) = \rho \cdot \Delta V_2^{\rho}(s_{\text{bad}}). \quad (14)$$

**Step 2: Formalizing the Value Gap  $\Delta V_2^{\rho}(s_{\text{bad}})$ .** The value gap is the expected difference in total future rewards. This difference is precisely the expected number of steps wasted in the trap. Note that the value of state  $s_{\text{bad}}$  at step  $h$  under a policy  $\pi'$  is the expected sum of future rewards. Let  $\tau = \tau(\pi')$  be the random variable for the number of steps to escape (i.e., play  $\mathbf{a}^*$ ), starting from step  $h$ . Let  $C = H - h + 1$  be the number of steps remaining in the episode, then the total reward

collected from  $h = 2$  is  $V_{1,2}^{\pi', \rho}(s_{bad}) = \mathbb{E}[\mathbb{I}[\tau \leq C] \cdot (C - \tau + 2)]$  as it will always receive  $r = 1$  when at  $s_{good}$ .

Moreover, note that the total number of available rewards is  $C$ , and since  $C = \min(\tau - 1, C) + \mathbb{I}[\tau \leq C](C - \tau + 1)$ , the value can therefore be expressed as  $V_{1,2}^{\pi', \rho}(s_{bad}) = C - \mathbb{E}[\min(\tau - 1, C)]$ .

Therefore, the value gap is the difference in the expected number of wasted steps:

$$\begin{aligned} \Delta V_2^\rho(s_{bad}) &= (C - \mathbb{E}[\min(\tau^* - 1, C)]) - (C - \mathbb{E}[\min(\tau - 1, C)]) \\ &= \mathbb{E}[\min(\tau - 1, C)] - \mathbb{E}[\min(\tau^* - 1, C)]. \end{aligned}$$

where  $\tau^*$  is the escape probability of  $\pi^*$ . Since the best-response policy  $\pi_1^*$  plays  $a_1^*$  deterministically, so its escape time  $\tau^*$  depends only on the other agents' policies,  $\pi_{-1}$ . The algorithm's escape time  $\tau$  depends on its full policy  $\pi$ .

**Step 3: Lower Bounding the Value Gap.** The best response for Agent 1 is to play  $a_1^*$ , so  $\tau^*$  does not involve any search for Agent 1. In contrast,

However, the algorithm does not know  $a_1^*$  and must search. We are interested in the worst-case regret over the choice of  $a^*$ . The expected wasted steps for the algorithm is  $\mathbb{E}[\min(\tau - 1, C)]$ . Let  $p_1 = \Pr_{\pi_1}(a_1 = a_1^*)$  and  $p_{-1} = \Pr_{\pi_{-1}}(a_{-1} = a_{-1}^*)$ . The algorithm's one-step escape probability is  $p_1 \cdot p_{-1}$ . Its expected escape time is  $\mathbb{E}[\tau] = 1/(p_1 \cdot p_{-1})$ . The expected wasted steps is lower-bounded by:

$$\mathbb{E}[\min(\tau - 1, C)] \geq \Omega(\min(\mathbb{E}[\tau - 1], C)) = \Omega(\min(1/(p_1 \cdot p_{-1}), H - 1)),$$

where the inequality is due to Lemma 2.

In the worst case over the unknown  $a^*$ , the probabilities  $p_1$  and  $p_{-1}$  are minimized:

$$\inf_{a_1^*} p_1 \leq 1/A_1 \quad \text{and} \quad \inf_{a_{-1}^*} p_{-1} \leq 1 / \left( \prod_{i=2}^N A_i \right).$$

The best-response policy suffers much less waste. Thus, the value gap  $\Delta V_2^\rho(s_{bad})$  is dominated by the algorithm's large number of wasted steps.

$$\sup_{a^*} \Delta V_2^\rho(s_{bad}) \geq \Omega \left( \min \left\{ 1 / \left( (1/A_1) \cdot (1 / \left( \prod_{i=2}^N A_i \right)) \right), H \right\} \right) = \Omega \left( \min \left\{ \prod_{i=1}^N A_i, H \right\} \right).$$

**Step 4: Finalizing the Bound.** Substituting this back into the per-episode regret expression from Step 1:

$$\sup_{a^*} \mathbb{E}[\text{Regret}_1^k] \geq \rho \cdot \Omega \left( \min \left\{ \prod_{i=1}^N A_i, H \right\} \right).$$

This per-episode regret is incurred because the information bottleneck prevents the algorithm from learning  $a^*$ . Summing over  $K$  episodes gives the final total regret bound:

$$\inf_{\mathcal{ALG}} \sup_{a^*} \mathbb{E}[\text{Regret}_1(K)] = \sum_{k=1}^K \sup_{a^*} \mathbb{E}[\text{Regret}_1^k] \geq \Omega \left( \rho K \cdot \min \left\{ \prod_{i=1}^N A_i, H \right\} \right).$$

This completes the proof.  $\square$

**Lemma 2.** Let  $\tau$  be the random variable for the escape time from the trap state, and let  $C = H - 1$  be the number of steps remaining in the episode. The true expected number of wasted steps,  $\mathbb{E}[\min(\tau - 1, C)]$ , has the following asymptotic lower bound:

$$\mathbb{E}[\min(\tau - 1, C)] \geq \Omega(\min(\mathbb{E}[\tau - 1], C)).$$

*Proof.* Note that  $\tau$  follows a Geometric distribution  $\tau \sim \text{Geo}(p)$  and have the probability mass function  $P(\tau = k) = (1 - p)^{k-1}p$  for  $k \in \{1, 2, 3, \dots\}$ . The random variable  $\tau - 1$  represents the number of failures before the first success. Its expectation is  $\mathbb{E}[\tau - 1] = \frac{1-p}{p}$ .

We first derive an expression for  $\mathbb{E}[\min(\tau - 1, C)]$ . We use the tail sum formula for the expectation of a non-negative, integer-valued random variable  $X$ , which states  $\mathbb{E}[X] = \sum_{k=0}^{\infty} P(X > k)$ .

Let  $X = \min(\tau - 1, C)$ . The event  $\{X > k\}$  is equivalent to the event  $\{\tau - 1 > k \text{ and } C > k\}$ .

- If  $k \geq C$ , then  $P(X > k) = 0$ .
- If  $k < C$ , then  $P(X > k) = P(\tau - 1 > k)$ .

The event  $\{\tau - 1 > k\}$  means the first  $k + 1$  trials resulted in failure, so its probability is  $P(\tau > k + 1) = (1 - p)^{k+1}$ .

The expectation is therefore the sum over the non-zero probabilities:

$$\begin{aligned}\mathbb{E}[\min(\tau - 1, C)] &= \sum_{k=0}^{\infty} P(\min(\tau - 1, C) > k) \\ &= \sum_{k=0}^{C-1} P(\tau - 1 > k) = \sum_{k=0}^{C-1} (1 - p)^{k+1}.\end{aligned}$$

Letting  $q = 1 - p$ , this is a finite geometric series:

$$\sum_{j=1}^C q^j = q \frac{1 - q^C}{1 - q} = \frac{q(1 - q^C)}{p}.$$

Substituting  $q = 1 - p$  back, we express the expectation in terms of  $\mathbb{E}[\tau - 1]$ :

$$\mathbb{E}[\min(\tau - 1, C)] = \frac{1 - p}{p} (1 - (1 - p)^C) = \mathbb{E}[\tau - 1] (1 - (1 - p)^C).$$

Let  $\mu = \mathbb{E}[\tau - 1] = \frac{1-p}{p}$ . We want to show that there exists a universal constant  $k > 0$  such that:

$$\mu(1 - (1 - p)^C) \geq k \cdot \min(\mu, C).$$

We proceed with a case analysis based on the relationship between  $\mu$  and  $C$ .

**Case 1:**  $\mu \leq C$ : In this case,  $\min(\mu, C) = \mu$ . We need to show that  $\mu(1 - (1 - p)^C) \geq k \cdot \mu$ , which simplifies to proving that  $1 - (1 - p)^C \geq k$ .

The condition  $\mu \leq C$  implies a lower bound on  $p$ :

$$\frac{1 - p}{p} \leq C \implies 1 - p \leq Cp \implies 1 \leq (C + 1)p \implies p \geq \frac{1}{C + 1}.$$

Using the standard inequality  $1 - x \leq e^{-x}$ , we have  $(1 - p)^C \leq e^{-pC}$ . Thus,

$$1 - (1 - p)^C \geq 1 - e^{-pC}.$$

Since  $p \geq \frac{1}{C+1}$ , we have  $pC \geq \frac{C}{C+1}$ . As the function  $f(x) = 1 - e^{-x}$  is increasing for  $x > 0$ ,

$$1 - e^{-pC} \geq 1 - e^{-C/(C+1)}.$$

The function  $g(C) = \frac{C}{C+1}$  is increasing for  $C \geq 1$ , with a minimum value of  $g(1) = 1/2$ . Therefore, for any integer  $C \geq 1$ ,

$$1 - (1 - p)^C \geq 1 - e^{-1/2}.$$

Thus, the inequality holds in this case with the constant  $k_1 = 1 - e^{-1/2} \approx 0.393$ .

**Case 2:**  $\mu > C$ : In this case,  $\min(\mu, C) = C$ . We need to show that  $\mu(1 - (1 - p)^C) \geq kC$ .

The condition  $\mu > C$  implies an upper bound on  $p$ :

$$\frac{1 - p}{p} > C \implies 1 - p > Cp \implies 1 > (C + 1)p \implies p < \frac{1}{C + 1}.$$

From our calculation of the expectation, we have a sum of  $C$  positive, decreasing terms:

$$\mathbb{E}[\min(\tau - 1, C)] = \sum_{k=0}^{C-1} (1 - p)^{k+1}.$$

This sum is greater than  $C$  times its smallest term, which is  $(1-p)^C$ :

$$\mathbb{E}[\min(\tau - 1, C)] > C(1-p)^C.$$

From the condition  $p < \frac{1}{C+1}$ , it follows that  $1-p > 1 - \frac{1}{C+1} = \frac{C}{C+1}$ . Therefore,

$$\mathbb{E}[\min(\tau - 1, C)] > C \left( \frac{C}{C+1} \right)^C = C \left( 1 - \frac{1}{C+1} \right)^C.$$

The sequence  $a_C = \left( 1 - \frac{1}{C+1} \right)^C$  is decreasing for  $C \geq 1$ , and its limit as  $C \rightarrow \infty$  is  $1/e$ . Hence, for all  $C \geq 1$ , the sequence is bounded below by its limit:

$$\left( 1 - \frac{1}{C+1} \right)^C \geq \lim_{n \rightarrow \infty} \left( 1 - \frac{1}{n+1} \right)^n = \frac{1}{e}.$$

This gives the lower bound:

$$\mathbb{E}[\min(\tau - 1, C)] > C \cdot \frac{1}{e}.$$

So, the inequality holds in this case with the constant  $k_2 = 1/e \approx 0.368$ . By combining the two cases, the inequality is shown to hold for a universal constant  $k = \min(k_1, k_2) = \min(1 - e^{-1/2}, 1/e) = 1/e$ .

Therefore, for all  $p \in (0, 1)$  and integers  $C \geq 1$ , we have established that:

$$\mathbb{E}[\min(\tau - 1, C)] \geq \frac{1}{e} \min(\mathbb{E}[\tau - 1], C) = \Omega(\min(\mathbb{E}[\tau - 1], C)),$$

which hence completes the proof.  $\square$

## D.2 HARDNESS WITHOUT SUPPORT SHIFT

**Example 2** (The ‘‘Robust Corrupted Bandit’’ Game). *Consider a class of  $N$ -agent DRMGs,  $\{M_\theta\}_{\theta \in \mathcal{A}}$ , where each game is parameterized by a secret ‘‘best’’ joint action  $\theta \in \mathcal{A}$ .*

- **States and Horizon:** A single state,  $s$ , and horizon  $H = 1$ . This reduces the problem to a one-shot game, equivalent to a multi-armed bandit setting where each episode corresponds to a single step or arm pull.
- **Action Spaces:** The joint action space  $\mathcal{A}$  is the set of arms, with size  $|\mathcal{A}| = \prod_{i=1}^N A_i$ .
- **Reward Function** ( $R \in \{0, 1\}$ ): The rewards are stochastic. Let  $\epsilon \in (0, 1/2)$  be a small constant. The nominal model  $M_\theta$  defines the following Bernoulli reward distributions for any agent  $i$ :

$$\mathbb{E}[R_i(s, \mathbf{a}) | M_\theta] = \begin{cases} 1/2 + \epsilon, & \text{if } \mathbf{a} = \theta \\ 1/2, & \text{if } \mathbf{a} \neq \theta. \end{cases}$$

- **KL-Divergence Uncertainty Set:** The true reward distribution for an action  $\mathbf{a}$ , denoted  $\tilde{P}(\cdot | \mathbf{a})$ , can be any distribution that is close to the nominal one  $P^*(\cdot | \mathbf{a})$ :

$$\mathcal{P}_{i,h,KL}^{\rho_i}(\cdot, \mathbf{a}) = \left\{ \tilde{P} : \text{KL}(\tilde{P}(\cdot | \mathbf{a}) \| P_{M_\theta}(\cdot | \mathbf{a})) \leq \rho_i, \forall \mathbf{a} \in \mathcal{A} \right\}.$$

*This uncertainty set does not have a support shift.*

The learning problem is to identify the best arm  $\theta$  by observing noisy rewards that are actively corrupted by an adversary.

**Theorem 4** (Lower Bound for Robust Learning without Support Shift). For the ‘‘Robust Corrupted Bandit’’ game, any learning algorithm suffers the following cumulative regret lower bound over  $K$  episodes (steps):

$$\inf_{\mathcal{ALG}} \sup_{\theta \in \mathcal{A}} \mathbb{E}[\text{Regret}_i(K)] \geq \Omega \left( \sqrt{\prod_{i=1}^N A_i K} \right).$$

*Proof.* The proof proceeds by a formal reduction to the classic multi-armed bandit (MAB) problem.

Let  $\mathcal{M}_\rho = \{M_{\theta,\rho}\}_{\theta \in \mathcal{A}}$  denote the class of robust game instances from our example, with uncertainty radius  $\rho > 0$ . Let  $\mathcal{M}_0 = \{M_{\theta,0}\}_{\theta \in \mathcal{A}}$  be the corresponding class of non-robust instances, where the uncertainty radius is zero and the rewards are always drawn from the nominal distributions.

Note that since the horizon  $H = 1$ , the robust problem reduces to a non-robust one, and thus the worst-case regret over the robust class  $\mathcal{M}_\rho$  must be at least as high as the worst-case regret over the non-robust class  $\mathcal{M}_0$ :

$$\mathbb{E}[\text{Regret}(K; M_{\theta,\rho})] \geq \mathbb{E}[\text{Regret}(K; M_{\theta,0})].$$

And thus

$$\inf_{\mathcal{ALG}} \sup_{\theta \in \mathcal{A}} \mathbb{E}[\text{Regret}(K; M_{\theta,\rho})] \geq \inf_{\mathcal{ALG}} \sup_{\theta \in \mathcal{A}} \mathbb{E}[\text{Regret}(K; M_{\theta,0})]. \quad (15)$$

Therefore, we can establish a lower bound for the robust problem by proving one for the simpler non-robust case.

The non-robust problem instance,  $\mathcal{M}_0$ , is a classic stochastic multi-armed bandit problem with  $M = |\mathcal{A}|$  arms. A foundational result in this area provides a strong lower bound on regret.

Note that following standard lemma:

**Lemma 3.** (Auer et al., 2002) *For any integer  $M \geq 2$  and  $K > M$ , and for any bandit algorithm, there exists a multi-armed bandit problem instance with  $M$  arms whose reward distributions are supported on  $[0, 1]$ , such that the expected cumulative regret after  $K$  steps is lower-bounded by:*

$$\mathbb{E}[\text{Regret}(K)] \geq \Omega(\sqrt{MK}).$$

We apply the lemma to our non-robust problem instance  $\mathcal{M}_0$ .

- The number of arms,  $M$ , is the size of the joint action space,  $|\mathcal{A}|$ .
- The number of steps is  $K$ .
- The reward distributions (Bernoulli) are supported on  $[0, 1]$ .

The conditions of the lemma are met. Therefore, for the class of problems  $\mathcal{M}_0$ , the worst-case regret is lower-bounded:

$$\inf_{\mathcal{ALG}} \sup_{\theta \in \mathcal{A}} \mathbb{E}[\text{Regret}(K; M_{\theta,0})] \geq \Omega\left(\sqrt{\prod_{i=1}^N A_i K}\right). \quad (16)$$

Combining the regret dominance principle from eq. 15 with the specific lower bound from eq. 16, we arrive at the final result for our robust problem:

$$\inf_{\mathcal{ALG}} \sup_{\theta \in \mathcal{A}} \mathbb{E}[\text{Regret}_i(K; M_{\theta,\rho})] \geq \Omega\left(\sqrt{\prod_{i=1}^N A_i K}\right). \quad (17)$$

This completes the formal proof by reduction. □

## E PROOF OF REGRET BOUND OF TV-MORNAVI

In this section, we prove our regret bound for TV-DRMG. Before presenting all the proofs, we first denote  $\pi^\dagger$  as the joint robust best responses over the agents, and is given by

$$\pi^\dagger = \pi_1^{\dagger,\rho_1}(\pi_{-1}) \times \cdots \times \pi_m^{\dagger,\rho_m}(\pi_{-m}). \quad (18)$$

We will use the notation of  $\pi^\dagger$  later on our proof-lines. In addition, we leverage Assumption 1, which generalizes to the case where the minimal value vanishes, i.e.,  $\min_{s \in \mathcal{S}} V(s) = 0$ , to address the support shift or extrapolation challenge arising in interactive data collection, as discussed in Remark

B.3 of (Lu et al., 2024). Consequently, this allows us to eliminate the  $\min_{s \in \mathcal{S}} V(s)$  term in the dual formulation of the TV-DRMG optimization problem, as shown in 10.

We now recall the bonus term used in TV-MORNAVI for agent  $i$  in episode  $k$  at step  $h$ , as follows:

$$\beta_{i,h}^k(s, \mathbf{a}) = \sqrt{\frac{c_1 \iota \text{Var}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \left( \frac{\bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i}}{2} \right) \right]}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{2\mathbb{E}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i} - \underline{V}_{i,h+1}^{k, \rho_i}]}{H} + \frac{c_2 H^2 S \iota}{\sqrt{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{1}{\sqrt{K}}, \quad (19)$$

where  $\iota = \log \left( S^2 (\prod_{i=1}^m A_i) H^2 K^{3/2} / \delta \right)$  and  $c_1, c_2$  are absolute constants.

We begin by defining the high-probability event  $\mathcal{E}_{TV}$ , stated in the next lemma. Our proof outline is inspired by (Lu et al., 2024) and (Ghosh et al., 2025).

**Lemma 4** (Uniform Concentration Bound of event  $\mathcal{E}_{TV}$ ). *Let  $\mathcal{E}_{TV}$  be the event in which, for all  $(s, \mathbf{a}, s', h, k) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H] \times [K]$ , and for all  $\eta$  in a  $1/(S\sqrt{K})$ -cover of  $[0, H]$ , and is defined as*

$$\mathcal{E}_{TV} := \left\{ \left| \left[ \mathbb{E}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} - \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \right] \left( \eta - V_{i,h+1}^{\dagger, \pi_{-i}^{k, \rho_i}} \right)_+ \right| \leq \sqrt{\frac{c_1 \iota \text{Var}_{\hat{P}_h^k} \left( \eta - V_{i,h+1}^{\dagger, \pi_{-i}^{k, \rho_i}} \right)_+}{N_h^k(s, \mathbf{a}) \vee 1}} + \frac{c_2 H \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}, \right. \\ \left| \hat{P}_h^k(s' | s, \mathbf{a}) - P_h^*(s' | s, \mathbf{a}) \right| \leq \sqrt{\frac{c_1 \min \{P_h^*(s' | s, \mathbf{a}), \hat{P}_h^k(s' | s, \mathbf{a})\}}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{c_2 \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}, \\ \forall (s, \mathbf{a}, s', h, k) \in \mathcal{M} \times \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H] \times [K], \forall \eta \in \mathcal{N}_{1/(S\sqrt{K})}([0, H]) \left. \right\}, \quad (20)$$

where  $\iota = \log \left( S^3 (\prod_{i=1}^m A_i) H^2 K^{3/2} / \delta \right)$ ,  $c_1, c_2 > 0$  are two absolute constants,  $\mathcal{N}_{1/(S\sqrt{K})}([0, H])$  denotes an  $1/S\sqrt{K}$ -cover of the interval  $[0, H]$ .

Then, this event  $\mathcal{E}_{TV}$  occurs with high probability, i.e.,  $\Pr(\mathcal{E}_{TV}) \geq 1 - \delta$ .

*Proof.* This proof builds upon standard techniques by applying classical concentration inequalities and a union bound. To simplify our analysis, we first consider a fixed state-action-time tuple  $(s, \mathbf{a}, h)$  within a given episode  $k$ . We can then construct an equivalent stochastic process:

- (i) Before the agents' interaction, the environment draws a sequence of next states  $\{s^{(1)}, s^{(2)}, \dots, s^{(k-1)}\}$  independently from the nominal distribution  $P_h^*(\cdot|s, \mathbf{a})$ , where  $s^{(i)} \in \mathcal{S}$  represents the state sampled in episode  $i$ .
- (ii) When the agents visit the  $(s, \mathbf{a})$  tuple at time step  $h$  for the  $i$ -th time, the environment causes a transition to the pre-sampled next state  $s^{(i)}$ .

The randomness of this constructed process is identical to that of our original, interactive learning environment. Consequently, the probability of any event is the same in both contexts. This allows us to prove the required concentration inequalities within this more tractable, simplified setting.

Leveraging this fact, we directly apply Lemma 27, which presents a variant of Bernstein's inequality and its empirical counterpart from (Maurer & Pontil, 2009). To establish a uniform bound, we apply a union bound across all tuples  $(h, s, \mathbf{a}, s', k, \eta) \in [H] \times \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [K] \times \mathcal{N}_{1/(S\sqrt{K})}([0, H])$ .

The size of this  $\epsilon$ -cover,  $\mathcal{N}_{1/(S\sqrt{K})}([0, H])$ , is on the order of  $\mathcal{O}(SH\sqrt{K})$ .  $\square$



## E.1 PROOF OF THEOREM 1 (TV-DRMG SETTING)

*Proof.* By leveraging Lemma 8, we can establish an upper bound on the regret by considering the difference between the optimistic and pessimistic value functions:

$$\text{Regret}_{\text{NASH}}(K) = \sum_{k=1}^K \max_{i \in \mathcal{M}} \left( V_{i,1}^{\dagger, \pi_{-i}^k, \rho_i} - V_{i,1}^{\pi^k, \rho_i} \right) (s_1^k) \leq \sum_{k=1}^K \max_{i \in \mathcal{M}} \left( \bar{V}_{i,1}^{k, \rho_i} - \underline{V}_{i,1}^{k, \rho_i} \right) (s_1^k). \quad (21)$$

For the TV-divergence uncertainty set, we begin by analyzing the difference between the upper and lower Q-values. Given our definitions for  $\bar{Q}_h^k, Q_{i,h}^{k, \rho_i}, \bar{V}_{i,h}^{k, \rho_i}$ , and  $\underline{V}_{i,h}^{k, \rho_i}$  (from eq. 5- 8), along with the bonus term  $\beta_{i,h}^k(s, \mathbf{a})$  defined in eq. 19, we can establish a bound on this difference for any  $(h, k) \in [H] \times [K]$  and  $(s, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$ :

$$\bar{Q}_h^k(s, \mathbf{a}) - Q_h^k(s, \mathbf{a}) \leq \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i}] - \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [\underline{V}_{i,h+1}^{k, \rho_i}] + 2\beta_{i,h}^k(s, \mathbf{a}). \quad (22)$$

We introduce two key terms,  $A$  and  $B$ , to simplify this expression:

$$A := \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i}] + \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [\underline{V}_{i,h+1}^{k, \rho_i}] - \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [\underline{V}_{i,h+1}^{k, \rho_i}]. \quad (23)$$

$$B := \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [\underline{V}_{i,h+1}^{k, \rho_i}]. \quad (24)$$

By substituting these definitions into eq. 22, we obtain a new bound:

$$\bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) - Q_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq A + B + 2\beta_{i,h}^k(s, \mathbf{a}). \quad (25)$$

We then proceed to bound each of these terms. A concentration bound argument tailored for TV robust expectations in Lemma 6 shows that  $A \leq 2\beta_{i,h}^k(s, \mathbf{a})$ . For term  $B$ , we use the dual representation of  $\sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}[V]$  from eq. 10 and Assumption 1 to first establish that  $B \leq \sup_{\eta \in [0, H]} \{ \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} [\eta - \bar{V}_{i,h+1}^{k, \rho_i}] + - \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} [\eta - \underline{V}_{i,h+1}^{k, \rho_i}] \}$ . Since  $\bar{V}_{i,h+1}^{k, \rho_i} \geq \underline{V}_{i,h+1}^{k, \rho_i}$  (by Lemma 8), we can simplify this further to  $B \leq \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i} - \underline{V}_{i,h+1}^{k, \rho_i}]$ .

By substituting the bounds for  $A$  and  $B$  back into eq. 25, we arrive at the following inequality:

$$\bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) - Q_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i} - \underline{V}_{i,h+1}^{k, \rho_i}] + 4\beta_{i,h}^k(s, \mathbf{a}). \quad (26)$$

Using Lemma 7 to upper bound the bonus term, and rearranging the terms, we obtain:

$$\begin{aligned} \bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) - Q_{i,h}^{k, \rho_i}(s, \mathbf{a}) &\leq \left( 1 + \frac{20}{H} \right) \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i} - \underline{V}_{i,h+1}^{k, \rho_i}] \\ &\quad + 4 \sqrt{\frac{c_1 \ell \text{Var}_{P_h^*(\cdot|s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i}]}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} \\ &\quad + \frac{4c_2 H^2 S \ell}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \sqrt{\frac{4}{K}}, \end{aligned} \quad (27)$$

where  $c_1, c_2 > 0$  are absolute constants. From the definitions in eq. 8, the difference in V-functions is given by:

$$\bar{V}_{i,h}^{k, \rho_i}(s) - \underline{V}_{i,h}^{k, \rho_i}(s) = \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} \left[ \bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) - Q_{i,h}^{k, \rho_i}(s, \mathbf{a}) \right]. \quad (28)$$

Now, let's define a new recursive value function  $\tilde{V}_h^{k, \rho_{\min}}$  and a corresponding Q-function  $\tilde{Q}_h^{k, \rho_{\min}}$  with  $\tilde{V}_{H+1}^{k, \rho_{\min}} = 0$ , where  $\rho_{\min} = \min_{i \in \mathcal{M}} \rho_i$ :

$$\begin{aligned} \tilde{Q}_h^{k, \rho_{\min}}(s, \mathbf{a}) = & \left(1 + \frac{20}{H}\right) \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} [\tilde{V}_{h+1}^{k, \rho_{\min}}] + 4\sqrt{\frac{c_1 \iota \text{Var}_{P_h^*(\cdot|s, \mathbf{a})} [\tilde{V}_{h+1}^{k, \rho_{\min}}]}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} \\ & + \frac{4c_2 H^2 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \sqrt{\frac{4}{K}}, \end{aligned} \quad (29)$$

$$\tilde{V}_h^{k, \rho_{\min}}(s) = \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [\tilde{Q}_h^{k, \rho_{\min}}(s, \mathbf{a})]. \quad (30)$$

It is a well-known property of robust value functions under TV-divergence that they become more conservative as the uncertainty radius  $\rho_i$  decreases (e.g., (Iyengar, 2005; Nilim & El Ghaoui, 2005)). Given that  $\rho_{\min} \leq \rho_i$  for all agents  $i \in \mathcal{M}$ , it follows that for every next state  $s' \in \mathcal{S}$ :

$$V_{i, h+1}^{k, \rho_i}(s') \leq V_{h+1}^{k, \rho_{\min}}(s') \quad \forall i \in \mathcal{M} \text{ and } s \in \mathcal{S}.$$

We can inductively prove that for any  $(i, h, s, \mathbf{a}) \in \mathcal{M} \times [H] \times \mathcal{S} \times \mathcal{A}$ :

$$\max_{i \in \mathcal{M}} \left( \tilde{Q}_{i, h}^{k, \rho_i}(s, \mathbf{a}) - Q_{i, h}^{k, \rho_i}(s, \mathbf{a}) \right) \leq \tilde{Q}_h^{k, \rho_{\min}}(s, \mathbf{a}), \quad (31)$$

$$\max_{i \in \mathcal{M}} \left( \tilde{V}_{i, h}^{k, \rho_i}(s) - V_{i, h}^{k, \rho_i}(s) \right) \leq \tilde{V}_h^{k, \rho_{\min}}(s). \quad (32)$$

Therefore, we only need to upper bound the sum  $\sum_{k=1}^K \tilde{V}_1^{k, \rho_{\min}}(s_1^k)$ . For simplicity, we define the following notations for the differences at any  $(h, k) \in [H] \times [K]$ :

$$\Delta_h^k := \tilde{V}_h^{k, \rho_{\min}}(s_h^k), \quad (33)$$

$$\zeta_h^k := \Delta_h^k - \tilde{Q}_h^{k, \rho_{\min}}(s_h^k, \mathbf{a}_h^k), \quad (34)$$

$$\xi_h^k := \mathbb{E}_{P_h^*(\cdot|s_h^k, \mathbf{a}_h^k)} [\tilde{V}_{h+1}^{k, \rho_{\min}}] - \Delta_{h+1}^k. \quad (35)$$

We can confirm that  $\{\zeta_h^k\}_{(h,k)}$  and  $\{\xi_h^k\}_{(h,k)}$  are martingale difference sequences with respect to their respective filtrations. By substituting eq. 29 into eq. 34, we get:

$$\begin{aligned} \Delta_h^k &= \zeta_h^k + \tilde{Q}_h^{k, \rho_{\min}}(s_h^k, \mathbf{a}_h^k) \\ &\leq \zeta_h^k + \left(1 + \frac{20}{H}\right) \mathbb{E}_{P_h^*(\cdot|s_h^k, \mathbf{a}_h^k)} [\tilde{V}_{h+1}^{k, \rho_{\min}}] + 4\sqrt{\frac{c_1 \iota \text{Var}_{P_h^*(\cdot|s_h^k, \mathbf{a}_h^k)} [\tilde{V}_{h+1}^{k, \rho_{\min}}]}{\{N_h^k(s_h^k, \mathbf{a}_h^k) \vee 1\}}} \\ &\quad + \frac{4c_2 H^2 S \iota}{\{N_h^k(s_h^k, \mathbf{a}_h^k) \vee 1\}} + \sqrt{\frac{4}{K}} \\ &= \zeta_h^k + \left(1 + \frac{20}{H}\right) \xi_h^k + \left(1 + \frac{20}{H}\right) \Delta_{h+1}^k + 4\sqrt{\frac{c_1 \iota \text{Var}_{P_h^*(\cdot|s, \mathbf{a})} [\tilde{V}_{h+1}^{k, \rho_{\min}}]}{\{N_h^k(s_h^k, \mathbf{a}_h^k) \vee 1\}}} \\ &\quad + \frac{4c_2 H^2 S \iota}{\{N_h^k(s_h^k, \mathbf{a}_h^k) \vee 1\}} + \sqrt{\frac{4}{K}}. \end{aligned} \quad (36)$$

By recursively applying eq. 36 and noting that  $\left(1 + \frac{20}{H}\right)^h \leq \left(1 + \frac{20}{H}\right)^H \leq c$  for some constant  $c \geq 0$ , we can upper bound the right-hand side of eq. 21 as:

$$\begin{aligned} \text{Regret}_{\text{NASH}}(K) &\leq \sum_{k=1}^K \Delta_1^k \leq c \sum_{k=1}^K \sum_{h=1}^H \left\{ (\zeta_h^k + \xi_h^k) \right. \\ &\quad + \left( 4\sqrt{\frac{c_1 \iota \text{Var}_{P_h^*(\cdot|s, \mathbf{a})} [\tilde{V}_{h+1}^{k, \rho_{\min}}]}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{4c_2 H^2 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} \right) \\ &\quad \left. + \sqrt{\frac{4}{K}} \right\}. \end{aligned} \quad (37)$$

The first term, a sum of martingale differences, is bounded using the Azuma-Hoeffding inequality from Lemma 26, yielding:

$$\sum_{k=1}^K \sum_{h=1}^H (\zeta_h^k + \xi_h^k) \leq c_1 \min \left\{ \frac{1}{\rho_{\min}}, H \right\} \sqrt{HK\iota}, \quad (38)$$

where  $c_1 > 0$  is an absolute constant. For the second term, we apply the Cauchy-Schwarz inequality to the summation of the variance terms:

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \sqrt{\frac{\text{Var}_{P_h^*}(\cdot | s_h^k, \mathbf{a}_h^k) [V_{h+1}^{\pi^k, \rho_{\min}}]}{N_h^k(s_h^k, \mathbf{a}_h^k) \vee 1}} &\leq \sqrt{\left( \sum_{k=1}^K \sum_{h=1}^H \text{Var}_{P_h^*}(\cdot | s_h^k, \mathbf{a}_h^k) [V_{h+1}^{\pi^k, \rho_{\min}}] \right)} \\ &\quad \sqrt{\left( \sum_{k=1}^K \sum_{h=1}^H \frac{1}{N_h^k(s_h^k, \mathbf{a}_h^k) \vee 1} \right)}. \end{aligned} \quad (39)$$

The second factor on the right-hand side is bounded by  $c_2 HS(\prod_{i=1}^m A_i)\iota$ , as shown in (Liu et al., 2021, Theorem 3), while the first factor is bounded using the Law of Total Variation and standard martingale concentration arguments (from (Jin et al., 2018) and (Lu et al., 2024)):

$$\sum_{k=1}^K \sum_{h=1}^H \text{Var}_{P_h^*}(\cdot | s_h^k, \mathbf{a}_h^k) [V_{h+1}^{\pi^k, \rho_{\min}}] \leq c_3 \cdot \left( \min \left\{ \frac{1}{\rho_{\min}}, H \right\} HK + \min \left\{ \frac{1}{\rho_{\min}}, H \right\}^3 H\iota \right). \quad (40)$$

By combining these bounds and substituting them into eq. 39, we can obtain a final bound for the second term. The third term,  $\sum_{k=1}^K \sum_{h=1}^H \sqrt{\frac{4}{K}}$ , is straightforwardly bounded by  $c_5 \sqrt{H^2 K}$ . By combining the bounds for all three terms, we arrive at the final regret bound for  $\text{Regret}_{\text{Nash}}(K)$ :

$$\text{Regret}_{\text{NASH}}(K) = \mathcal{O} \left( \sqrt{\min \left\{ \frac{1}{\rho_{\min}}, H \right\} H^2 SK \left( \prod_{i \in \mathcal{M}} A_i \right) \iota'} \right), \quad (41)$$

where  $\iota' = \log^2 \left( \frac{SHK \prod_{i \in \mathcal{M}} A_i}{\delta} \right)$ . This completes the proof of Theorem 1.  $\square$

**Remark 1.** The methodology for bounding the regret for Correlated Equilibrium (CE) and Coarse Correlated Equilibrium (CCE) settings mirrors the approach outlined here for the Nash equilibrium in the TV-DRMG context. The proofs leverage Lemma 9 and Lemma 10, respectively.

## E.2 KEY LEMMAS FOR TV-DRMG

**Lemma 5** (Gap between maximum and minimum (Lu et al., 2024)). Consider any RMG  $\mathcal{MG}_{\text{rob}} = \{\mathcal{S}, \mathcal{A}, H, \{\mathcal{P}_{TV}^{\rho_i}(P^*)\}_{i=1}^m, r\}$ . The robust value function  $V_{i,h}^{\pi, \rho_i}$  for all  $i \in \mathcal{M}$  and  $h \in [H]$  associated with any joint policy  $\pi$  satisfies

$$\forall (i, h) \in \mathcal{M} \times [H] : \max_{s \in \mathcal{S}} V_{i,h}^{\pi, \rho_i}(s) - \min_{s \in \mathcal{S}} V_{i,h}^{\pi, \rho_i}(s) \leq \nu_H^{\rho_i},$$

where  $\nu_H^{\rho_i} := \min \left\{ \frac{1}{\rho_i}, H - h + 1 \right\} \leq \min \left\{ \frac{1}{\rho_i}, H \right\}$ .

*Proof.* Refer to the proof-lines of Lemma 3 in (Shi et al., 2024b).  $\square$

**Lemma 6** (Bound of optimistic and pessimistic value estimators with bonus for TV-DRMG). Under the typical event  $\mathcal{E}_{TV}$  defined in eq. 20 and by setting the bonus  $\beta_{i,h}^k$  as in eq. 19, it holds that

$$\begin{aligned} \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [\overline{V}_{i,h+1}^{k, \rho_i}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [\overline{V}_{i,h+1}^{k, \rho_i}] \\ + \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [\underline{V}_{i,h+1}^{k, \rho_i}] - \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [\underline{V}_{i,h+1}^{k, \rho_i}] \leq 2\beta_{i,h}^k(s, \mathbf{a}). \end{aligned}$$

*Proof.* Let's denote the term to be bounded as  $A$ .

$$A := \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} \right] \\ + \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ \underline{V}_{i,h+1}^{k,\rho_i} \right] - \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ \underline{V}_{i,h+1}^{k,\rho_i} \right]. \quad (42)$$

Under the high-probability event  $\mathcal{E}_{TV}$  (as defined in eq. 20), we can apply the concentration inequality from Lemma 12 to upper bound  $A$  as follows:

$$A \leq 2\sqrt{\frac{c_1 \text{Var}_{\widehat{\mathcal{P}}_h^k} \left( \overline{V}_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right) \iota}{N_h^k(s, \mathbf{a}) \vee 1}} + \frac{2 \mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right]}{H} \\ + \frac{2c'_2 H^2 S \iota}{N_h^k(s, \mathbf{a}) \vee 1} + \frac{2}{\sqrt{K}}. \quad (43)$$

where  $\iota = \log(S^2(\prod_{i=1}^m A_i)H^2K^{3/2}/\delta)$  and  $c_1, c'_2 > 0$  are absolute constants. By applying the result from Lemma 14 to the variance term in eq. 43, we obtain the required bound presented in the lemma statement. This concludes the proof.  $\square$

**Lemma 7** (Bound of the bonus term for TV-DRMG). *Under the typical event  $\mathcal{E}_{TV}$ , the bonus term defined in 19 is bounded by*

$$\beta_{i,h}^k(s, \mathbf{a}) \leq \sqrt{\frac{c_1 \iota \text{Var}_{P_h^*(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{\pi^k, \rho_i} \right]}{N_h^k(s, \mathbf{a}) \vee 1}} + \frac{5 \mathbb{E}_{P_h^*(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right]}{H} \\ + \frac{c_2 H^2 S \iota}{N_h^k(s, \mathbf{a}) \vee 1} + \sqrt{\frac{1}{K}}.$$

where  $\iota = \log(S^3(\prod_{i=1}^m A_i)H^2K^{3/2}/\delta)$  and  $c_1, c_2 > 0$  are constants.

*Proof.* The proof-lines are similar to (Lu et al., 2024, Lemma E.4) or (Ghosh et al., 2025, Lemma K.3). Recall the bonus term defined in eq. 19. We need to bound the first and second term of eq. 19. We first bound the second term of  $\beta_{i,h}^k(s, \mathbf{a})$  by using Lemma 13, and we get

$$\frac{2 \mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right]}{H} \leq \left( \frac{2}{H} + \frac{2}{H^2} \right) \mathbb{E}_{P_h^*(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right] + \frac{c'_2 H S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} \\ \leq \frac{4 \mathbb{E}_{P_h^*(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right]}{H} + \frac{c'_2 H S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}, \quad (44)$$

where the second inequality is from  $H \geq 1$ . We now bound the first term (variance term) of eq. 19 by using Lemma 15, which gives

$$\sqrt{\frac{c_1 \iota \text{Var}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \frac{\overline{V}_{i,h+1}^{k,\rho_i} + \underline{V}_{i,h+1}^{k,\rho_i}}{2} \right]}{N_h^k(s, \mathbf{a}) \vee 1}} \leq \sqrt{\frac{c'_1 \iota \text{Var}_{P_h^*(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{\pi^k, \rho_i} \right]}{N_h^k(s, \mathbf{a}) \vee 1}} \\ + \frac{\mathbb{E}_{P_h^*(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right]}{H} \\ + \frac{c_3 H^2 S \iota}{N_h^k(s, \mathbf{a}) \vee 1}. \quad (45)$$

where  $c_3 > 0$  is an absolutely constant. Thus by combining eq. 44 and eq. 45 with the choice of bonus term in eq. 19, we can conclude the proof of Lemma 7.  $\square$

NE VERSION: OPTIMISTIC AND PESSIMISTIC ESTIMATION OF THE ROBUST VALUES FOR TV-DRMG.

Here we will proof the optimistic estimations are indeed upper bounds of the corresponding robust V-value and robust Q-value functions fro NE version.

**Lemma 8** (Optimistic and pessimistic estimation of the robust values for TV-DRMG for NE version). *By setting the bonus term  $\beta_{i,h}^k$  as in eq. 19, with probability  $1 - \delta$ , for any  $(s, \mathbf{a}, h, i)$  and  $k \in [K]$ , it holds that*

$$Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s, \mathbf{a}) \leq \overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}), \quad \underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi_{-i}^k, \rho_i}(s, \mathbf{a}), \quad (46)$$

$$V_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s) \leq \overline{V}_{i,h}^{k, \rho_i}(s), \quad \underline{V}_{i,h}^{k, \rho_i}(s) \leq V_{i,h}^{\pi_{-i}^k, \rho_i}(s). \quad (47)$$

*Proof.* The proof-lines are similar to (Ghosh et al., 2025) adapted to the multi-agent case. We will run a proof for each inequality outlined in Lemma 8.

- **Ineq. 1:** To prove  $Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s, \mathbf{a}) \leq \overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a})$ .
- **Ineq. 2:** To prove  $\underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi_{-i}^k, \rho_i}(s, \mathbf{a})$ .

We know that, at step  $h = H + 1$ ,  $\overline{V}_{i,H+1}^{k, \rho_i}(s) = V_{i,H+1}^{\dagger, \pi_{-i}^k, \rho_i}(s) = 0$ . Now, we assume that both eq. 46 and eq. 47 hold at the  $(h + 1)$ -th step.

- **Proof of Ineq. 1:** We first consider robust  $Q$  at the  $h$ -th step. Then, by Proposition 1 (Robust Bellman Equation) and eq. 5, we have that

$$\begin{aligned} \overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) - Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s, \mathbf{a}) &= \min \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \overline{V}_{i,h+1}^{k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right] \right. \\ &\quad \left. + \beta_{i,h}^k(s, \mathbf{a}), \nu_H^{\rho_i} - Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s, \mathbf{a}) \right\} \\ &\geq \min \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right] \right. \\ &\quad \left. + \beta_{i,h}^k(s, \mathbf{a}), 0 \right\}, \end{aligned} \quad (48)$$

where the second inequality follows from the induction of  $V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \leq \overline{V}_{i,h+1}^{k, \rho_i}$  at the  $h + 1$ -th step and the fact that  $Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i} \leq \nu_H^{\rho_i}$  by Lemma 5. By Lemma 11, we get

$$\begin{aligned} \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right] &\leq \sqrt{\frac{c_1 \text{Var}_{\widehat{\mathcal{P}}_h^k} \left( V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} \\ &\quad + \frac{c_2 H \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}. \end{aligned} \quad (49)$$

Now by further applying Lemma 14 to the variance term in the above inequality, we can obtain that

$$\begin{aligned}
& \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] \\
& \leq \sqrt{\frac{c_1 \left( \text{Var}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \left( \frac{\overline{V}_{i,h+1}^{k,\rho_i} + \underline{V}_{i,h+1}^{k,\rho_i}}{2} \right) \right] + 4H \mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right] \right) \iota}{\{N_h^k(s,\mathbf{a}) \vee 1\}}} \\
& \quad + \frac{c_2 H \iota}{\{N_h^k(s,\mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}} \\
& \stackrel{(i)}{\leq} \sqrt{\frac{c_1 \iota \text{Var}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \left( \frac{\overline{V}_{i,h+1}^{k,\rho_i} + \underline{V}_{i,h+1}^{k,\rho_i}}{2} \right) \right]}{\{N_h^k(s,\mathbf{a}) \vee 1\}}} + \sqrt{\frac{4H c_1 \iota \mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right]}{\{N_h^k(s,\mathbf{a}) \vee 1\}}} \\
& \quad + \frac{c_2 H \iota}{\{N_h^k(s,\mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}} \\
& \stackrel{(ii)}{\leq} \sqrt{\frac{c_1 \iota \text{Var}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \left( \frac{\overline{V}_{i,h+1}^{k,\rho_i} + \underline{V}_{i,h+1}^{k,\rho_i}}{2} \right) \right]}{\{N_h^k(s,\mathbf{a}) \vee 1\}}} + \frac{\mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right]}{H} \\
& \quad + \frac{H^2 c'_2 \iota}{\{N_h^k(s,\mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}, \tag{50}
\end{aligned}$$

where the inequality (i) is due to  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ , and the last inequality (ii) is from  $\sqrt{ab} \leq a+b$  where  $c'_2 > 0$  is an absolute constant. Therefore, combining eqns. 48, 49, 50, and the choice of bonus in 19, we can conclude that  $\overline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}) - Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s,\mathbf{a}) \geq 0$ .

• **Proof of Ineq. 2:** By Proposition 1 (Robust Bellman Equation) and eq. 6, we have that

$$\begin{aligned}
\overline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}) - Q_{i,h}^{\pi^k,\rho_i}(s,\mathbf{a}) &= \max \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ \underline{V}_{i,h+1}^{k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] \right. \\
&\quad \left. - \beta_{i,h}^k(s,\mathbf{a}), 0 - Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s,\mathbf{a}) \right\}, \\
&\leq \max \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] \right. \\
&\quad \left. - \beta_{i,h}^k(s,\mathbf{a}), 0 \right\}, \tag{51}
\end{aligned}$$

where the second inequality follows from the induction of  $V_{i,h+1}^{\pi^k,\rho_i} \geq \underline{V}_{i,h+1}^{k,\rho_i}$  at the  $h+1$ -th step and the fact that  $Q_{i,h}^{\pi^k,\rho_i} \geq 0$ . By Lemma 11, we can confirm that

$$\begin{aligned}
\sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] &\leq \sqrt{\frac{c_1 \text{Var}_{\widehat{\mathcal{P}}_h^k} \left( V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right) \cdot \iota}{\{N_h^k(s,\mathbf{a}) \vee 1\}}} \\
&\quad + \frac{\mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right]}{H} \\
&\quad + \frac{c'_2 H^2 S \iota}{\{N_h^k(s,\mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}. \tag{52}
\end{aligned}$$

Now by further applying Lemma 14 to the variance term in the above inequality, with an argument similar to eq. 49 we can obtain that

$$\begin{aligned} \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k, \rho_i} \right] &\leq \sqrt{\frac{c_1 \text{Var}_{\widehat{\mathcal{P}}_h^k} \left( V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} \\ &\quad + \frac{\mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} \left[ \overline{V}_{i,h+1}^{k, \rho_i} - V_{i,h+1}^{k, \rho_i} \right]}{H} \\ &\quad + \frac{c_2'' H^2 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}. \end{aligned} \quad (53)$$

where  $c_2'' > 0$  is an absolute constant. Therefore, combining eqns. 51, 52, 53, and the choice of bonus in 19,  $\underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) - Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}) \leq 0$ .

Therefore, by eq. 50 and eq. 53, we have proved that at step  $h$ , it holds that

$$Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s, \mathbf{a}) \leq \overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}), \quad \underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}). \quad (54)$$

We now assume that eq. 46 hold for  $h$ -th step. Then, by the definition of robust value function as given by robust Bellman equation (Proposition 1), and eq. 8, and NASH Equilibrium, we get

$$\overline{V}_{i,h}^{k, \rho_i}(s) = \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} \left[ \overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \right] = \max_{\pi_i'} \mathbb{E}_{\mathbf{a} \sim \pi_i' \times \pi_{-i}^k(\cdot|s)} \left[ \overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \right]. \quad (55)$$

By the definition of  $V_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s)$  in eq. 3, we get

$$V_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s) = \max_{\pi_i'} \mathbb{E}_{\mathbf{a} \sim \pi_i' \times \pi_{-i}^k(\cdot|s)} \left[ Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s, \mathbf{a}) \right]. \quad (56)$$

Since by induction, for any  $(s, \mathbf{a})$ ,  $\overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \geq Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s, \mathbf{a})$ . As a result, we also have  $\overline{V}_{i,h}^{k, \rho_i}(s) \geq V_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s)$ , which is eq. 47 for  $h$ -th step. Similarly, we can show that

$$\begin{aligned} \underline{V}_{i,h}^{k, \rho_i}(s) &= \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} \left[ \underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \right], \\ &\stackrel{(i)}{\leq} \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} \left[ Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}) \right], \\ &\stackrel{(ii)}{=} V_{i,h}^{\pi^k, \rho_i}(s), \end{aligned} \quad (57)$$

where (i) is due to the fact that  $\underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a})$  and (ii) is by definition of  $V_{i,h}^{\pi^k, \rho_i}(s)$  as given by Bellman equation in Proposition 1.  $\square$

#### CCE VERSION: OPTIMISTIC AND PESSIMISTIC ESTIMATION OF THE ROBUST VALUES FOR TV-DRMG.

Here we will proof the optimistic estimations are indeed upper bounds of the corresponding robust V-value and robust Q-value functions for CCE version.

**Lemma 9** (Optimistic and pessimistic estimation of the robust values for TV-DRMG for CCE version). *By setting the bonus term  $\beta_{i,h}^k$  as in eq. 19, with probability  $1 - \delta$ , for any  $(s, \mathbf{a}, h, i)$  and  $k \in [K]$ , it holds that*

$$\max_{\phi \in \Phi_i} Q_{i,h}^{\phi \diamond \pi^k, \rho_i}(s, \mathbf{a}) \leq \overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}), \quad \underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}), \quad (58)$$

$$\max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k, \rho_i}(s) \leq \overline{V}_{i,h}^{k, \rho_i}(s), \quad \underline{V}_{i,h}^{k, \rho_i}(s) \leq V_{i,h}^{\pi^k, \rho_i}(s). \quad (59)$$

*Proof.* The proof-lines are similar to (Ghosh et al., 2025) adapted to the multi-agent case. We will run a proof for each inequality outlined in Lemma 9.

- **Ineq. 1:** To prove  $Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})$ .
- **Ineq. 2:** To prove  $\underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi_{-i}^k,\rho_i}(s, \mathbf{a})$ .

We know that, at step  $h = H + 1$ ,  $\bar{V}_{i,H+1}^{k,\rho_i}(s) = V_{i,H+1}^{\dagger,\pi_{-i}^k,\rho_i}(s) = 0$ . Now, we assume that both eq. 58 and eq. 59 hold at the  $(h + 1)$ -th step.

- **Proof of Ineq. 1:** We first consider robust  $Q$  at the  $h$ -th step. Then, by Proposition 1 (Robust Bellman Equation) and eq. 5, we have that

$$\begin{aligned} \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) - Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) &= \min \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \bar{V}_{i,h+1}^{k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] \right. \\ &\quad \left. + \beta_{i,h}^k(s, \mathbf{a}), \nu_H^{\rho_i} - Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) \right\}, \\ &\geq \min \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] \right. \\ &\quad \left. + \beta_{i,h}^k(s, \mathbf{a}), 0 \right\}, \end{aligned} \quad (60)$$

where the second inequality follows from the induction of  $V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \leq \bar{V}_{i,h+1}^{k,\rho_i}$  at the  $h+1$ -th step and the fact that  $Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i} \leq \nu_H^{\rho_i}$  by Lemma 5. By Lemma 11, we get

$$\begin{aligned} \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] &\leq \sqrt{\frac{c_1 \text{Var}_{\widehat{\mathcal{P}}_h^k} \left( V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} \\ &\quad + \frac{c_2 H \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}. \end{aligned} \quad (61)$$

Now by further applying Lemma 14 to the variance term in the above inequality, we can obtain that

$$\begin{aligned} &\sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] \\ &\leq \sqrt{\frac{c_1 \left( \text{Var}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} \left[ \left( \frac{\bar{V}_{i,h+1}^{k,\rho_i} + V_{i,h+1}^{k,\rho_i}}{2} \right) \right] + 4H \mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} \left[ \bar{V}_{i,h+1}^{k,\rho_i} - V_{i,h+1}^{k,\rho_i} \right] \right) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} \\ &\quad + \frac{c_2 H \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}} \\ &\stackrel{(i)}{\leq} \sqrt{\frac{c_1 \iota \text{Var}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} \left[ \left( \frac{\bar{V}_{i,h+1}^{k,\rho_i} + V_{i,h+1}^{k,\rho_i}}{2} \right) \right]}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \sqrt{\frac{4H c_1 \iota \mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} \left[ \bar{V}_{i,h+1}^{k,\rho_i} - V_{i,h+1}^{k,\rho_i} \right]}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} \\ &\quad + \frac{c_2 H \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}} \\ &\stackrel{(ii)}{\leq} \sqrt{\frac{c_1 \iota \text{Var}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} \left[ \left( \frac{\bar{V}_{i,h+1}^{k,\rho_i} + V_{i,h+1}^{k,\rho_i}}{2} \right) \right]}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{\mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} \left[ \bar{V}_{i,h+1}^{k,\rho_i} - V_{i,h+1}^{k,\rho_i} \right]}{H} \\ &\quad + \frac{H^2 c_2' \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}, \end{aligned} \quad (62)$$



where the inequality (i) is due to  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ , and the last inequality (ii) is from  $\sqrt{ab} \leq a+b$  where  $c'_2 > 0$  is an absolute constant. Therefore, combining eqns. 60, 61, 62, and the choice of bonus in 19, we can conclude that  $\bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) - Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s, \mathbf{a}) \geq 0$ .

• **Proof of Ineq. 2:** By Proposition 1 (Robust Bellman Equation) and eq. 6, we have that

$$\begin{aligned} \underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) - Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}) &= \max \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{\pi^k, \rho_i} \right] \right. \\ &\quad \left. - \beta_{i,h}^k(s, \mathbf{a}), 0 - Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s, \mathbf{a}) \right\}, \\ &\leq \max \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{\pi^k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{\pi^k, \rho_i} \right] \right. \\ &\quad \left. - \beta_{i,h}^k(s, \mathbf{a}), 0 \right\}, \end{aligned} \quad (63)$$

where the second inequality follows from the induction of  $\underline{V}_{i,h+1}^{\pi^k, \rho_i} \geq \underline{V}_{i,h+1}^{k,\rho_i}$  at the  $h+1$ -th step and the fact that  $Q_{i,h}^{\pi^k, \rho_i} \geq 0$ . By Lemma 11, we can confirm that

$$\begin{aligned} \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{\pi^k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{\pi^k, \rho_i} \right] &\leq \sqrt{\frac{c_1 \text{Var}_{\widehat{\mathcal{P}}_h^k} \left( \underline{V}_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} \\ &\quad + \frac{\mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{\pi^k, \rho_i} \right]}{H} \\ &\quad + \frac{c'_2 H^2 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}. \end{aligned} \quad (64)$$

Now by further applying Lemma 14 to the variance term in the above inequality, with an argument similar to eq. 61 we can obtain that

$$\begin{aligned} \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{\pi^k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{\pi^k, \rho_i} \right] &\leq \sqrt{\frac{c_1 \text{Var}_{\widehat{\mathcal{P}}_h^k} \left( \underline{V}_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} \\ &\quad + \frac{\mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{\pi^k, \rho_i} \right]}{H} \\ &\quad + \frac{c''_2 H^2 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}. \end{aligned} \quad (65)$$

where  $c''_2 > 0$  is an absolute constant. Therefore, combining eqns. 63, 64, 65, and the choice of bonus in 19,  $\underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) - Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}) \leq 0$ .

Therefore, by eq. 62 and eq. 65, we have proved that at step  $h$ , it holds that

$$Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}), \quad \underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}). \quad (66)$$

We now assume that eq. 58 hold for  $h$ -th step. Then, by the definition of robust value function as given by robust Bellman equation (Proposition 1), eq. 8, and CCE Equilibrium, we get

$$\bar{V}_{i,h}^{k,\rho_i}(s) = \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} \left[ \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \right] \geq \max_{\pi'_i} \mathbb{E}_{\mathbf{a} \sim \pi'_i \times \pi_{-i}^k(\cdot|s)} \left[ \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \right], \quad (67)$$

By the definition of  $V_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s)$  in eq. 3, we get

$$V_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s) = \max_{\pi'_i} \mathbb{E}_{\mathbf{a} \sim \pi'_i \times \pi_{-i}^k(\cdot|s)} \left[ Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s, \mathbf{a}) \right]. \quad (68)$$

Since by induction, for any  $(s, \mathbf{a})$ ,  $\bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \geq Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a})$ . As a result, we also have  $\bar{V}_{i,h}^{k,\rho_i}(s) \geq V_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s)$ , which is eq. 59 for  $h$ -th step. Similarly, we can show that

$$\begin{aligned} \underline{V}_{i,h}^{k,\rho_i}(s) &= \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} \left[ \underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \right], \\ &\stackrel{(i)}{\leq} \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} \left[ Q_{i,h}^{\pi^k,\rho_i}(s, \mathbf{a}) \right], \\ &\stackrel{(ii)}{=} V_{i,h}^{\pi^k,\rho_i}(s), \end{aligned} \quad (69)$$

where (i) is due to the fact that  $Q_{i,h}^{k,\rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k,\rho_i}(s, \mathbf{a})$  and (ii) is by definition of  $V_{i,h}^{\pi^k,\rho_i}(s)$  as given by Bellman equation in Proposition 1.  $\square$

#### CE VERSION: OPTIMISTIC AND PESSIMISTIC ESTIMATION OF THE ROBUST VALUES FOR TV-DRMG.

Here we will proof the optimistic estimations are indeed upper bounds of the corresponding robust V-value and robust Q-value functions for CE version.

**Lemma 10** (Optimistic and pessimistic estimation of the robust values for TV-DRMG for CE version). *By setting the bonus term  $\beta_{i,h}^k$  as in eq. 19, with probability  $1 - \delta$ , for any  $(s, \mathbf{a}, h, i)$  and  $k \in [K]$ , it holds that*

$$Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}), \quad \underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k,\rho_i}(s, \mathbf{a}), \quad (70)$$

$$V_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s) \leq \bar{V}_{i,h}^{k,\rho_i}(s), \quad \underline{V}_{i,h}^{k,\rho_i}(s) \leq V_{i,h}^{\pi^k,\rho_i}(s). \quad (71)$$

*Proof.* The proof-lines are similar to (Ghosh et al., 2025) adapted to the multi-agent case. We will run a proof for each inequality outlined in Lemma 10.

- **Ineq. 1:** To prove  $Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})$ .
- **Ineq. 2:** To prove  $\underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k,\rho_i}(s, \mathbf{a})$ .

We know that, at step  $h = H + 1$ ,  $\bar{V}_{i,H+1}^{k,\rho_i}(s) = V_{i,H+1}^{\dagger,\pi_{-i}^k,\rho_i}(s) = 0$ . Now, we assume that both eq. 70 and eq. 71 hold at the  $(h + 1)$ -th step.

- **Proof of Ineq. 1:** We first consider robust  $Q$  at the  $h$ -th step. Then, by Proposition 1 (Robust Bellman Equation) and eq. 5, we have that

$$\begin{aligned} \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) - Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) &= \min \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \bar{V}_{i,h+1}^{k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] \right. \\ &\quad \left. + \beta_{i,h}^k(s, \mathbf{a}), \nu_H^{\rho_i} - Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) \right\}, \\ &\geq \min \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] \right. \\ &\quad \left. + \beta_{i,h}^k(s, \mathbf{a}), 0 \right\}. \end{aligned} \quad (72)$$

where the second inequality follows from the induction of  $V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \leq \bar{V}_{i,h+1}^{k,\rho_i}$  at the  $h+1$ -th step and the fact that  $Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i} \leq \nu_H^{\rho_i}$  by Lemma 5. By Lemma 11, we get

$$\begin{aligned} \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] &\leq \sqrt{\frac{c_1 \text{Var}_{\widehat{\mathcal{P}}_h^k} \left( V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} \\ &\quad + \frac{c_2 H \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}. \end{aligned} \quad (73)$$

Now by further applying Lemma 14 to the variance term in the above inequality, we can obtain that

$$\begin{aligned}
& \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] \\
& \leq \sqrt{\frac{c_1 \left( \text{Var}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \left( \frac{\overline{V}_{i,h+1}^{k,\rho_i} + \underline{V}_{i,h+1}^{k,\rho_i}}{2} \right) \right] + 4H \mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right] \right) \iota}{\{N_h^k(s,\mathbf{a}) \vee 1\}}} \\
& + \frac{c_2 H \iota}{\{N_h^k(s,\mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}} \\
& \stackrel{(i)}{\leq} \sqrt{\frac{c_1 \iota \text{Var}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \left( \frac{\overline{V}_{i,h+1}^{k,\rho_i} + \underline{V}_{i,h+1}^{k,\rho_i}}{2} \right) \right]}{\{N_h^k(s,\mathbf{a}) \vee 1\}}} + \sqrt{\frac{4H c_1 \iota \mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right]}{\{N_h^k(s,\mathbf{a}) \vee 1\}}} \\
& + \frac{c_2 H \iota}{\{N_h^k(s,\mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}} \\
& \stackrel{(ii)}{\leq} \sqrt{\frac{c_1 \iota \text{Var}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \left( \frac{\overline{V}_{i,h+1}^{k,\rho_i} + \underline{V}_{i,h+1}^{k,\rho_i}}{2} \right) \right]}{\{N_h^k(s,\mathbf{a}) \vee 1\}}} + \frac{\mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right]}{H} \\
& + \frac{H^2 c'_2 \iota}{\{N_h^k(s,\mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}, \tag{74}
\end{aligned}$$

where the inequality (i) is due to  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ , and the last inequality (ii) is from  $\sqrt{ab} \leq a+b$  where  $c'_2 > 0$  is an absolute constant. Therefore, combining eqns. 72, 73, 74, and the choice of bonus in 19, we can conclude that  $\underline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}) - Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s,\mathbf{a}) \geq 0$ .

• **Proof of Ineq. 2:** By Proposition 1 (Robust Bellman Equation) and eq. 6, we have that

$$\begin{aligned}
\underline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}) - Q_{i,h}^{\pi^k,\rho_i}(s,\mathbf{a}) &= \max \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ \underline{V}_{i,h+1}^{k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] \right. \\
&\quad \left. - \beta_{i,h}^k(s,\mathbf{a}), 0 - Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s,\mathbf{a}) \right\}, \\
&\leq \max \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] \right. \\
&\quad \left. - \beta_{i,h}^k(s,\mathbf{a}), 0 \right\}, \tag{75}
\end{aligned}$$

where the second inequality follows from the induction of  $V_{i,h+1}^{\pi^k,\rho_i} \geq \underline{V}_{i,h+1}^{k,\rho_i}$  at the  $h+1$ -th step and the fact that  $Q_{i,h}^{\pi^k,\rho_i} \geq 0$ . By Lemma 11, we can confirm that

$$\begin{aligned}
\sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] &\leq \sqrt{\frac{c_1 \text{Var}_{\widehat{\mathcal{P}}_h^k} \left( V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right) \cdot \iota}{\{N_h^k(s,\mathbf{a}) \vee 1\}}} \\
&\quad + \frac{\mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s,\mathbf{a})} \left[ \overline{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i} \right]}{H} \\
&\quad + \frac{c'_2 H^2 S \iota}{\{N_h^k(s,\mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}. \tag{76}
\end{aligned}$$

Now by further applying Lemma 14 to the variance term in the above inequality, with an argument similar to eq. 73 we can obtain that

$$\begin{aligned} \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\pi^k, \rho_i}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\pi^k, \rho_i}] &\leq \sqrt{\frac{c_1 \text{Var}_{\widehat{\mathcal{P}}_h^k} (V_{i,h+1}^{\pi^k, \rho_i}) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} \\ &\quad + \frac{\mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} [\overline{V}_{i,h+1}^{k, \rho_i} - V_{i,h+1}^{k, \rho_i}]}{H} \\ &\quad + \frac{c_2'' H^2 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}, \end{aligned} \quad (77)$$

where  $c_2'' > 0$  is an absolute constant. Therefore, combining eqns. 75, 76, 77, and the choice of bonus in 19,  $\underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) - Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}) \leq 0$ .

Therefore, by eq. 74 and eq. 77, we have proved that at step  $h$ , it holds that

$$Q_{i,h}^{\dagger, \pi^k, \rho_i}(s, \mathbf{a}) \leq \overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}), \quad \underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}). \quad (78)$$

We now assume that eq. 70 hold for  $h$ -th step. Then, by the definition of robust value function as given by robust Bellman equation (Proposition 1), eq. 8, and CE Equilibrium, we get

$$\overline{V}_{i,h}^{k, \rho_i}(s) = \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [\overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a})] = \max_{\phi \in \Phi_i} \mathbb{E}_{\mathbf{a} \sim \phi \diamond \pi^k(\cdot|s)} [\overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a})]. \quad (79)$$

By the definition of  $\max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k, \rho_i}(s)$  in eq. 3, we get

$$\max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k, \rho_i}(s) = \max_{\phi \in \Phi_i} \mathbb{E}_{\mathbf{a} \sim \phi \diamond \pi^k(\cdot|s)} \left[ \max_{\phi'} Q_{i,h}^{\phi' \diamond \pi^k, \rho_i}(s, \mathbf{a}) \right]. \quad (80)$$

Since by induction, for any  $(s, \mathbf{a})$ ,  $\overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \geq \max_{\phi \in \Phi_i} Q_{i,h}^{\phi \diamond \pi^k, \rho_i}(s, \mathbf{a})$ . As a result, we also have

$\overline{V}_{i,h}^{k, \rho_i}(s) \geq \max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k, \rho_i}(s)$ , which is eq. 161 for  $h$ -th step. Similarly, we can show that

$$\begin{aligned} \underline{V}_{i,h}^{k, \rho_i}(s) &= \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [\underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a})], \\ &\stackrel{(i)}{\leq} \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a})], \\ &\stackrel{(ii)}{=} V_{i,h}^{\pi^k, \rho_i}(s), \end{aligned} \quad (81)$$

where (i) is due to the fact that  $\underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a})$  and (ii) is by definition of  $V_{i,h}^{\pi^k, \rho_i}(s)$  as given by Bellman equation in Proposition 1.  $\square$

### E.3 AUXILIARY LEMMAS FOR TV-DRMG

**Lemma 11** (Bernstein bound for TV-DRMG and the robust value functions of  $\pi^k$  and  $\pi^\dagger$ ). *Under event  $\mathcal{E}_{TV}$  in eq. 20 and definition of  $\pi^\dagger$  as given in eq. 18, we assume that for any  $\text{EQUILIBRIUM} \in \{\text{NASH}, \text{CE}, \text{CCE}\}$  the optimism and pessimism inequalities holds at  $(h+1, k)$ , where these inequalities can correspond to any of the following cases of EQUILIBRIUM:*

- **NE:** Lemma 8 using eq. 46 and eq. 47,
- **CCE:** Lemma 9 using eq. 58 and eq. 59,
- **CE:** Lemma 10 using eq. 70 and eq. 71,

Then, it holds that

$$\left| \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}}(s, \mathbf{a})[V_{i,h+1}^{\pi^k, \rho_i}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a})[V_{i,h+1}^{\pi^k, \rho_i}] \right| \leq \begin{cases} \sqrt{\frac{c_1 \text{Var}_{\widehat{P}_h^k}(V_{i,h+1}^{\pi^k, \rho_i}) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{c_2 H \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}, & \text{if } \pi^k = \pi^\dagger \\ \sqrt{\frac{c_1 \text{Var}_{\widehat{P}_h^k}(V_{i,h+1}^{\pi^k, \rho_i}) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{\mathbb{E}_{\widehat{P}_h^k(\cdot | s, \mathbf{a})}[\bar{V}_{i,h+1}^{\pi^k, \rho_i} - V_{i,h+1}^{\pi^k, \rho_i}]}{H} + \frac{c'_2 H^2 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}, & \text{otherwise,} \end{cases}$$

where  $\iota = \log \left( \frac{S^2 (\prod_{i=1}^m A_i) H^2 K^{3/2}}{\delta} \right)$  and  $c_1, c'_2 > 0$  are absolute constants.

*Proof.* By our definition of the operator  $\sigma_{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a})[V_{i,h+1}^{\pi^k, \rho_i}]$  in eq. 10, we can arrive at,

$$\begin{aligned} \left| \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}}(s, \mathbf{a})[V_{i,h+1}^{\pi^k, \rho_i}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a})[V_{i,h+1}^{\pi^k, \rho_i}] \right| &\leq \sup_{\eta \in [0, H]} \left| \left\{ \mathbb{E}_{\widehat{P}_h^k(\cdot | s, \mathbf{a})}[(\eta - V_{i,h+1}^{\pi^k, \rho_i})_+] \right. \right. \\ &\quad \left. \left. - \mathbb{E}_{P_h^*(\cdot | s, \mathbf{a})}[(\eta - V_{i,h+1}^{\pi^k, \rho_i})_+] \right\} \right| \\ &= \text{Term (i)} + \text{Term (ii)}. \end{aligned} \quad (82)$$

where we denote

$$\begin{aligned} \text{Term (i)} &:= \sup_{\eta \in [0, H]} \left| \left\{ \mathbb{E}_{\widehat{P}_h^k(\cdot | s, \mathbf{a})}[(\eta - V_{i,h+1}^{\pi^k, \rho_i})_+] \right. \right. \\ &\quad \left. \left. - \mathbb{E}_{P_h^*(\cdot | s, \mathbf{a})}[(\eta - V_{i,h+1}^{\pi^k, \rho_i})_+] \right\} \right| \end{aligned} \quad (83)$$

$$\begin{aligned} \text{Term (ii)} &:= \sup_{\eta \in [0, H]} \left| \left\{ \mathbb{E}_{\widehat{P}_h^k(\cdot | s, \mathbf{a})} \left[ \left( \eta - V_{i,h+1}^{\pi^k, \rho_i} \right)_+ - \left( \eta - V_{i,h+1}^{\pi_{-i}^k, \rho_i} \right)_+ \right] \right. \right. \\ &\quad \left. \left. - \mathbb{E}_{P_h^*(\cdot | s, \mathbf{a})} \left[ \left( \eta - V_{i,h+1}^{\pi^k, \rho_i} \right)_+ - \left( \eta - V_{i,h+1}^{\pi_{-i}^k, \rho_i} \right)_+ \right] \right\} \right|. \end{aligned} \quad (84)$$

We deal with Term (i) and Term (ii) respectively.

**Bound for Term (i):** Term (i) is referred to Bernstein bound for Bernstein bound for TV-DRMG and the robust value function of the robust best response  $\pi_i^{\dagger, \rho_i}(\pi_{-i})$ . More specifically, we find the Bernstein bound on the gap  $\left| \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}}(s, \mathbf{a})[V_{i,h+1}^{\pi_i^{\dagger, \rho_i}}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a})[V_{i,h+1}^{\pi_i^{\dagger, \rho_i}}] \right|$ . Therefore, by the definition of the operator  $\sigma_{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a})[V_{i,h+1}^{\pi_i^{\dagger, \rho_i}}]$  in eq. 10, we can arrive at,

$$\begin{aligned} &\left| \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}}(s, \mathbf{a})[V_{i,h+1}^{\pi_i^{\dagger, \rho_i}}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a})[V_{i,h+1}^{\pi_i^{\dagger, \rho_i}}] \right| \\ &\leq \sup_{\eta \in [0, H]} \left| \left\{ \mathbb{E}_{\widehat{P}_h^k(\cdot | s, \mathbf{a})} \left[ \left( \eta - V_{i,h+1}^{\pi_i^{\dagger, \rho_i}} \right)_+ \right] - \mathbb{E}_{P_h^*(\cdot | s, \mathbf{a})} \left[ \left( \eta - V_{i,h+1}^{\pi_i^{\dagger, \rho_i}} \right)_+ \right] \right\} \right| \\ &= \text{Term (i)}. \end{aligned} \quad (85)$$

By now according to the first inequality of event  $\mathcal{E}$  in eq. 20, we can bound eq. 85 as

$$\begin{aligned} \text{Term (i)} &\leq \sqrt{\frac{c_1 \text{Var}_{\hat{P}_h^k} \left( \eta - V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right)_+ \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{c_2 H \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} \\ &\leq \sqrt{\frac{c_1 \text{Var}_{\hat{P}_h^k} \left( V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{c_2 H \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}, \end{aligned} \quad (86)$$

for any  $\eta \in \mathcal{N}_{1/(S\sqrt{K})}([0, H])$ . Here the second inequality is because  $\text{Var}[(a - X)_+] \leq \text{Var}[X]$ . Therefore, by applying the covering argument in eq. 86, for any  $\eta \in [0, H]$ , it holds that

$$\text{Term (i)} \leq \sqrt{\frac{c_1 \text{Var}_{\hat{P}_h^k} \left( V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{c_2 H \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}. \quad (87)$$

**Bound for Term (ii):** For Term (ii), we apply the second inequality of event  $\mathcal{E}$  in eq. 20, and we obtain that

$$\begin{aligned} \text{Term (ii)} &\leq \sup_{\eta \in [0, H]} \left\{ \sum_{s' \in \mathcal{S}} \left( \sqrt{\frac{c_1 \min \{P_h^*(s' | s, \mathbf{a}), P_h^k(s' | s, \mathbf{a})\} \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{c_2 \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} \right) \right. \\ &\quad \left. \times \left| \left( \eta - V_{i,h+1}^{\pi_{-i}^k, \rho_i} \right)_+ - \left( \eta - V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right)_+ \right| \right\}. \end{aligned} \quad (88)$$

Now by assuming that eq. 47 holds at  $(h+1, k)$ , we can upper bound the absolute value above by

$$\begin{aligned} \left| \left( \eta - V_{i,h+1}^{\pi_{-i}^k, \rho_i} \right)_+ - \left( \eta - V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right)_+ \right| &\stackrel{(i)}{\leq} \left| V_{i,h+1}^{\pi_{-i}^k, \rho_i} - V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right| \\ &\stackrel{(ii)}{\leq} \bar{V}_{i,h+1}^{k, \rho_i}(s') - \underline{V}_{i,h+1}^{k, \rho_i}(s'), \end{aligned} \quad (89)$$

where the first inequality (i) is due to the 1-Lipschitz continuity of  $\psi_\eta(x) = (\eta - x)_+$ , and the second inequality (ii) is due to eq. 47. Thus combining eq. 88 and eq. 89, we get

$$\begin{aligned} \text{Term (ii)} &\leq \sum_{s' \in \mathcal{S}} \left( \sqrt{\frac{c_1 \hat{P}_h^k(s' | s, \mathbf{a}) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{c_2 \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} \right) \cdot \left( \bar{V}_{i,h+1}^{k, \rho_i}(s') - \underline{V}_{i,h+1}^{k, \rho_i}(s') \right) \\ &\stackrel{(i)}{\leq} \sum_{s' \in \mathcal{S}} \left( \frac{\hat{P}_h^k(s' | s, \mathbf{a})}{H} + \frac{c_1 H \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{c_2 \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} \right) \\ &\quad \cdot \left( \bar{V}_{i,h+1}^{k, \rho_i}(s') - \underline{V}_{i,h+1}^{k, \rho_i}(s') \right) \\ &\stackrel{(ii)}{\leq} \frac{\mathbb{E}_{\hat{P}_h^k(\cdot | s, \mathbf{a})} \left[ \bar{V}_{i,h+1}^{k, \rho_i} - \underline{V}_{i,h+1}^{k, \rho_i} \right]}{H} + \frac{c_2' H^2 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}, \end{aligned} \quad (90)$$

where  $c_2' > 0$  is an absolute constant. The first inequality (i) is by  $\sqrt{ab} \leq a + b$  and the second inequality (ii) is due to  $\bar{V}_{i,h+1}^{k, \rho_i}, \underline{V}_{i,h+1}^{k, \rho_i} \in [0, H]$ . Finally, by combining eq. 87 and eq. 90 and applying in eq. 82, we get the required bound as

$$\begin{aligned} \text{Term (ii)} &\leq \sqrt{\frac{c_1 \text{Var}_{\hat{P}_h^k} \left( V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{\mathbb{E}_{\hat{P}_h^k(\cdot | s, \mathbf{a})} \left[ \bar{V}_{i,h+1}^{k, \rho_i} - \underline{V}_{i,h+1}^{k, \rho_i} \right]}{H} + \frac{c_2' H^2 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} \\ &\quad + \frac{1}{\sqrt{K}}. \end{aligned} \quad (91)$$

This concludes the proof of Lemma 11.  $\square$

**Lemma 12** (Bernstein bound for TV-DRMG and optimistic and pessimistic robust value estimators). *Under event  $\mathcal{E}_{TV}$  in eq. 20 and definition of  $\pi^\dagger$  as given in eq. 18, we assume that for any  $\text{EQUILIBRIUM} \in \{\text{NASH}, \text{CE}, \text{CCE}\}$  the optimism and pessimism inequalities holds at  $(h+1, k)$ , where these inequalities can correspond to any of the following cases of  $\text{EQUILIBRIUM}$ :*

- **NE:** Lemma 8 using eq. 46 and eq. 47,
- **CCE:** Lemma 9 using eq. 58 and eq. 59,
- **CE:** Lemma 10 using eq. 70 and eq. 71,

Then, it holds that

$$\begin{aligned} & \max \left\{ \left| \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i}] \right|, \left| \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{k, \rho_i}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{k, \rho_i}] \right| \right\} \\ & \leq \sqrt{\frac{c_1 \text{Var}_{\widehat{\mathcal{P}}_h^k} (V_{i,h+1}^{\dagger, \pi^k_{i, \rho_i}}) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{\mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i} - V_{i,h+1}^{k, \rho_i}]}{H} + \frac{c'_2 H^2 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + \frac{1}{\sqrt{K}}, \end{aligned}$$

where  $\iota = \log \left( \frac{S^2 (\prod_{i=1}^m A_i) H^2 K^{3/2}}{\delta} \right)$  and  $c_1, c'_2 > 0$  are absolute constants.

*Proof.* This follows from the same proof as Lemma 11 and is thus omitted.  $\square$

**Lemma 13** (Non-robust Concentration for TV-DRMG). *Under event  $\mathcal{E}_{TV}$  in eq. 20 and definition of  $\pi^\dagger$  as given in eq. 18, we assume that for any  $\text{EQUILIBRIUM} \in \{\text{NASH}, \text{CE}, \text{CCE}\}$  the optimism and pessimism inequalities holds at  $(h+1, k)$ , where these inequalities can correspond to any of the following cases of  $\text{EQUILIBRIUM}$ :*

- **NE:** Lemma 8 using eq. 46 and eq. 47,
- **CCE:** Lemma 9 using eq. 58 and eq. 59,
- **CE:** Lemma 10 using eq. 70 and eq. 71,

Then, it holds that

$$\begin{aligned} \left| \mathbb{E}_{P_h^\star(\cdot|s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i} - V_{i,h+1}^{k, \rho_i}] - \mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i} - V_{i,h+1}^{k, \rho_i}] \right| & \leq \frac{\mathbb{E}_{\widehat{\mathcal{P}}_h^k(\cdot|s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i} - V_{i,h+1}^{k, \rho_i}]}{H} \\ & + \frac{c'_2 H^2 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}, \end{aligned}$$

where  $\iota = \log \left( \frac{S^2 (\prod_{i=1}^m A_i) H^2 K^{3/2}}{\delta} \right)$  and  $c'_2 > 0$  are absolute constants.

*Proof.* Assuming that eq. 47 holds for  $(h+1, k)$ , we apply the second inequality of event  $\mathcal{E}$  in eq. 20 to get the required bound Lemma 13.  $\square$

**Lemma 14** (Variance analysis for  $\pi^\dagger$  for TV-DRMG). *Under the definition of  $\pi^\dagger$  as given in eq. 18, we assume that for any  $\text{EQUILIBRIUM} \in \{\text{NASH}, \text{CE}, \text{CCE}\}$  the optimism and pessimism inequalities holds at  $(h+1, k)$ , where these inequalities can correspond to any of the following cases of  $\text{EQUILIBRIUM}$ :*

- **NE:** Lemma 8 using eq. 46 and eq. 47,
- **CCE:** Lemma 9 using eq. 58 and eq. 59,
- **CE:** Lemma 10 using eq. 70 and eq. 71,

Then, it holds that

$$\left| \text{Var}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \frac{\bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i}}{2} \right] - \text{Var}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right] \right| \leq 4H \mathbb{E}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \bar{V}_{h+1}^{k, \rho_i} - \underline{V}_{h+1}^{k, \rho_i} \right].$$

*Proof.* Our proof closely follows the lines of Lemma 22 in (Liu et al., 2021) and Lemma E.11 in (Lu et al., 2024), with detailed elaboration on each step for clarity. The left hand side of the inequality in Lemma 14 can be upper bounded by the following

$$\begin{aligned} & \left| \text{Var}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \left( \frac{\bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i}}{2} \right) \right] - \text{Var}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right] \right| \\ & \leq \left| \mathbb{E}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \left( \frac{\bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i}}{2} \right)^2 \right] - \mathbb{E}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \left( V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right)^2 \right] \right| \\ & \quad + \left| \left( \mathbb{E}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \left( \frac{\bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i}}{2} \right) \right] \right)^2 - \left( \mathbb{E}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right] \right)^2 \right|. \end{aligned} \quad (92)$$

By applying eq. 47 and the facts that  $\bar{V}_{i,h+1}^{k, \rho_i}$  and  $\underline{V}_{i,h+1}^{k, \rho_i}$ ,  $\bar{V}_{i,h+1}^{k, \rho_i}$ ,  $\underline{V}_{i,h+1}^{k, \rho_i}$ ,  $V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \in [0, H]$ , we can further upper bound eq. 92 as

$$\begin{aligned} & \left| \text{Var}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \left( \frac{\bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i}}{2} \right) \right] - \text{Var}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right] \right| \\ & \leq 4H \mathbb{E}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \left| \frac{\bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i}}{2} - V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i} \right| \right] \leq 4H \mathbb{E}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \bar{V}_{i,h+1}^{k, \rho_i} - \underline{V}_{i,h+1}^{k, \rho_i} \right]. \end{aligned} \quad (93)$$

This concludes the proof of Lemma 14.  $\square$

**Lemma 15** (Variance analysis for any robust joint policy  $\pi^k$  for TV-DRMG). *Under event  $\mathcal{E}_{TV}$  in eq. 20 and definition of  $\pi^\dagger$  as given in eq. 18, we assume that for any EQUILIBRIUM  $\in \{\text{NASH}, \text{CE}, \text{CCE}\}$  the optimism and pessimism inequalities holds at  $(h+1, k)$ , where these inequalities can correspond to any of the following cases of EQUILIBRIUM:*

- **NE:** Lemma 8 using eq. 46 and eq. 47,
- **CCE:** Lemma 9 using eq. 58 and eq. 59,
- **CE:** Lemma 10 using eq. 70 and eq. 71,

Then, then the following inequality holds,

$$\begin{aligned} & \left| \text{Var}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \left( \frac{\bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i}}{2} \right) \right] - \text{Var}_{P_h^*(\cdot|s, \mathbf{a})} \left[ V_{i,h+1}^{\pi^k, \rho_i} \right] \right| \\ & \leq 4H \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \bar{V}_{h+1}^{k, \rho_i} - \underline{V}_{h+1}^{k, \rho_i} \right] + \frac{c_2' H^4 S_t}{\{N_h^k(s, \mathbf{a}) \vee 1\}} + 1. \end{aligned}$$

*Proof.* We follow the proof-lines of Lemma 23 in (Liu et al., 2021) and Lemma E.12 of (Lu et al., 2024). We present a detailed derivation as follows. We first relate the variance on  $\hat{P}_h^k$  to the variance on  $P_h^*$ . Specifically, we have

$$\left| \text{Var}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \left( \frac{\bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i}}{2} \right) \right] - \text{Var}_{P_h^*(\cdot|s, \mathbf{a})} \left[ V_{i,h+1}^{\pi^k, \rho_i} \right] \right| \leq \text{Term (i)} + \text{Term (ii)}, \quad (94)$$



where we denote

$$\text{Term (i)} := \left| \text{Var}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \frac{\bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i}}{2} \right] - \text{Var}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \frac{\bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i}}{2} \right] \right|. \quad (95)$$

$$\text{Term (ii)} := \left| \text{Var}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \left( \frac{\bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i}}{2} \right) \right] - \text{Var}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ V_{i,h+1}^{\pi^k, \rho_i} \right] \right|. \quad (96)$$

We will now bound Term (i) and Term (ii) respectively.

- **Term (i):** By applying the fact  $\left( \bar{V}_{i,h+1}^{k, \rho_i} + \underline{V}_{i,h+1}^{k, \rho_i} \right) / 2 \in [0, H]$  in the variance terms on Term (i), we can upper bound Term (i) as

$$\begin{aligned} \text{Term (i)} &\leq H^2 \sum_{s' \in \mathcal{S}} \left| P_h^*(s'|s, \mathbf{a}) - \hat{P}_h^k(s'|s, \mathbf{a}) \right| \\ &\stackrel{(i)}{\leq} H^2 \sum_{s' \in \mathcal{S}} \left( \sqrt{\frac{c_1 \hat{P}_h^k(s'|s, \mathbf{a}) \cdot \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{c_2 \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} \right) \\ &\stackrel{(ii)}{\leq} H^2 \left( \sqrt{\frac{c_1 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}} + \frac{c_2 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}} \right) \\ &\stackrel{(iii)}{\leq} 1 + \frac{c_2' H^4 S \iota}{\{N_h^k(s, \mathbf{a}) \vee 1\}}, \end{aligned} \quad (97)$$

where the inequality (i) is by the second inequality in event  $\mathcal{E}$  in eq. 20, the inequality (ii) is by Cauchy-Schwartz inequality and the probability distribution sums up to 1, and the last inequality (iii) is from the fact  $\sqrt{ab} \leq a + b$ .

- **Term (ii):** By using the proof-lines of Lemma 14 and assuming that the optimism and pessimism inequality eq. 47 holds for  $(h+1, k)$ , we can bound Term (ii) as

$$\text{Term (ii)} \leq 4H \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \bar{V}_{h+1}^{k, \rho_i} - \underline{V}_{h+1}^{k, \rho_i} \right]. \quad (98)$$

Applying eq. 97 and eq. 98, we get the required bound in Lemma 15.  $\square$

## F PROOF OF REGRET BOUND OF KL-MORNAVI

Similar to (Ghosh et al., 2025), we consider the following definitions:

$$\hat{P}_{\min, h}^k(s, \mathbf{a}) := \min_{s' \in \mathcal{S}} \left\{ \hat{P}_h^k(s'|s, \mathbf{a}) : \hat{P}_h^k(s'|s, \mathbf{a}) > 0 \right\}, \quad (99)$$

$$P_{\min, h}^*(s, \mathbf{a}) := \min_{s' \in \mathcal{S}} \left\{ P_h^*(s'|s, \mathbf{a}) : P_h^*(s'|s, \mathbf{a}) > 0 \right\}, \quad (100)$$

$$P_{\min}^* := \min_{(h, s) \in [H] \times \mathcal{S}} P_{\min, h}^*(s, \pi_h^*(s)), \quad (101)$$

where the following inequality is satisfied:  $P_h^*(s'|s, \mathbf{a}) \geq P_{\min, h}^*(s, \pi_h^*(s)) \geq P_{\min}^*$ .

We now recall the bonus term of KL-MORNAVI for agent  $i$  in episode  $k$  at step  $h$ , as follows:

$$\beta_{i, h}^k(s, \mathbf{a}) = \frac{2c_f H}{\sigma_i} \sqrt{\frac{\iota}{(N_h^k(s, \mathbf{a}) \vee 1) \hat{P}_{\min, h}^k(s, \mathbf{a})}} + \sqrt{\frac{1}{K}}, \quad (102)$$

where  $\hat{P}_{\min, h}^k(s, \mathbf{a}) = \min_{s' \in \mathcal{S}} \{ \hat{P}_h^k(s'|s, \mathbf{a}) : \hat{P}_h^k(s'|s, \mathbf{a}) > 0 \}$ ,  $\iota = \log \left( S^2 (\prod_{i=1}^m A_i) H^2 K^{3/2} / \delta \right)$ , and  $c_f$  is an absolute constant.

Before proceeding to all key lemmas, we introduce the high-probability “typical” event  $\mathcal{E}_{\text{KL}}$  in the lemma below. The proof strategy follows (Lu et al., 2024) and (Ghosh et al., 2025).

**Lemma 16** (Uniform Concentration Bound of event  $\mathcal{E}_{KL}$ ). *Let  $\mathcal{E}_{KL}$  be the event in which, for all  $(s, \mathbf{a}, s', h, k) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H] \times [K]$ , and for all  $\eta$  in a  $\frac{1}{\rho_{\min} S \sqrt{K}}$ -cover of  $[0, H/\rho_{\min}]$ , and is defined as*

$$\begin{aligned} \mathcal{E}_{KL} = & \left\{ \left| \log \left( \mathbb{E}_{\hat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{V_{h+1}}{\eta} \right\} \right] \right) - \log \left( \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{V_{h+1}}{\eta} \right\} \right] \right) \right| \right. \\ & \left. \leq c_1 \sqrt{\frac{\iota}{\{N_h^k(s, \mathbf{a}) \vee 1\} \hat{P}_{\min, h}^k(s, \mathbf{a})}} \right. \\ & \left. \forall (h, s, \mathbf{a}, s', k) \in [H] \times \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [K], \forall \eta \in \mathcal{N}_{\frac{1}{\rho_{\min} S \sqrt{K}}} \left( \left[ 0, \frac{H}{\rho_{\min}} \right] \right) \right\}, \quad (103) \end{aligned}$$

where  $\hat{P}_{\min, h}^k(s, \mathbf{a})$  is defined in eq. 99,  $\iota = \log \left( S^3 \left( \prod_{i=1}^m A_i \right) H^2 K^{3/2} / \delta \right)$ ,  $c_1 > 0$  is an absolute constant and  $\eta \in \mathcal{N}_{\frac{1}{\rho_{\min} S \sqrt{K}}}([0, H/\rho_{\min}])$ , where  $\rho_{\min} = \min_{i \in \mathcal{M}} \rho_i$  and  $\mathcal{N}_{\frac{1}{\rho_{\min} S \sqrt{K}}}([0, H/\rho_{\min}])$  denotes an  $1/(\rho_{\min} S \sqrt{K})$ -cover of the interval  $[0, H/\rho_{\min}]$ .

Then, this event  $\mathcal{E}_{KL}$  occurs with high probability, i.e.,  $\Pr(\mathcal{E}_{KL}) \geq 1 - \delta$ .

*Proof.* The proof follows standard techniques: we apply classical concentration inequalities followed by a union bound. Consider a fixed tuple  $(s, \mathbf{a}, h)$  for a fixed episode  $k$ . Now we consider the following equivalent random process: (i) before the agents starts, the environment samples  $\{s^{(1)}, s^{(2)}, \dots, s^{(k-1)}\}$  independently from  $P_h^*(\cdot|s, \mathbf{a})$ , where  $s^{(i)} \in \mathcal{S}$  denotes the state sampled at episode  $i$ ; (ii) during the interaction between the agents and the environment, the  $i$ -th time the state and joint actions  $(s, \mathbf{a})$  tuple is visited at step  $h$ , the environment will make the agents transit to the next state  $s^{(i)}$ . Note that the randomness induced by this interaction procedure is exactly the same as the original one, which means the probability of any event in this context is the same as in the original problem. Therefore, it suffices to prove the target concentration inequality in this context.

Based on the above fact, we directly apply (Wang et al., 2024e, Lemma 16). To extend the bound uniformly, we apply a union bound over all tuples  $(h, s, \mathbf{a}, s', k, \eta) \in [H] \times \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [K] \times \mathcal{N}_{1/(\rho_{\min} S \sqrt{K})}([0, H/\rho_{\min}])$ . Note that the  $\eta$ -cover for each agent  $i$  lies in the interval  $[0, H/\rho_i] \leq [0, H/\rho_{\min}]$  for all  $i \in \mathcal{M}$ , and this cover contains a valid  $\frac{1}{\rho_i S \sqrt{K}}$ -cover for each agent-specific interval  $[0, \frac{H}{\rho_i}]$ . Therefore, we define the common  $\eta$ -cover as  $\eta \in \mathcal{N}_{\frac{1}{\rho_{\min} S \sqrt{K}}} \left( \left[ 0, \frac{H}{\rho_{\min}} \right] \right)$ , where  $\mathcal{N}_{\frac{1}{\rho_{\min} S \sqrt{K}}} \left( \left[ 0, \frac{H}{\rho_{\min}} \right] \right)$  denotes a  $\frac{1}{\rho_{\min} S \sqrt{K}}$ -cover of the interval  $\left[ 0, \frac{H}{\rho_{\min}} \right]$ .  $\square$

## PROOF OF THEOREM 2 (KL-DRMG SETTING)

*Proof.* With Lemma 19, we can establish an upper bound on the regret by considering the difference between our optimistic and pessimistic value functions:

$$\text{Regret}_{\text{NASH}}(K) = \sum_{k=1}^K \max_{i \in \mathcal{M}} (V_{i,1}^{\dagger, \pi_{-i}^k, \rho_i} - V_{i,1}^{k, \rho_i})(s_1^k) \leq \sum_{k=1}^K \max_{i \in \mathcal{M}} (\bar{V}_{i,1}^{k, \rho_i} - \underline{V}_{i,1}^{k, \rho_i})(s_1^k). \quad (104)$$

For the KL-divergence uncertainty set, we will refer to the bonus term as  $\beta_{i,h}^k(s, \mathbf{a})$ , as given in eq. 102. Our first step is to establish a bound on the difference between the upper and lower Q-values. Given our definitions for  $\bar{Q}_{i,h}^{k, \rho_i}$ ,  $\underline{Q}_{i,h}^{k, \rho_i}$ ,  $\bar{V}_{i,h}^{k, \rho_i}$ ,  $\underline{V}_{i,h}^{k, \rho_i}$ , and the bonus term  $\beta_{i,h}^{k, \rho_i}(s, \mathbf{a})$  as defined in eq. 5 through eq. 102, for any  $(i, h, k, s, \mathbf{a}) \in \mathcal{M} \times [H] \times [K] \times \mathcal{S} \times \mathcal{A}$ , we have

$$\bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) - \underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq \sigma_{\widehat{\mathcal{P}}_{i,h}^{k, \rho_i}(s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i}] - \sigma_{\widehat{\mathcal{P}}_{i,h}^{k, \rho_i}(s, \mathbf{a})} [\underline{V}_{i,h+1}^{k, \rho_i}] + 2\beta_{i,h}^{k, \rho_i}(s, \mathbf{a}). \quad (105)$$

We define the following terms,  $A$  and  $B$ , to simplify our analysis:

$$A := \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a})} \left[ \overline{V}_{i,h+1}^{k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \overline{V}_{i,h+1}^{k, \rho_i} \right] + \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{k, \rho_i} \right] - \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{k, \rho_i} \right]. \quad (106)$$

$$B := \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \overline{V}_{i,h+1}^{k, \rho_i} \right] - \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{k, \rho_i} \right]. \quad (107)$$

By applying eq. 106 and eq. 107 to eq. 105, we obtain:

$$\overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) - \underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq A + B + 2\beta_{i,h}^{k, \rho_i}(s, \mathbf{a}). \quad (108)$$

We can upper bound term  $A$  using a concentration argument tailored for KL robust expectations from Lemma 17, which shows that

$$A \leq 2\beta_{i,h}^{k, \rho_i}(s, \mathbf{a}). \quad (109)$$

For term  $B$ , we use the definition of  $\mathbb{E}_{\mathcal{P}_h^{\rho}(s, \mathbf{a})}[V]$  from eq. 11 to establish the following bound:

$$\begin{aligned} B &= \sup_{\eta \in [0, \frac{H}{\rho_i}]} \left\{ -\eta \log \left( \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{\overline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} \right] \right) - \eta \rho_i \right\} \\ &\quad - \sup_{\eta \in [0, \frac{H}{\rho_i}]} \left\{ -\eta \log \left( \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{\underline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} \right] \right) - \eta \rho_i \right\} \\ &\leq \sup_{\eta \in [0, H/\rho_i]} \eta \left\{ \log \left( \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{\underline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} \right] \right) \right. \\ &\quad \left. - \log \left( \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{\overline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} \right] \right) \right\} \\ &= \sup_{\eta \in [0, H/\rho_i]} \eta \log \left( \frac{\mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{\underline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} \right]}{\mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{\overline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} \right]} \right) \\ &= \sup_{\eta \in [0, H/\rho_i]} \eta \log \left( 1 + \frac{\mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{\underline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} - \exp \left\{ -\frac{\overline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} \right]}{\mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{\overline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} \right]} \right) \\ &\stackrel{(a)}{\leq} \sup_{\eta \in [0, H/\rho_i]} \eta \frac{\mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{\underline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} - \exp \left\{ -\frac{\overline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} \right]}{\mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{\overline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} \right]} \\ &\stackrel{(b)}{\leq} \sup_{\eta \in [\underline{\eta}, H/\rho_i]} \eta \exp \left\{ \frac{H}{\underline{\eta}} \right\} \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{\underline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} - \exp \left\{ -\frac{\overline{V}_{i,h+1}^{k, \rho_i}}{\eta} \right\} \right] \\ &\stackrel{(c)}{\leq} \exp \left\{ \frac{H}{\underline{\eta}} \right\} \mathbb{E}_{P_h^*(s, \mathbf{a})} \left[ \overline{V}_{i,h+1}^{k, \rho_i} - \underline{V}_{i,h+1}^{k, \rho_i} \right], \end{aligned} \quad (110)$$

where inequality (a) uses the fact that  $\log(1+x) \leq x$ , inequality (b) holds because  $0 \leq \overline{V}_{i,h+1}^{k, \rho_i} \leq H$  and  $\eta \in [\underline{\eta}, H/\rho_i]$ , and inequality (c) is due to the  $\frac{1}{\eta}$ -Lipschitz continuity of  $\phi_\eta(x) = \exp \left\{ -\frac{x}{\eta} \right\}$  for  $x \geq 0$ , as well as  $\underline{V}_{i,h+1}^{k, \rho_i} \leq \overline{V}_{i,h+1}^{k, \rho_i}$ .

By applying the bounds for  $A$  and  $B$  to eq. 108, we get

$$\overline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) - \underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq \exp \left\{ \frac{H}{\underline{\eta}} \right\} \mathbb{E}_{P_h^*(s, \mathbf{a})} \left[ \overline{V}_{i,h+1}^{k, \rho_i} - \underline{V}_{i,h+1}^{k, \rho_i} \right] + 4\beta_h^{k, \rho_i}(s, \mathbf{a}). \quad (111)$$

Using Lemma 18 to upper bound the bonus term, and rearranging the terms, we further obtain:

$$\begin{aligned} \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) - \underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) &\leq \exp\left\{\frac{H}{\underline{\eta}}\right\} \mathbb{E}_{P_h^*(s, \mathbf{a})} [\bar{V}_{i,h+1}^{k,\rho_i} - \underline{V}_{i,h+1}^{k,\rho_i}] \\ &\quad + \frac{4c_1 H}{\rho_{\min}} \sqrt{\frac{\iota^2}{\{N_h^k(s, \mathbf{a}) \vee 1\} P_{\min}^*}} + \sqrt{\frac{4}{K}}, \end{aligned} \quad (112)$$

where  $c_1 > 0$  is an absolute constant. From the definitions in eq. 8, the difference in V-functions is given by:

$$\bar{V}_{i,h}^{k,\rho_i}(s) - \underline{V}_{i,h}^{k,\rho_i}(s) = \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [\bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) - \underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})]. \quad (113)$$

We now define a new recursive value function  $\tilde{V}_h^{k,\rho_{\min}}$  and a corresponding Q-function  $\tilde{Q}_h^{k,\rho_{\min}}$  with  $\tilde{V}_{H+1}^{k,\rho_{\min}} = 0$ , where  $\rho_{\min} = \min_{i \in \mathcal{M}} \rho_i$ :

$$\tilde{Q}_h^{k,\rho_{\min}}(s, \mathbf{a}) = \exp\left\{\frac{H}{\underline{\eta}}\right\} \mathbb{E}_{P_h^*(s, \mathbf{a})} [\tilde{V}_{h+1}^{k,\rho_{\min}}] + \frac{4c_1 H}{\rho_{\min}} \sqrt{\frac{\iota^2}{\{N_h^k(s, \mathbf{a}) \vee 1\} P_{\min}^*}} + \sqrt{\frac{4}{K}}. \quad (114)$$

$$\tilde{V}_h^{k,\rho_{\min}}(s) = \mathbb{E}_{\mathbf{a} \sim \pi_h^k(\cdot|s)} [\tilde{Q}_{i,h}^{k,\rho_{\min}}(s, \mathbf{a})]. \quad (115)$$

By an inductive proof, we can show that for any  $(i, h, s, \mathbf{a}) \in \mathcal{M} \times [H] \times \mathcal{S} \times \mathcal{A}$ , the following bounds hold:

$$\max_{i \in \mathcal{M}} (\bar{Q}_{i,h}^{k,\rho_i} - \underline{Q}_{i,h}^{k,\rho_i})(s, \mathbf{a}) \leq \tilde{Q}_h^{k,\rho_{\min}}(s, \mathbf{a}), \quad (116)$$

$$\max_{i \in \mathcal{M}} (\bar{V}_{i,h}^{k,\rho_i} - \underline{V}_{i,h}^{k,\rho_i})(s) \leq \tilde{V}_h^{k,\rho_{\min}}(s). \quad (117)$$

Therefore, our analysis can focus on bounding the sum  $\sum_{k=1}^K \tilde{V}_1^{k,\rho_{\min}}(s_1^k)$ . For simplicity, we introduce the following notations for the differences at any  $(h, k) \in [H] \times [K]$ :

$$\Delta_h^k := \tilde{V}_h^{k,\rho_{\min}}(s_h^k), \quad (118)$$

$$\zeta_h^k := \Delta_h^k - \tilde{Q}_h^{k,\rho_{\min}}(s_h^k, \mathbf{a}_h^k), \quad (119)$$

$$\xi_h^k := \mathbb{E}_{P_h^*(\cdot|s_h^k, \mathbf{a}_h^k)} [\tilde{V}_{h+1}^{k,\rho_{\min}}] - \Delta_{h+1}^k. \quad (120)$$

We can confirm that  $\{\zeta_h^k\}_{(h,k)}$  and  $\{\xi_h^k\}_{(h,k)}$  are martingale difference sequences with respect to their respective filtrations. By substituting eq. 114 into eq. 119, we obtain the recursive relationship:

$$\begin{aligned} \Delta_{i,h}^k &= \zeta_{i,h}^k + \tilde{Q}_h^{k,\rho_{\min}}(s_h^k, \mathbf{a}_h^k) \\ &\leq \zeta_{i,h}^k + \exp\left\{\frac{H}{\underline{\eta}}\right\} \mathbb{E}_{P_h^*(s, \mathbf{a})} [\tilde{V}_{h+1}^{k,\rho_{\min}}] + \frac{4c_1 H}{\rho_{\min}} \sqrt{\frac{\iota^2}{\{N_h^k(s, \mathbf{a}) \vee 1\} P_{\min}^*}} + \sqrt{\frac{4}{K}} \\ &= \zeta_{i,h}^k + \exp\left\{\frac{H}{\underline{\eta}}\right\} \xi_{i,h}^k + \exp\left\{\frac{H}{\underline{\eta}}\right\} \Delta_{i,h+1}^k + \frac{4c_1 H}{\rho_{\min}} \sqrt{\frac{\iota^2}{\{N_h^k(s, \mathbf{a}) \vee 1\} P_{\min}^*}} \\ &\quad + \sqrt{\frac{4}{K}}. \end{aligned} \quad (121)$$

By recursively applying eq. 121 and noting that  $1 \leq \left(\exp\left\{\frac{H}{\underline{\eta}}\right\}\right)^h \leq \left(\exp\left\{\frac{H}{\underline{\eta}}\right\}\right)^H := d_H$ , we can upper bound the right hand side of eq. 104 as:

$$\begin{aligned} \text{Regret}_{\text{NASH}}(K) &\leq \sum_{k=1}^K \Delta_1^k \leq c' d_H \sum_{k=1}^K \sum_{h=1}^H \left\{ (\zeta_h^k + \xi_h^k) \right. \\ &\quad \left. + \left( \frac{4c_1 H}{\rho_{\min}} \sqrt{\frac{\iota^2}{\{N_h^k(s, \mathbf{a}) \vee 1\} P_{\min}^*}} + \sqrt{\frac{4}{K}} \right) \right\}. \end{aligned} \quad (122)$$

Next, we bound each of these two main terms. The first term, a sum of martingale differences, is bounded using the Azuma-Hoeffding inequality from Lemma 26, yielding:

$$\sum_{k=1}^K \sum_{h=1}^H (\zeta_{i,h}^k + \xi_{i,h}^k) \leq c'_1 \sqrt{H^3 K L}, \quad (123)$$

where  $c'_1 > 0$  is an absolute constant. For the second term, we apply the proof lines of (Liu et al., 2021, Theorem 3) to bound the sum of the inverse counts:

$$\sum_{k=1}^K \sum_{h=1}^H \sqrt{\frac{1}{\{N_h^k(s_h^k, \mathbf{a}_h^k) \vee 1\}}} \leq c'_2 \left( \sqrt{H^2 K S \prod_{i \in \mathcal{M}} A_i} + HS \prod_{i \in \mathcal{M}} A_i \right). \quad (124)$$

By applying eq. 124 to the second term of eq. 122, we get the following:

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \left( \frac{4c_1 H}{\rho_{\min}} \sqrt{\frac{\iota^2}{\{N_h^k(s, \mathbf{a}) \vee 1\} P_{\min}^*}} + \sqrt{\frac{4}{K}} \right) &\leq c'_2 \left( \sqrt{\frac{H^4 K S (\prod_{i \in \mathcal{M}} A_i) \iota^2}{\rho_{\min}^2 P_{\min}^*}} \right. \\ &\quad \left. + \frac{H^2 S (\prod_{i \in \mathcal{M}} A_i) \iota}{\rho_{\min} \sqrt{P_{\min}^*}} + \sqrt{H^2 K} \right). \end{aligned} \quad (125)$$

By combining the bounds for both terms in eq. 122, we can upper bound the final regret as follows:

$$\begin{aligned} \text{Regret}_{\text{NASH}}(K) &\leq c' d_H \left( \sqrt{\frac{H^4 K S (\prod_{i \in \mathcal{M}} A_i) \iota^2}{\rho_{\min}^2 P_{\min}^*}} \right) \\ &= \mathcal{O} \left( \sqrt{\frac{H^4 \exp(2H^2) K S (\prod_{i \in \mathcal{M}} A_i) (\iota')^3}{\rho_{\min}^2 P_{\min}^*}} \right). \end{aligned} \quad (126)$$

This completes the proof of Theorem 2.  $\square$

**Remark 2.** The proof techniques for bounding  $\text{Regret}_{\text{CCE}}(K)$  and  $\text{Regret}_{\text{CE}}(K)$  follow the same lines of proof for  $\text{Regret}_{\text{NASH}}(K)$ , leveraging Lemma 20 and Lemma 21, respectively, in the context of KL-DRMG.

## F.1 KEY LEMMAS FOR KL-DRMG

**Lemma 17** (Concentration Bound for Robust Value Estimators in KL-DRMG). *Let  $\mathcal{E}_{KL}$  be the typical event and let the bonus term  $\beta_{i,h}^k$  be set defined in eq. 102. Then, the following inequality holds:*

$$\begin{aligned} \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [\bar{V}_{i,h+1}^{k, \rho_i}] + \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{k, \rho_i}] - \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{k, \rho_i}] \\ \leq \frac{2c_1 H}{\rho_{\min}} \sqrt{\frac{\iota}{\{N_h^k(s, \mathbf{a}) \vee 1\} \widehat{P}_{\min, h}^k(s, \mathbf{a})}} + \sqrt{\frac{2}{K}}, \end{aligned} \quad (127)$$

where  $\iota = \log(S^3 (\prod_{i=1}^m A_i) H^2 K^{3/2} / \delta)$ , and  $c_1 > 0$  is an absolute constant.

*Proof.* We begin by defining the term that we need to bound. Let's denote this term by  $A$ :

$$A := \sigma_{\widehat{\mathcal{P}}_h^{\rho}(s, \mathbf{a})} [\bar{V}_{h+1}^k] - \sigma_{\mathcal{P}_h^{\rho}(s, \mathbf{a})} [\bar{V}_{h+1}^k] + \sigma_{\mathcal{P}_h^{\rho}(s, \mathbf{a})} [V_{h+1}^k] - \sigma_{\widehat{\mathcal{P}}_h^{\rho}(s, \mathbf{a})} [V_{h+1}^k]. \quad (128)$$

Under the high-probability event  $\mathcal{E}_{KL}$ , we can directly apply the concentration inequality given in Lemma 24. This allows us to upper bound  $A$  as follows:

$$A \leq \frac{2c_1 H}{\rho_{\min}} \sqrt{\frac{\iota}{\{N_h^k(s, \mathbf{a}) \vee 1\} \widehat{P}_{\min, h}^k(s, \mathbf{a})}} + \sqrt{\frac{2}{K}}, \quad (129)$$

where  $c_1 > 0$  is an absolute constant and  $\iota = \log(S^3 (\prod_{i=1}^m A_i) H^2 K^{3/2} / \delta)$ . This bound is exactly the bonus term multiplied by a constant. Therefore, based on our choice of  $\beta_{i,h}^k(s, \mathbf{a})$  as defined in eq. 102, the inequality in eq. 127 holds. This completes the proof of Lemma 17.  $\square$

**Lemma 18** (Bound of the bonus term for KL-DRMG). *Let  $\mathcal{E}_{KL}$  be the typical event, the bonus term  $\beta_{i,h}^k$  in eq. 102 is bounded by*

$$\beta_{i,h}^k(s, \mathbf{a}) \leq \frac{c_1 H}{\rho_{\min}} \sqrt{\frac{\iota^2}{\{N_h^k(s, \mathbf{a}) \vee 1\} P_{\min}^*}} + \sqrt{\frac{1}{K}}, \quad (130)$$

where  $\iota = \log(S^3 (\prod_{i=1}^m A_i) H^2 K^{3/2} / \delta)$ , and  $c_1 > 0$  is an absolute constant.

*Proof.* The proof-lines are similar to (Ghosh et al., 2025, Lemma K.7). We recall the choice of  $\beta_{i,h}^k$  as given in eq. 102, i.e.

$$\beta_{i,h}^k(s, \mathbf{a}) = \frac{2c_f H}{\rho_i} \sqrt{\frac{\iota}{\{N_h^k(s, \mathbf{a}) \vee 1\} \hat{P}_{\min,h}^k(s, \mathbf{a})}} + \sqrt{\frac{1}{K}}, \quad (131)$$

where  $\iota = \log(S^3 (\prod_{i=1}^m A_i) H^2 K^{3/2} / \delta)$ ,  $\hat{P}_{\min,h}^k(s, \mathbf{a})$  is defined in eq. 99, and  $c_f > 0$  is an absolute constant.

By Lemma 25 and the union bound, it holds that with probability at least  $1 - \delta$  that for all  $(h, s, \mathbf{a}) \in [H] \times \mathcal{S} \times \mathcal{A}$ , we get

$$\forall s' \in \mathcal{S} : P_h^*(s' | s, \mathbf{a}) \geq \frac{\hat{P}_h^k(s' | s, \mathbf{a})}{e^2} \geq \frac{P_h^*(s' | s, \mathbf{a})}{8e^{2\iota}}. \quad (132)$$

To characterize the relation between  $P_{\min,h}^*(s, \mathbf{a})$  and  $\hat{P}_{\min,h}^k(s, \mathbf{a})$  for any  $(h, s, \mathbf{a}) \in [H] \times \mathcal{S} \times \mathcal{A}$ , we suppose—without loss of generality—that  $P_{\min,h}^*(s, \mathbf{a}) = P_h^*(s_1 | s, \mathbf{a})$  and  $\hat{P}_{\min,h}^k(s, \mathbf{a}) = \hat{P}_h^k(s_2 | s, \mathbf{a})$  for some  $s_1, s_2 \in \mathcal{S}$ . Then, it follows that

$$\begin{aligned} P_{\min,h}^*(s, \mathbf{a}) &= P_h^*(s_1 | s, \mathbf{a}) \\ &\stackrel{(i)}{\geq} \frac{\hat{P}_h^k(s_1 | s, \mathbf{a})}{e^2} \geq \frac{\hat{P}_{\min,h}^k(s, \mathbf{a})}{e^2} \\ &= \frac{\hat{P}_h^k(s_2 | s, \mathbf{a})}{e^2} \stackrel{(ii)}{\geq} \frac{P_h^*(s_2 | s, \mathbf{a})}{8e^{2\iota}} \\ &\geq \frac{P_{\min,h}^*(s, \mathbf{a})}{8e^{2\iota}} \stackrel{(iii)}{\geq} \frac{P_{\min}^*}{8e^{2\iota}}. \end{aligned} \quad (133)$$

where the inequalities (i) and (ii) follow from eq. 132, and inequality (iii) follows by eq. 101.

By applying eq. 133 in eq. 131, we get

$$\begin{aligned} \beta_{i,h}^k(s, \mathbf{a}) &\leq \frac{2c_f H}{\rho_i} \sqrt{\frac{\iota^2}{\{N_h^k(s, \mathbf{a}) \vee 1\} P_{\min}^*}} + \sqrt{\frac{1}{K}} \leq \frac{c_1 H}{\rho_{\min}} \sqrt{\frac{\iota^2}{\{N_h^k(s, \mathbf{a}) \vee 1\} P_{\min}^*}} \\ &\quad + \sqrt{\frac{1}{K}}. \end{aligned} \quad (134)$$

This concludes the proof of Lemma 18.  $\square$

#### NE VERSION: OPTIMISTIC AND PESSIMISTIC ESTIMATION OF THE ROBUST VALUES FOR KL-DRMG.

Here we will proof the optimistic estimations are indeed upper bounds of the corresponding robust V-value and robust Q-value functions fro NE version.

**Lemma 19** (Optimistic and pessimistic estimation of the robust values for KL-DRMG for NE Version). *Under the event  $\mathcal{E}_{KL}$  and by setting the bonus term  $\beta_{i,h}^k$  as in eq. 102, it holds that*

$$Q_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}), \quad \underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}), \quad (135)$$

$$V_{i,h}^{\dagger, \pi_{-i}^k, \rho_i}(s) \leq \bar{V}_{i,h}^{k, \rho_i}(s), \quad \underline{V}_{i,h}^{k, \rho_i}(s) \leq V_{i,h}^{\pi^k, \rho_i}(s). \quad (136)$$

*Proof.* The proof-lines are similar to (Ghosh et al., 2025) adapted to the multi-agent case. We will run a proof for each inequality outlined in Lemma 19

- **Ineq. 1:** To prove  $Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})$ .
- **Ineq. 2:** To prove  $\underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k,\rho_i}(s, \mathbf{a})$ .

Assume that both eq. 135 and eq. 136 hold at the  $(h+1)$ -th step.

- **Proof of Ineq. 1:** We first consider robust  $Q$  at the  $h$ -th step. Then, by Proposition 1 (Robust Bellman Equation) and eq. 5, we have that

$$\begin{aligned} Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) - \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) &= \max \left\{ \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \bar{V}_{i,h+1}^{k,\rho_i} \right] \right. \\ &\quad \left. - \beta_{i,h}^k(s, \mathbf{a}), Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) - H \right\}, \\ &\leq \max \left\{ \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] \right. \\ &\quad \left. - \beta_{i,h}^k(s, \mathbf{a}), 0 \right\}, \end{aligned} \quad (137)$$

where the second inequality follows from the induction of  $V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \leq \bar{V}_{i,h+1}^{k,\rho_i}$  at the  $h+1$ -th step and the fact that  $Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i} \leq H$ . By Lemma 22 and by the definition of  $\hat{P}_{\min,h}^k(s, \mathbf{a})$  as given in eq. 99, we have that

$$\begin{aligned} \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] &\leq \frac{c_1 H}{\rho_i} \sqrt{\frac{L}{\{N_h^k(s, \mathbf{a}) \vee 1\} \hat{P}_{\min,h}^k(s, \mathbf{a})}} \\ &\quad + \sqrt{\frac{1}{K}}. \end{aligned} \quad (138)$$

By the choice of  $\beta_{i,h}^k$  in eq. 102 and eq. 138 and applying in eq. 137, we conclude that

$$Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}). \quad (139)$$

- **Proof of Ineq. 2:** By using Proposition 1 (Robust Bellman Equation) and eq. 6, we have that

$$\begin{aligned} \underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) - Q_{i,h}^{\pi^k,\rho_i}(s, \mathbf{a}) &= \max \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \underline{V}_{i,h+1}^{k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] - \beta_{i,h}^k(s, \mathbf{a}), \right. \\ &\quad \left. 0 - Q_{i,h}^{\pi^k,\rho_i}(s, \mathbf{a}) \right\} \end{aligned} \quad (140)$$

$$\begin{aligned} &\leq \max \left\{ \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] - \beta_{i,h}^k(s, \mathbf{a}), \right. \\ &\quad \left. 0 \right\}, \end{aligned} \quad (141)$$

where the second inequality follows from the induction of  $\underline{V}_{i,h+1}^{k,\rho_i} \leq V_{i,h+1}^{\pi^k,\rho_i}$  at the  $(h+1)$ -th step and the fact that  $Q_{i,h}^{\pi^k,\rho_i} \geq 0$ . By Lemma 23, we get

$$\begin{aligned} \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s,\mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] &\leq \frac{c_1 H}{\rho_i} \sqrt{\frac{L}{\{N_h^k(s,\mathbf{a}) \vee 1\} \widehat{P}_{\min,h}^k(s,\mathbf{a})}} \\ &\quad + \sqrt{\frac{1}{K}}. \end{aligned} \quad (142)$$

By the choice of  $\beta_{i,h}^k$  in eq. 102 and eq. 142 and applying in eq. 141, we conclude that

$$Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s,\mathbf{a}) \leq \overline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}). \quad (143)$$

Therefore, by eq. 139 and eq. 143, we have proved that at step  $h$ , it holds that

$$Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s,\mathbf{a}) \leq \overline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}), \quad \underline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}) \leq Q_{i,h}^{\pi^k,\rho_i}(s,\mathbf{a}). \quad (144)$$

We now assume that eq. 135 hold for  $h$ -th step. Then, by the definition of robust value function as given by robust Bellman equation (Proposition 1), eq. 8, and NASH Equilibrium, we get

$$\overline{V}_{i,h}^{k,\rho_i}(s) = \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} \left[ \overline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}) \right] = \max_{\pi'_i} \mathbb{E}_{\mathbf{a} \sim \pi'_i \times \pi_{-i}^k(\cdot|s)} \left[ \overline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}) \right]. \quad (145)$$

By the definition of  $V_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s)$  in eq. 3, we get

$$V_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s) = \max_{\pi'_i} \mathbb{E}_{\mathbf{a} \sim \pi'_i \times \pi_{-i}^k(\cdot|s)} \left[ Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s,\mathbf{a}) \right]. \quad (146)$$

Sine by induction, for any  $(s,\mathbf{a})$ ,  $\overline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}) \geq Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s,\mathbf{a})$ . As a result, we also have  $\overline{V}_{i,h}^{k,\rho_i}(s) \geq V_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s)$ , which is eq. 136 for  $h$ -th step. Similarly, we can show that

$$\begin{aligned} \underline{V}_{i,h}^{k,\rho_i}(s) &= \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} \left[ \underline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}) \right], \\ &\stackrel{(i)}{\leq} \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} \left[ Q_{i,h}^{\pi^k,\rho_i}(s,\mathbf{a}) \right], \\ &\stackrel{(ii)}{=} V_{i,h}^{\pi^k,\rho_i}(s), \end{aligned} \quad (147)$$

where (i) is due to the fact that  $\underline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}) \leq Q_{i,h}^{\pi^k,\rho_i}(s,\mathbf{a})$  and (ii) is by definition of  $V_{i,h}^{\pi^k,\rho_i}(s)$  as given by Bellman equation in Proposition 1.  $\square$

#### CCE VERSION: OPTIMISTIC AND PESSIMISTIC ESTIMATION OF THE ROBUST VALUES FOR KL-DRMG.

Here we will proof the optimistic estimations are indeed upper bounds of the corresponding robust V-value and robust Q-value functions fro CCE version.

**Lemma 20** (Optimistic and pessimistic estimation of the robust values for KL-DRMG for CCE Version). *Under the event  $\mathcal{E}_{KL}$  and by setting the bonus term  $\beta_{i,h}^k$  as in eq. 102, it holds that*

$$Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s,\mathbf{a}) \leq \overline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}), \quad \underline{Q}_{i,h}^{k,\rho_i}(s,\mathbf{a}) \leq Q_{i,h}^{\pi^k,\rho_i}(s,\mathbf{a}), \quad (148)$$

$$V_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s) \leq \overline{V}_{i,h}^{k,\rho_i}(s), \quad \underline{V}_{i,h}^{k,\rho_i}(s) \leq V_{i,h}^{\pi^k,\rho_i}(s). \quad (149)$$

*Proof.* The proof-lines are similar to (Ghosh et al., 2025) adapted to the multi-agent case. We will run a proof for each inequality outlined in Lemma 20



- **Ineq. 1:** To prove  $Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})$ .
- **Ineq. 2:** To prove  $\underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k,\rho_i}(s, \mathbf{a})$ .

Assume that both eq. 148 and eq. 149 hold at the  $(h+1)$ -th step.

- **Proof of Ineq. 1:** We first consider robust  $Q$  at the  $h$ -th step. Then, by Proposition 1 (Robust Bellman Equation) and eq. 5, we have that

$$\begin{aligned} Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) - \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) &= \max \left\{ \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}} \left[ \bar{V}_{i,h+1}^{k,\rho_i} \right] \right. \\ &\quad \left. - \beta_{i,h}^k(s, \mathbf{a}), Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) - H \right\}, \\ &\leq \max \left\{ \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] \right. \\ &\quad \left. - \beta_{i,h}^k(s, \mathbf{a}), 0 \right\}, \end{aligned} \quad (150)$$

where the second inequality follows from the induction of  $V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \leq \bar{V}_{i,h+1}^{k,\rho_i}$  at the  $h+1$ -th step and the fact that  $Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i} \leq H$ . By Lemma 22 and by the definition of  $\hat{P}_{\min,h}^k(s, \mathbf{a})$  as given in eq. 99, we have that

$$\begin{aligned} \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] - \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}} \left[ V_{i,h+1}^{\dagger,\pi_{-i}^k,\rho_i} \right] &\leq \frac{c_1 H}{\rho_i} \sqrt{\frac{L}{\{N_h^k(s, \mathbf{a}) \vee 1\} \hat{P}_{\min,h}^k(s, \mathbf{a})}} \\ &\quad + \sqrt{\frac{1}{K}}. \end{aligned} \quad (151)$$

By the choice of  $\beta_{i,h}^k$  in eq. 102 and eq. 151 and applying in eq. 150, we conclude that

$$Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}). \quad (152)$$

- **Proof of Ineq. 2:** By using Proposition 1 (Robust Bellman Equation) and eq. 6, we have that

$$\begin{aligned} \underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) - Q_{i,h}^{\pi^k,\rho_i}(s, \mathbf{a}) &= \max \left\{ \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}} \left[ \underline{V}_{i,h+1}^{k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] \right. \\ &\quad \left. - \beta_{i,h}^k(s, \mathbf{a}), 0 - Q_{i,h}^{\pi^k,\rho_i}(s, \mathbf{a}) \right\} \\ &\leq \max \left\{ \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] \right. \\ &\quad \left. - \beta_{i,h}^k(s, \mathbf{a}), 0 \right\}, \end{aligned} \quad (153)$$

where the second inequality follows from the induction of  $\underline{V}_{i,h+1}^{k,\rho_i} \leq V_{i,h+1}^{\pi^k,\rho_i}$  at the  $(h+1)$ -th step and the fact that  $Q_{i,h}^{\pi^k,\rho_i} \geq 0$ . By Lemma 23, we get

$$\begin{aligned} \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\pi^k,\rho_i} \right] &\leq \frac{c_1 H}{\rho_i} \sqrt{\frac{L}{\{N_h^k(s, \mathbf{a}) \vee 1\} \hat{P}_{\min,h}^k(s, \mathbf{a})}} \\ &\quad + \sqrt{\frac{1}{K}}. \end{aligned} \quad (154)$$

By the choice of  $\beta_{i,h}^k$  in eq. 102 and eq. 154 and applying in eq. 153, we conclude that

$$Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}). \quad (155)$$

Therefore, by eq. 152 and eq. 155, we have proved that at step  $h$ , it holds that

$$Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}), \quad \underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi_{-i}^k,\rho_i}(s, \mathbf{a}). \quad (156)$$

We now assume that eq. 148 hold for  $h$ -th step. Then, by the definition of robust value function as given by robust Bellman equation (Proposition 1), eq. 8, and CCE Equilibrium, we get

$$\bar{V}_{i,h}^{k,\rho_i}(s) = \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [\bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})] \geq \max_{\pi'_i} \mathbb{E}_{\mathbf{a} \sim \pi'_i \times \pi_{-i}^k(\cdot|s)} [\bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})]. \quad (157)$$

By the definition of  $V_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s)$  in eq. 3, we get

$$V_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s) = \max_{\pi'_i} \mathbb{E}_{\mathbf{a} \sim \pi'_i \times \pi_{-i}^k(\cdot|s)} [Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a})]. \quad (158)$$

Sine by induction, for any  $(s, \mathbf{a})$ ,  $\bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \geq Q_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s, \mathbf{a})$ . As a result, we also have  $\bar{V}_{i,h}^{k,\rho_i}(s) \geq V_{i,h}^{\dagger,\pi_{-i}^k,\rho_i}(s)$ , which is eq. 149 for  $h$ -th step. Similarly, we can show that

$$\begin{aligned} \underline{V}_{i,h}^{k,\rho_i}(s) &= \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [\underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})], \\ &\stackrel{(i)}{\leq} \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [Q_{i,h}^{\pi_{-i}^k,\rho_i}(s, \mathbf{a})], \\ &\stackrel{(ii)}{=} V_{i,h}^{\pi_{-i}^k,\rho_i}(s), \end{aligned} \quad (159)$$

where (i) is due to the fact that  $\underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi_{-i}^k,\rho_i}(s, \mathbf{a})$  and (ii) is by definition of  $V_{i,h}^{\pi_{-i}^k,\rho_i}(s)$  as given by Bellman equation in Proposition 1.  $\square$

CE VERSION: OPTIMISTIC AND PESSIMISTIC ESTIMATION OF THE ROBUST VALUES FOR KL-DRMG.

Here we will proof the optimistic estimations are indeed upper bounds of the corresponding robust V-value and robust Q-value functions fro CE version.

**Lemma 21** (Optimistic and pessimistic estimation of the robust values for KL-DRMG for CE version). *By setting the bonus term  $\beta_{i,h}^k$  as in eq. 102, with probability  $1 - \delta$ , for any  $(s, \mathbf{a}, h, i)$  and  $k \in [K]$ , it holds that*

$$\max_{\phi \in \Phi_i} Q_{i,h}^{\phi \diamond \pi^k,\rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}), \quad \underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi_{-i}^k,\rho_i}(s, \mathbf{a}), \quad (160)$$

$$\max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k,\rho_i}(s) \leq \bar{V}_{i,h}^{k,\rho_i}(s), \quad \underline{V}_{i,h}^{k,\rho_i}(s) \leq V_{i,h}^{\pi_{-i}^k,\rho_i}(s). \quad (161)$$

*Proof.* The proof-lines are similar to (Ghosh et al., 2025) adapted to the multi-agent case. We will run a proof for each inequality outlined in Lemma 21

- **Ineq. 1:** To prove  $\max_{\phi \in \Phi_i} Q_{i,h}^{\phi \diamond \pi^k,\rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a})$ .
- **Ineq. 2:** To prove  $\underline{Q}_{i,h}^{k,\rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi_{-i}^k,\rho_i}(s, \mathbf{a})$ .

Assume that both eq. 160 and eq. 161 hold at the  $(h+1)$ -th step.

- **Proof of Ineq. 1:** We first consider robust  $Q$  at the  $h$ -th step. Then, by Proposition 1 (Robust Bellman Equation) and eq. 5, we have that

$$\begin{aligned}
& \max_{\phi \in \Phi_i} Q_{i,h}^{\phi \diamond \pi^k, \rho_i}(s, \mathbf{a}) - \bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \\
&= \max \left\{ \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k, \rho_i} \right] - \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}} \left[ \bar{V}_{i,h+1}^{k, \rho_i} \right] - \beta_{i,h}^k(s, \mathbf{a}), \right. \\
&\quad \left. \max_{\phi \in \Phi_i} Q_{i,h}^{\phi \diamond \pi^k, \rho_i}(s, \mathbf{a}) - H \right\} \\
&\leq \max \left\{ \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k, \rho_i} \right] - \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}} \left[ \max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k, \rho_i} \right] - \beta_{i,h}^k(s, \mathbf{a}), \right. \\
&\quad \left. 0 \right\}, \tag{162}
\end{aligned}$$

where the second inequality follows from the induction of  $\max_{\phi \in \Phi_i} V_{i,h+1}^{\phi \diamond \pi^k, \rho_i}(s) \leq \bar{V}_{i,h+1}^{k, \rho_i}(s)$  at the  $h+1$ -th step and the fact that  $\max_{\phi \in \Phi_i} Q_{i,h}^{\phi \diamond \pi^k, \rho_i}(s, \mathbf{a}) \leq H$ . By Lemma 22 and by the definition of  $\hat{P}_{\min,h}^k(s, \mathbf{a})$  as given in eq. 99, we have that

$$\begin{aligned}
& \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ \max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k, \rho_i}(s) \right] - \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}} \left[ \max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k, \rho_i}(s) \right] \\
&\leq \frac{c_1 H}{\rho_i} \sqrt{\frac{L}{\{N_h^k(s, \mathbf{a}) \vee 1\} \hat{P}_{\min,h}^k(s, \mathbf{a})}} + \sqrt{\frac{1}{K}}. \tag{163}
\end{aligned}$$

By the choice of  $\beta_{i,h}^k$  in eq. 102 and eq. 163 and applying in eq. 162, we conclude that

$$\max_{\phi \in \Phi_i} Q_{i,h}^{\phi \diamond \pi^k, \rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}). \tag{164}$$

- **Proof of Ineq. 2:** By using Proposition 1 (Robust Bellman Equation) and eq. 6, we have that

$$\begin{aligned}
& \underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) - Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}) \\
&= \max \left\{ \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}} \left[ \underline{V}_{i,h+1}^{k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\pi^k, \rho_i} \right] - \beta_{i,h}^k(s, \mathbf{a}), 0 - Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}) \right\}, \\
&\leq \max \left\{ \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}} \left[ \underline{V}_{i,h+1}^{k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\pi^k, \rho_i} \right] - \beta_{i,h}^k(s, \mathbf{a}), 0 \right\}, \tag{165}
\end{aligned}$$

where the second inequality follows from the induction of  $\underline{V}_{i,h+1}^{k, \rho_i} \leq V_{i,h+1}^{\pi^k, \rho_i}$  at the  $(h+1)$ -th step and the fact that  $Q_{i,h}^{\pi^k, \rho_i} \geq 0$ . By Lemma 23, we get

$$\begin{aligned}
& \sigma_{\widehat{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})}} \left[ \underline{V}_{i,h+1}^{k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} \left[ V_{i,h+1}^{\pi^k, \rho_i} \right] \leq \frac{c_1 H}{\rho_i} \sqrt{\frac{L}{\{N_h^k(s, \mathbf{a}) \vee 1\} \hat{P}_{\min,h}^k(s, \mathbf{a})}} \\
&\quad + \sqrt{\frac{1}{K}}. \tag{166}
\end{aligned}$$

By the choice of  $\beta_{i,h}^k$  in eq. 102 and eq. 166 and applying in eq. 165, we conclude that

$$\underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}). \tag{167}$$

Therefore, by eq. 164 and eq. 167, we have proved that at step  $h$ , it holds that

$$\max_{\phi \in \Phi_i} Q_{i,h}^{\phi \diamond \pi, \rho_i}(s, \mathbf{a}) \leq \bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}), \quad \underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a}). \quad (168)$$

We now assume that eq. 160 hold for  $h$ -th step. Then, by the definition of robust value function as given by robust Bellman equation (Proposition 1), eq. 8, and CE Equilibrium, we get

$$\bar{V}_{i,h}^{k, \rho_i}(s) = \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [\bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a})] = \max_{\phi \in \Phi_i} \mathbb{E}_{\mathbf{a} \sim \phi \diamond \pi^k(\cdot|s)} [\bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a})]. \quad (169)$$

By the definition of  $\max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k, \rho_i}(s)$  in eq. 3, we get

$$\max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k, \rho_i}(s) = \max_{\phi \in \Phi_i} \mathbb{E}_{\mathbf{a} \sim \phi \diamond \pi^k(\cdot|s)} \left[ \max_{\phi'} Q_{i,h}^{\phi' \diamond \pi^k, \rho_i}(s, \mathbf{a}) \right]. \quad (170)$$

Since by induction, for any  $(s, \mathbf{a})$ ,  $\bar{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \geq \max_{\phi \in \Phi_i} Q_{i,h}^{\phi \diamond \pi^k, \rho_i}(s, \mathbf{a})$ . As a result, we also have

$\bar{V}_{i,h}^{k, \rho_i}(s) \geq \max_{\phi \in \Phi_i} V_{i,h}^{\phi \diamond \pi^k, \rho_i}(s)$ , which is eq. 161 for  $h$ -th step. Similarly, we can show that

$$\begin{aligned} \underline{V}_{i,h}^{k, \rho_i}(s) &= \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [\underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a})], \\ &\stackrel{(i)}{\leq} \mathbb{E}_{\mathbf{a} \sim \pi^k(\cdot|s)} [Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a})], \\ &\stackrel{(ii)}{=} V_{i,h}^{\pi^k, \rho_i}(s), \end{aligned} \quad (171)$$

where (i) is due to the fact that  $\underline{Q}_{i,h}^{k, \rho_i}(s, \mathbf{a}) \leq Q_{i,h}^{\pi^k, \rho_i}(s, \mathbf{a})$  and (ii) is by definition of  $V_{i,h}^{\pi^k, \rho_i}(s)$  as given by Bellman equation in Proposition 1.  $\square$

## F.2 AUXILIARY LEMMAS FOR KL-DRMG

**Lemma 22** (Concentration of Value Function in KL-DRMG). *Under the typical event  $\mathcal{E}_{KL}$  as defined in eq. 103, the following concentration bound holds with probability at least  $1 - \delta$ :*

$$\left| \sigma_{\widehat{\mathcal{P}}_h^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i}] - \sigma_{\mathcal{P}_h^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i}] \right| \leq \frac{c_1 H}{\rho_i} \sqrt{\frac{L}{\{N_h^k(s, \mathbf{a}) \vee 1\} \widehat{P}_{\min, h}^k(s, \mathbf{a})}} + \frac{1}{\sqrt{K}},$$

where  $\iota = \log \left( S^3 \left( \prod_{i=1}^m A_i \right) H^2 K^{3/2} / \delta \right)$  and  $c_1$  is an absolute constant.

*Proof.* This proof establishes a concentration bound for the difference between the empirical and true robust value functions. We use the definition of the KL-divergence operator  $\sigma_{\mathcal{P}_h^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i}]$  from eq. 11 and the empirical minimum probability  $\widehat{P}_{\min, h}^k(s, \mathbf{a})$  from eq. 99 to express this difference as a supremum:

$$\begin{aligned} &\left| \sigma_{\widehat{\mathcal{P}}_h^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i}] - \sigma_{\mathcal{P}_h^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i}] \right| \\ &\leq \sup_{\eta \in [\underline{\eta}, H/\rho_i]} \eta \left| \log \left( \mathbb{E}_{\widehat{P}_h^k(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i}}{\eta} \right\} \right] \right) \right. \\ &\quad \left. - \log \left( \mathbb{E}_{P_h^*(\cdot|s, \mathbf{a})} \left[ \exp \left\{ -\frac{V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i}}{\eta} \right\} \right] \right) \right|. \end{aligned} \quad (172)$$

Under the high-probability event  $\mathcal{E}_{KL}$  (defined in eq. 103), we apply a known concentration inequality from (Wang et al., 2024e, Lemma 16) to bound this expression:

$$\left| \sigma_{\widehat{\mathcal{P}}_h^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i}] - \sigma_{\mathcal{P}_h^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\dagger, \pi_{-i}^k, \rho_i}] \right| \leq \frac{c_1 H}{\rho_i} \sqrt{\frac{L}{\{N_h^k(s, \mathbf{a}) \vee 1\} \widehat{P}_{\min, h}^k(s, \mathbf{a})}}, \quad (173)$$

This bound holds for any  $\eta$  within a fine-grained cover of the interval  $[0, H/\rho_{\min}]$ . By applying a standard covering argument, we extend this bound to hold for all  $\eta \in [0, H/\rho_{\min}]$ , thereby concluding the proof of Lemma 22.  $\square$

**Lemma 23** (Bound for DRMG-KL and the robust value function of  $\pi^k$ ). *Under event  $\mathcal{E}_{KL}$  in eq. 103 and for any EQUILIBRIUM  $\in \{NASH, CE, CCE\}$ , we assume that the optimism and pessimism inequalities hold at  $(h+1, k)$ , where these inequalities can correspond to any of the following cases of EQUILIBRIUM:*

- **NE:** Lemma 19 using eq. 135 and eq. 136,
- **CCE:** Lemma 20 using eq. 148 and eq. 149,
- **CE:** Lemma 21 using eq. 160 and eq. 161.

Then the following bound holds:

$$\left| \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\pi^k, \rho_i}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\pi^k, \rho_i}] \right| \leq \frac{c_1 H}{\rho_i} \sqrt{\frac{L}{\{N_h^k(s, \mathbf{a}) \vee 1\} \widehat{P}_{\min, h}^k(s, \mathbf{a})}} + \frac{1}{\sqrt{K}},$$

where  $\iota = \log(S^3 (\prod_{i=1}^m A_i) H^2 K^{3/2} / \delta)$ , and  $c_1$  is an absolute constant.

*Proof.* This proof establishes a concentration bound for the difference between the empirical and true robust value functions under the KL-divergence. By using the definition of the robust operator  $\sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\pi^k, \rho_i}]$  from eq. 11 and the empirical minimum probability  $\widehat{P}_{\min, h}^k(s, \mathbf{a})$  from eq. 99, we can bound the absolute difference as follows:

$$\begin{aligned} & \left| \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\pi^k, \rho_i}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\pi^k, \rho_i}] \right| \\ & \leq \sup_{\eta \in [\underline{\eta}, H/\rho_i]} \eta \left| \log \left( \mathbb{E}_{\widehat{P}_h^k(\cdot | s, \mathbf{a})} \left[ \exp \left\{ -\frac{V_{i,h+1}^{\pi^k, \rho_i}}{\eta} \right\} \right] \right) \right. \\ & \quad \left. - \log \left( \mathbb{E}_{P_h^*(\cdot | s, \mathbf{a})} \left[ \exp \left\{ -\frac{V_{i,h+1}^{\pi^k, \rho_i}}{\eta} \right\} \right] \right) \right|. \end{aligned} \quad (174)$$

Under the high-probability event  $\mathcal{E}_{KL}$  (defined in eq. 103), and by applying a known concentration inequality from (Wang et al., 2024e, Lemma 17), we can establish a uniform bound on this difference:

$$\left| \sigma_{\widehat{\mathcal{P}}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\pi^k, \rho_i}] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}(s, \mathbf{a})} [V_{i,h+1}^{\pi^k, \rho_i}] \right| \leq \frac{c_1 H}{\rho_i} \sqrt{\frac{L}{\{N_h^k(s, \mathbf{a}) \vee 1\} \widehat{P}_{\min, h}^k(s, \mathbf{a})}}. \quad (175)$$

This inequality holds for any  $\eta$  in a fine-grained cover of the interval  $[0, H/\rho_{\min}]$ . We conclude the proof of Lemma 23 by using a standard covering argument to extend the bound to all  $\eta \in [0, H/\rho_{\min}]$ .  $\square$

**Lemma 24** (Bounds for RMG-KL and optimistic and pessimistic robust value estimators). *Under event  $\mathcal{E}_{KL}$  in eq. 103 and for any EQUILIBRIUM  $\in \{NASH, CE, CCE\}$ , we assume that the optimism and pessimism inequalities hold at  $(h+1, k)$ , where these inequalities can correspond to any of the following cases of EQUILIBRIUM:*

- **NE:** Lemma 19 using eq. 135 and eq. 136,
- **CCE:** Lemma 20 using eq. 148 and eq. 149,
- **CE:** Lemma 21 using eq. 160 and eq. 161.

Then the following bound holds:

$$\begin{aligned} & \max \left\{ \left| \widehat{\sigma_{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a})} \left[ \overline{V}_{i,h+1}^{k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a}) \left[ \overline{V}_{i,h+1}^{k, \rho_i} \right] \right|, \left| \widehat{\sigma_{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a})} \left[ V_{i,h+1}^{k, \rho_i} \right] - \sigma_{\mathcal{P}_{i,h}^{\rho_i}}(s, \mathbf{a}) \left[ V_{i,h+1}^{k, \rho_i} \right] \right| \right\} \\ & \leq \frac{c_1 H}{\rho_i} \sqrt{\frac{L}{\{N_h^k(s, \mathbf{a}) \vee 1\} \widehat{P}_{\min, h}^k(s, \mathbf{a})}} + \sqrt{\frac{1}{K}}, \end{aligned}$$

where  $\iota = \log(S^3(\prod_{i=1}^n A_i) H^2 K^{3/2}/\delta)$  and  $c_1$  is an absolute constant.

*Proof.* We follow the same proof lines as Lemma 23, and thereby we omit it.  $\square$

**Lemma 25** (Bound on Binomial random variable). *Suppose  $X \sim \text{Binomial}(n, p)$ , where  $n \geq 1$  and  $p \in [0, 1]$ . For any  $\delta \in (0, 1]$ , we have*

$$X \geq \frac{np}{8 \log(\frac{1}{\delta})}, \quad \text{if } np \geq 8 \log\left(\frac{1}{\delta}\right), \quad (176)$$

$$X \leq \begin{cases} e^2 np, & \text{if } np \geq \log\left(\frac{1}{\delta}\right), \\ 2e^2 \log\left(\frac{1}{\delta}\right), & \text{if } np \leq 2 \log\left(\frac{1}{\delta}\right), \end{cases} \quad (177)$$

hold with probability at least  $1 - 4\delta$ .

*Proof.* Refer to (Shi et al., 2023, Lemma 8) for details.  $\square$

## G OTHER TECHNICAL LEMMAS

Here, we present some auxiliary lemmas which are useful in the proof.

**Lemma 26** (Azuma Hoeffding’s Inequality). *Let  $\{Z_t\}_{t \in \mathbb{Z}_+}$  be a martingale with respect to the filtration  $\{\mathcal{F}_t\}_{t \in \mathbb{Z}_+}$ . Assume that there are predictable processes  $\{A_t\}_{t \in \mathbb{Z}_+}$  and  $\{B_t\}_{t \in \mathbb{Z}_+}$  with respect to  $\{\mathcal{F}_t\}_{t \in \mathbb{Z}_+}$ , i.e., for all  $t$ ,  $A_t$  and  $B_t$  are  $\mathcal{F}_{t-1}$ -measurable, and constants  $0 < c_1, c_2, \dots < +\infty$  such that  $A_t \leq Z_t - Z_{t-1} \leq B_t$  and  $B_t - A_t \leq c_t$  almost surely. Then, for all  $\beta > 0$*

$$\mathbb{P}\left(|Z_t - Z_0| \geq \beta\right) \leq \exp\left\{-\frac{2\beta^2}{\sum_{i \leq t} c_i^2}\right\}. \quad (178)$$

*Proof.* Refer to the proof of Theorem 5.1 of (Dubhashi & Panconesi, 2009).  $\square$

**Lemma 27** (Self-bounding variance inequality (Maurer & Pontil, 2009, Theorem 10)). *Let  $X_1, \dots, X_T$  be independent and identically distributed random variables with finite variance, that is,  $\text{Var}(X_1) < \infty$ . Assume that  $X_t \in [0, M]$  for every  $t$  with  $M > 0$ , and let*

$$S_T^2 = \frac{1}{T} \sum_{t=1}^T X_t^2 - \left(\frac{1}{T} \sum_{t=1}^T X_t\right)^2.$$

Then, for any  $\varepsilon > 0$ , we have

$$\mathbb{P}\left(\left|S_T - \sqrt{\text{Var}(X_1)}\right| \geq \varepsilon\right) \leq 2 \exp\left(-\frac{T\varepsilon^2}{2M^2}\right).$$

*Proof.* Refer to the proof of Lemma 7 of (Panaganti & Kalathil, 2022).  $\square$