

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 GAP-DEPENDENT BOUNDS FOR Q -LEARNING USING REFERENCE-ADVANTAGE DECOMPOSITION

Anonymous authors

Paper under double-blind review

ABSTRACT

We study the gap-dependent bounds of two important algorithms for on-policy Q -learning for finite-horizon episodic tabular Markov Decision Processes (MDPs): UCB-Advantage (Zhang et al. 2020) and Q-EarlySettled-Advantage (Li et al. 2021). UCB-Advantage and Q-EarlySettled-Advantage improve upon the results based on Hoeffding-type bonuses and achieve the almost optimal \sqrt{T} -type regret bound in the worst-case scenario, where T is the total number of steps. However, the benign structures of the MDPs such as a strictly positive suboptimality gap can significantly improve the regret. While gap-dependent regret bounds have been obtained for Q -learning with Hoeffding-type bonuses, it remains an open question to establish gap-dependent regret bounds for Q -learning using variance estimators in their bonuses and reference-advantage decomposition for variance reduction. We develop a novel error decomposition framework to prove gap-dependent regret bounds of UCB-Advantage and Q-EarlySettled-Advantage that are logarithmic in T and improve upon existing ones for Q -learning algorithms. Moreover, we establish the gap-dependent bound for the policy switching cost of UCB-Advantage and improve that under the worst-case MDPs. To our knowledge, this paper presents the first gap-dependent regret analysis for Q -learning using variance estimators and reference-advantage decomposition and also provides the first gap-dependent analysis on policy switching cost for Q -learning.

1 INTRODUCTION

Reinforcement Learning (RL) (Sutton & Barto, 2018) is a subfield of machine learning focused on sequential decision-making. Often modeled as a Markov Decision Process (MDP), RL tries to obtain an optimal policy through sequential interactions with the environment. It finds applications in various fields, such as games (Silver et al., 2016; 2017; 2018; Vinyals et al., 2019), robotics (Kober et al., 2013; Gu et al., 2017), and autonomous driving (Yurtsever et al., 2020).

In this paper, we focus on the on-policy RL tailored for episodic tabular MDPs with inhomogeneous transition kernels. Specifically, the agent interacts with an episodic MDP consisting of S states, A actions, and H steps per episode. The regret information bound for any MDP above and any learning algorithm with K episodes is $O(\sqrt{H^2SAT})$ where $T = KH$ denotes the total number of steps (Jin et al., 2018). Multiple RL algorithms in the literature (e.g. Zhang et al. (2020); Li et al. (2021); Zhang et al. (2024)) have reached a near-optimal \sqrt{T} -type regret that matches the information bound up to logarithmic factors, which acts as a worst-case guarantee.

In practice, RL algorithms often perform better than their worst-case guarantees, as such guarantees can be significantly improved under MDPs with benign structures (Zanette & Brunskill, 2019). This motivates the problem-dependent analysis for algorithms that exploit the benign MDPs (e.g., Wagenmaker et al. (2022a); Zhou et al. (2023); Zhang et al. (2024)). One of the benign structures is based on the dependency on the positive suboptimality gap: for every state, the best action outperforms others by a margin. It is important because nearly all non-degenerate environments with finite action sets satisfy some sub-optimality gap conditions (Yang et al., 2021). Recently, Simchowitz & Jamieson (2019) proved the log T -type regret if there exists a strictly positive suboptimality gap. Since then, the gap-dependent regret analysis has been widely studied, for example, Dann et al. (2021); Yang et al. (2021); Xu et al. (2021); Wang et al. (2022); He et al. (2021), etc.

Model-free RL algorithms, the focus of this paper, are also called Q -learning algorithms and directly learn the optimal action value function (Q -function) and state value function (V -function) to optimize the policy. It is widely used in practice due to its easy implementation (Jin et al., 2018) and the lower memory requirement that scales linearly in S while that for model-based algorithms scales quadratically. However, the literature on gap-dependent analysis for Q -learning is quite sparse. Yang et al. (2021) studied the gap-dependent regret of the Q-Hoeffding algorithm (Jin et al., 2018), the first model-free algorithm with a worst-case \sqrt{T} -type regret in the literature, and presented the first log T -type regret bound for model-free algorithms:

$$O\left(\frac{H^6 SA \log(SAT)}{\Delta_{\min}}\right). \quad (1)$$

where Δ_{\min} is defined as the minimum nonzero suboptimality gap for all the state-action-step triples.

Xu et al. (2021) proposed the multi-step bootstrapping algorithm and showed the same dependency on the minimum gap as Yang et al. (2021). Both papers used the simple Hoeffding-type bonuses for explorations in the algorithm design. However, their analysis frameworks based on Hoeffding-type bonuses cannot be directly applied to study two important Q -learning algorithms that improve the regrets of Jin et al. (2018) and achieve the almost optimal worst-case regret: UCB-Advantage (Zhang et al., 2020) and Q-EarlySettled-Advantage (Li et al., 2021). In particular, UCB-Advantage and Q-EarlySettled-Advantage use variance estimators in their bonuses and reference-advantage decomposition for variance reduction. It remains an important open question whether such techniques can improve gap-dependent regret:

Is it possible to establish a potentially improved gap-dependent regret bound for Q -learning using variance estimators in the bonuses and reference-advantage decomposition?

This is a challenging task due to several non-trivial difficulties. In particular, bounding the weighted sum of the errors of the estimated Q -functions is necessary to establish the gap-dependent regret bounds for UCB-Advantage and Q-EarlySettled-Advantage, which is very difficult as it involves the estimated reference and advantage functions and the bonuses that include variance estimators for both functions. However, the analysis framework of Xu et al. (2021) for their non-optimism algorithm cannot bound the weighted sum of such errors, and the analysis frameworks in all optimism-based model-free algorithms including Jin et al. (2018); Zhang et al. (2020); Li et al. (2021); Yang et al. (2021) can only bound the weighted sum under the simple Hoeffding-type bonus.

Besides the regret, the policy switching cost is also an important evaluation criterion for on-policy RL, especially in applications with restrictions on policy switching such as compiler optimization (Ashouri et al., 2018), hardware placements (Mirhoseini et al., 2017), database optimization (Krishnan et al., 2018), and material discovery (Nguyen et al., 2019). Under the worst-case MDPs, Bai et al. (2019) modified the algorithms in Jin et al. (2018) to reach a switching cost of $O(H^3 SA \log T)$, and UCB-Advantage (Zhang et al., 2020) reached an improved switching cost of $O(H^2 SA \log T)$ due to the stage design in Q -function update, both improving upon the cost of $\Theta(K)$ for regular Q -learning algorithms (e.g. Jin et al. (2018)). To our knowledge, none of existing works study gap-dependent switching costs for Q -learning algorithms, which remains open.

Summary of our contributions. In this paper, we give an affirmative answer to the open questions above by establishing gap-dependent regret bound for UCB-Advantage (Zhang et al., 2020) and Q-EarlySettled-Advantage (Li et al., 2021) as well as a gap-dependent policy switching cost for UCB-Advantage. For Q -learning, this paper provides the first gap-dependent regret analysis with both variance estimators and variance reduction and the first gap-dependent policy switching cost.

Our detailed contributions are summarized as follows.

- **Improved Gap-Dependent Regret.** Denote $\mathbb{Q}^* \in [0, H^2]$ as the *maximum conditional variance* for the MDP and $\beta \in (0, H]$ as the hyper-parameter to settle the reference function. We prove that UCB-Advantage guarantees a gap-dependent expected regret of

$$O\left(\frac{(\mathbb{Q}^* + \beta^2 H) H^3 SA \log(SAT)}{\Delta_{\min}} + \frac{H^8 S^2 A \log(SAT) \log(T)}{\beta^2}\right), \quad (2)$$

108 and Q-EarlySettled-Advantage guarantees a gap-dependent expected regret of
 109

$$110 \quad O\left(\frac{(\mathbb{Q}^* + \beta^2 H) H^3 S A \log(SAT)}{\Delta_{\min}} + \frac{H^7 S A \log^2(SAT)}{\beta}\right). \quad (3)$$

111
 112

113 These results are logarithmic in T and better than the worst-case \sqrt{T} -type regret in
 114 Zhang et al. (2020); Li et al. (2021). They also have a common gap-dependent term
 115 $\tilde{O}((\mathbb{Q}^* + \beta^2 H) H^3 S A / \Delta_{\min})$ where $\tilde{O}(\cdot)$ hides logarithmic factors. The other term in either
 116 Equation (2) or Equation (3) is gap-free. Our result is also better than Equation (1) for Yang
 117 et al. (2021); Xu et al. (2021) in the following ways. (a) Under the worst-case $\mathbb{Q}^* = \Theta(H^2)$
 118 and setting $\beta = O(1/\sqrt{H})$ as in Zhang et al. (2020) or $\beta = O(1)$ as in Li et al. (2021),
 119 $\tilde{O}((\mathbb{Q}^* + \beta^2 H) H^3 S A / \Delta_{\min})$ becomes $\tilde{O}(H^5 S A / \Delta_{\min})$, which is better than Equation (1) by
 120 a factor of H . (b) Under the best variance $\mathbb{Q}^* = 0$ which will happen when the MDP is deter-
 121 ministic, our regret in Equation (3) can linearly depend on $\tilde{O}(\Delta_{\min}^{-1/3})$, which is intrinsically better
 122 than the dependency on Δ_{\min}^{-1} in Equation (1). (c) Since our gap-free terms also logarithmically
 123 depend on T , they are smaller than Equation (1) when Δ_{\min} is sufficiently small.

124
 125 • **Gap-Dependent Policy Switching Cost.** We can prove that for any $\delta \in (0, 1)$, with probability
 126 at least $1 - \delta$, the policy switching cost for UCB-Advantage is at most

$$127 \quad O\left(H|D_{\text{opt}}| \log\left(\frac{T}{H|D_{\text{opt}}|} + 1\right) + H|D_{\text{opt}}^c| \log\left(\frac{H^4 S A^{\frac{1}{2}} \log(\frac{SAT}{\delta})}{\beta \sqrt{|D_{\text{opt}}^c| \Delta_{\min}}}\right)\right). \quad (4)$$

128
 129
 130

131 Here, D_{opt} is a subset of all state-action-step triples and represents all triples such that the action is
 132 optimal. D_{opt}^c is its complement, and $|\cdot|$ gives the cardinality of the set. Next, we compare Equation
 133 (4) with the worst-case costs of $O(H^3 S A \log T)$ in Bai et al. (2019) and $O(H^2 S A \log T)$
 134 in Zhang et al. (2020). Since $|D_{\text{opt}}| < H S A$ for non-degenerate MDPs, our first term in Equation
 135 (4) is better than the worst-case costs. Specifically, when each state has a unique optimal
 136 action so that $|D_{\text{opt}}| = H S$, it implies the improvement by removing a factor of A compared with
 137 $O(H^2 S A \log T)$. This improvement is significant in applications with a large action space (e.g.
 138 recommender systems (Covington et al., 2016) and text-based games (Bellemare et al., 2013)). For
 139 the second term where $|D_{\text{opt}}^c| < H S A$ in Equation (4), we also improve $\log T$ to $\log \log T$, and
 140 the significance of such improvement is pointed out by Qiao et al. (2022); Zhang et al. (2022b).

141 • **Technical Novelty and Contributions.**

142 For gap-dependent regret analysis, we develop an error decomposition framework that separates
 143 errors in reference estimations, advantage estimations, and reference settling. This helps bound
 144 the weighted sums mentioned above. We creatively handle the separated terms in the following
 145 way. (a) We relate the empirical errors and the bonus for reference estimations to \mathbb{Q}^* to avoid
 146 using their upper bounds $\Theta(H^2)$. This leverages the variance estimators. (b) When trying to
 147 bound the errors in reference and advantage estimations, we tackle the non-martingale difficulty,
 148 originating from the settled reference functions that depend on the whole learning process, with
 149 our novel surrogate reference functions so that the empirical estimations become martingale sums.
 150 To the best of our knowledge, we are the first to construct martingale surrogates in the literature
 151 for Q -learning using reference-advantage decomposition.

152 For the gap-dependent policy switching cost, we explore the unbalanced number of visits to states
 153 paired with optimal or suboptimal actions, which leads to the two terms in Equation (4).

154 **2 PRELIMINARIES**

155

156 Throughout this paper, we assume that $0/0 = 0$. For any $C \in \mathbb{N}$, we use $[C]$ to denote the set
 157 $\{1, 2, \dots, C\}$. We use $\mathbb{I}[x]$ to denote the indicator function, which equals 1 when the event x is true
 158 and 0 otherwise.

159

160 **Tabular episodic Markov decision process (MDP).** A tabular episodic MDP is denoted as $\mathcal{M} :=$
 161 $(\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r)$, where \mathcal{S} is the set of states with $|\mathcal{S}| = S$, \mathcal{A} is the set of actions with $|\mathcal{A}| = A$, H
 162 is the number of steps in each episode, $\mathbb{P} := \{\mathbb{P}_h\}_{h=1}^H$ is the transition kernel so that $\mathbb{P}_h(\cdot | s, a)$

162 characterizes the distribution over the next state given the state action pair (s, a) at step h , and
 163 $r := \{r_h\}_{h=1}^H$ are the deterministic reward functions with $r_h(s, a) \in [0, 1]$.
 164

165 In each episode, an initial state s_1 is selected arbitrarily by an adversary. Then, at each step $h \in [H]$,
 166 an agent observes a state $s_h \in \mathcal{S}$, picks an action $a_h \in \mathcal{A}$, receives the reward $r_h = r_h(s_h, a_h)$
 167 and then transits to the next state s_{h+1} . The episode ends when an absorbing state s_{H+1} is reached.
 168 Later on, for ease of presentation, when we describe s, a, h along with “any, each, all” or “ \forall ”, we
 169 will omit the sets $\mathcal{S}, \mathcal{A}, [H]$. We denote $\mathbb{P}_{s,a,h}f = \mathbb{E}_{s_{h+1} \sim \mathbb{P}_h(\cdot|s,a)}(f(s_{h+1})|s_h = s, a_h = a)$,
 170 $\mathbb{V}_{s,a,h}f = \mathbb{P}_{s,a,h}f^2 - (\mathbb{P}_{s,a,h}f)^2$ and $\mathbb{1}_s f = f(s), \forall(s, a, h)$ for any function $f : \mathcal{S} \rightarrow \mathbb{R}$.

171 **Policies, state value functions, and action value functions.** A policy π is a collection of H functions
 172 $\{\pi_h : \mathcal{S} \rightarrow \Delta^{\mathcal{A}}\}_{h \in [H]}$, where $\Delta^{\mathcal{A}}$ is the set of probability distributions over \mathcal{A} . A policy is
 173 deterministic if for any $s \in \mathcal{S}$, $\pi_h(s)$ concentrates all the probability mass on an action $a \in \mathcal{A}$. In
 174 this case, we denote $\pi_h(s) = a$. We use $V_h^\pi : \mathcal{S} \rightarrow \mathbb{R}$ to denote the state value function at step h
 175 under policy π . Mathematically, $V_h^\pi(s) := \sum_{h'=h}^H \mathbb{E}_{(s_{h'}, a_{h'}) \sim (\mathbb{P}, \pi)}[r_{h'}(s_{h'}, a_{h'}) | s_h = s]$. We also
 176 use $Q_h^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ to denote the action value function at step h , i.e., $Q_h^\pi(s, a) := r_h(s, a) +$
 177 $\sum_{h'=h+1}^H \mathbb{E}_{(s_{h'}, a_{h'}) \sim (\mathbb{P}, \pi)}[r_{h'}(s_{h'}, a_{h'}) | s_h = s, a_h = a]$. Azar et al. (2017) proved that there al-
 178 ways exists an optimal policy π^* that achieves the optimal value $V_h^*(s) = \sup_\pi V_h^\pi(s) = V_h^{\pi^*}(s)$
 179 for all $s \in \mathcal{S}$ and $h \in [H]$. The Bellman equation and the Bellman optimality equation are

$$\begin{cases} V_h^\pi(s) = \mathbb{E}_{a' \sim \pi_h(s)}[Q_h^\pi(s, a')] \\ Q_h^\pi(s, a) := r_h(s, a) + \mathbb{P}_{s,a,h}V_{h+1}^\pi \\ V_{H+1}^\pi(s) = 0, \forall(s, a, h) \end{cases} \quad \begin{cases} V_h^*(s) = \max_{a' \in \mathcal{A}} Q_h^*(s, a') \\ Q_h^*(s, a) := r_h(s, a) + \mathbb{P}_{s,a,h}V_{h+1}^* \\ V_{H+1}^*(s) = 0, \forall(s, a, h). \end{cases} \quad (5)$$

180 For any learning problem with K episodes, let π^k be the policy adopted in the k -th episode,
 181 and s_1^k be the corresponding initial state. The regret over $T = HK$ steps is $\text{Regret}(T) :=$
 182 $\sum_{k=1}^K (V_1^*(s_1^k) - V_1^{\pi^k}(s_1^k))$. Later, when we mention the episode index k with “any, each, all”
 183 or “ \forall ”, we will omit the set $[K]$.

184 **Suboptimality Gap.** For any given MDP, we can provide the following formal definition.

185 **Definition 2.1.** For any (s, a, h) , the suboptimality gap is defined as $\Delta_h(s, a) := V_h^*(s) - Q_h^*(s, a)$.

186 Equation (5) implies that $\Delta_h(s, a) \geq 0, \forall(s, a, h)$. Then it is natural to define the minimum gap,
 187 which is the minimum non-zero suboptimality gap with regard to all (s, a, h) .

188 **Definition 2.2.** We define the **minimum gap** as $\Delta_{\min} := \inf\{\Delta_h(s, a) : \Delta_h(s, a) > 0, (s, a, h) \in$
 189 $\mathcal{S} \times \mathcal{A} \times [H]\}$.

190 We remark that if $\{\Delta_h(s, a) : \Delta_h(s, a) > 0, (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]\} = \phi$, then all actions are
 191 optimal, leading to a degenerate MDP. Therefore, we assume that the set is nonempty and $\Delta_{\min} > 0$.
 192 Definitions 2.1 and 2.2 and the non-degeneration are standard in the literature on gap-dependent
 193 analysis (e.g. Simchowitz & Jamieson (2019); Xu et al. (2020)).

194 **Maximum Conditional Variance.** This quantity is formally defined as follows.

195 **Definition 2.3.** We define the **maximum conditional variance** as $\mathbb{Q}^* := \max_{s,a,h} \{\mathbb{V}_{s,a,h}(V_{h+1}^*)\}$.

196 Under our MDP with deterministic reward, Definition 2.3 coincides with that in (Zanette & Brun-
 197 skill, 2019) which performed variance-dependent regret analysis.

198 **Policy Switching Cost.** We provide the following definition for any algorithm with $K > 1$ episodes.

199 **Definition 2.4.** The policy switching cost for K episodes is defined as $N_{\text{switch}} :=$
 200 $\sum_{k=1}^{K-1} \tilde{N}_{\text{switch}}(\pi^{k+1}, \pi^k)$. Here, the $\tilde{N}_{\text{switch}}(\pi, \pi') := \sum_{s \in \mathcal{S}} \sum_{h=1}^H \mathbb{I}[\pi_h(s) \neq \pi'_h(s)]$ represents
 201 the local switching cost for any policies π and π' .

202 This definition is also used in Bai et al. (2019) and Zhang et al. (2020).

212 3 MAIN RESULTS

213 This section presents the gap-dependent regret for UCB-Advantage and Q-EarlySettled-Advantage
 214 in Subsection 3.1 and the gap-dependent policy switching cost for UCB-Advantage in Subsection
 215 3.3. We highlight a new technical tool for the gap-dependent regret bound in Subsection 3.2.

216 3.1 GAP-DEPENDENT REGRETS
217

218 UCB-Advantage (Zhang et al., 2020) is the first model-free algorithm that reaches an almost optimality worst-case regret, which is also reached by Q-EarlySettled-Advantage (Li et al., 2021). Both 219 algorithms are optimism-based, use upper confidence bounds (UCB) for exploration, and employ 220 variance estimators and reference-advantage decomposition. UCB-Advantage settles the reference 221 function at each (s, h) by comparing the number of visits to a threshold that relies on a hyper- 222 parameter $\beta \in (0, H]$. For readers' convenience, we provide UCB-Advantage without any modification 223 in Algorithm 1 of Appendix B.1.

225 Theorem 3.1 provides the expected regret upper bound of UCB-Advantage.

226 **Theorem 3.1.** *For UCB-Advantage (Algorithm 1 in Appendix B.1) with $\beta \in (0, H]$, $\mathbb{E}[\text{Regret}(T)]$ 227 is upper bounded by Equation (2).*

228 Q-EarlySettled-Advantage improved the burn-in cost of Zhang et al. (2020) for reaching the almost- 229 optimal worst-case regret by using both estimated upper and lower confidence bounds for V_h^* to 230 settle the reference function. In this paper, we slightly modify its reference settling condition. At 231 the end of k -th episode, for any (s, h) , the algorithm holds $V_h^{k+1}(s), V_h^{\text{LCB}, k+1}(s)$, the estimated 232 upper and lower bounds for $V_h^*(s)$, respectively. When $|V_h^{k+1}(s) - V_h^{\text{LCB}, k+1}(s)| \leq \beta$ holds for 233 the first time, it settles the reference function value $V_h^R(s)$ as $V_h^{k+1}(s)$. Li et al. (2021) set $\beta = 1$ 234 for worst-case MDPs. Our paper treats β as a hyper-parameter within $(0, H]$ to allow better control 235 over the learning process. Algorithms 2 and 3 provide our refined version. For the rest of this paper, 236 we still call it Q-EarlySettled-Advantage without special notice.

237 Theorem 3.2 provides the expected regret upper bound of Q-EarlySettled-Advantage.

238 **Theorem 3.2.** *For Q-EarlySettled-Advantage (Algorithms 2 and 3 in Appendix D.1) with $\beta \in (0, H]$, $\mathbb{E}[\text{Regret}(T)]$ is upper bounded by Equation (3).*

239 The proof sketch of Theorem 3.2 is presented in Section 4 to explain our technical contributions. The 240 complete proofs of Theorems 3.1 and 3.2 are provided in Appendix B and Appendix D, respectively.

241 Next, we compare the results of both theorems with the worst-case regrets in Zhang et al. (2020); Li 242 et al. (2021) and the gap-dependent regrets in Yang et al. (2021); Xu et al. (2021).

243 **Comparisons with Zhang et al. (2020); Li et al. (2021).** Since the regrets showed in Equations 244 (2) and (3) are logarithmic in T , they are better than the worst-case regret $\tilde{O}(\sqrt{H^2SAT})$ 245 when $T \geq \tilde{\Theta}(\text{poly}(HSA, \Delta_{\min}^{-1}, \beta^{-1}))$ where $\text{poly}(\cdot)$ represents some polynomial. In addition, our 246 results imply new guidance on setting the hyper-parameter β for the gap-dependent regret, which 247 is different from $\beta = 1/\sqrt{H}$ in Zhang et al. (2020) and $\beta = 1$ in Li et al. (2021), respectively. 248 When $\mathbb{Q}^* = 0$ which happens when the MDP is deterministic, if we set $\beta = \tilde{\Theta}(H(S\Delta_{\min})^{1/4})$ for 249 UCB-Advantage and $\beta = \tilde{\Theta}(H\Delta_{\min}^{1/3})$, the gap-dependent regrets will linearly depend on $\Delta_{\min}^{-1/2}$ and 250 $\Delta_{\min}^{-1/3}$, respectively. This provides new guidance on setting β when we have prior knowledge about 251 Δ_{\min} . When $0 < \mathbb{Q}^* \leq H^2$, the best available gap-dependent regret becomes $\tilde{\Theta}(\mathbb{Q}^*H^2SA)$ which 252 holds when $\beta \leq \sqrt{\mathbb{Q}^*/H}$. Knowing that the gap-free terms in Equations (2) and (3) monotonically 253 decrease in β , we will recommend setting $\beta = \tilde{O}(\sqrt{\mathbb{Q}^*/H})$ if prior knowledge on \mathbb{Q}^* is available.

254 **Comparisons with Yang et al. (2021); Xu et al. (2021).** The gap-dependent regret for Yang et al. 255 (2021) is provided in Equation (1). For the multi-step bootstrapping in Xu et al. (2021), their regret 256 bound is given by:

$$257 O \left(\left(\sum_{h=1}^H \sum_{s \in \mathcal{S}} \sum_{a \neq \pi_h^*(s)} \frac{1}{\Delta_h(s, a)} + \frac{|Z_{\text{mul}}|}{\Delta_{\min}} + SA \right) H^5 \log(K) \right), \quad (6)$$

258 where $Z_{\text{mul}} = \{(h, s, a) | \Delta_h(s, a) = 0 \wedge |Z_{\text{opt}}^h(s)| > 1\}$ and $Z_{\text{opt}}^h(s) = \{a | \Delta_h(s, a) = 0\}$. In MDPs 259 where $\Delta_h(s, a) = \Theta(\Delta_{\min})$ for $\Theta(HSA)$ state-action-step triples (e.g. the example in (Xu et al., 260 2021, Theorem 1.3)) or there are $\Theta(A)$ optimal actions for each state-step pair (s, h) , their regret 261 reduces to Equation (1), which is worse than ours.

Next, we compare Equations (2) and (3) with Equation (1). Under the worst-case variance $\mathbb{Q}^* = \Theta(H^2)$ and letting β be $\Theta(1/\sqrt{H})$ or $\Theta(1)$ which are the recommendations in Zhang et al. (2020); Li et al. (2021) respectively for the worst-case MDPs, the common gap-dependent term Equations (2) and (3) becomes $\tilde{O}(H^5SA/\Delta_{\min})$, which is better than Equation (1) by a factor of H . Under the best variance $\mathbb{Q}^* = 0$, the gap-dependent term becomes $\tilde{O}(\beta^2 H^3 SA/\Delta_{\min})$, which is better than Equation (1) for any $\beta \in (0, H]$. In addition, our best possible gap-dependent regret that is sublinear in Δ_{\min}^{-1} is also intrinsically better. Here, we remark that the proof in Yang et al. (2021); Xu et al. (2021) cannot benefit from $\mathbb{Q}^* = 0$ due to their use of Hoeffding-type bonuses.

We also comment on the gap-free terms in Equations (2) and (3). They are dominated by the gap-dependent term as long as $\Delta_{\min} \leq \tilde{O}(\text{poly}((HSA)^{-1}, \beta))$ for some polynomial $\text{poly}(\cdot)$. In addition, the gap-free term in Equation (3) is linear in S , which is better than that for Equation (2) thanks to Q-EarlySettled-Advantage algorithm. It utilizes both upper confidence bounds and lower confidence bounds for V -functions to settle the reference function.

3.2 OUR TECHNICAL TOOL: SURROGATE REFERENCE FUNCTIONS

We develop a new technical tool in the proofs of both Theorems 3.1 and 3.2: the surrogate reference functions. In this subsection, we explain it with the notations in the proof of Theorem 3.2 (Appendix D.1) for Q-EarlySettled-Advantage while all the ideas also apply to UCB-Advantage. A more detailed proof sketch will be provided in the next section. For a comprehensive explanation of Q-EarlySettled-Advantage, we refer readers to Appendix D.1, and for a detailed mathematical explanation of the surrogate function, please see Appendix G.

Before introducing the surrogate reference function, we provide a brief overview of the key steps of Q-EarlySettled-Advantage. Denote the estimated Q -function, the estimated V -function, and the reference function before the start of episode k as $Q_h^k(s, a), V_h^k(s), V_h^{R,k}(s)$ and episode k as $\{(s_h^k, a_h^k)\}_{h=1}^H$. Let $N_h^k(s, a)$ be the number of visits to (s, a, h) before the start of episode k . Let N_h^{k+1} be short for $N_h^{k+1}(s_h^k, a_h^k)$ and k^n be the episode index for the n -th visit to (s_h^k, a_h^k, h) . While remaining unchanged for the unvisited triples, the estimated Q -function is updated on the visited ones:

$$Q_h^{k+1}(s_h^k, a_h^k) = \min\{Q_h^{\text{UCB}, k+1}(s_h^k, a_h^k), Q_h^{\text{R}, k+1}(s_h^k, a_h^k), Q_h^k(s_h^k, a_h^k)\}, h \in [H]. \quad (7)$$

Here, $Q_h^{\text{UCB}, k+1}$ represents the Hoeffding-type estimation similar to Jin et al. (2018), and $Q_h^{\text{R}, k+1}(s_h^k, a_h^k)$ represents the reference-advantage type estimation as follows:

$$Q_h^{\text{R}, k+1}(s_h^k, a_h^k) = r_h^k(s_h^k, a_h^k) + \sum_{n=1}^{N_h^{k+1}} \left(\eta_n^{N_h^{k+1}} (V_{h+1}^{k^n} - V_{h+1}^{R,k^n}) + u_n^{N_h^{k+1}} V_{h+1}^{R,k^n} \right) (s_{h+1}^{k^n}) + \tilde{R}^{h, k+1}. \quad (8)$$

In Equation (8), $V_{h+1}^{k^n} - V_{h+1}^{R,k^n}$ represents the running estimation of the advantage function, and $\{\eta_n^{N_h^{k+1}}\}_{n=1}^{N_h^{k+1}}$ are the corresponding nonnegative weights that sum to 1. $\{u_n^{N_h^{k+1}}\}_{n=1}^{N_h^{k+1}}$ that sum to 1 are nonnegative weights for the reference function. $\tilde{R}^{h, k+1}$ is the cumulative bonus that dominates the variances in the two weighted sums. Next, the estimated V -function and the reference function are also updated. For any (s, h) , when some reference settling condition related to β is triggered at the end of episode k , the reference function will be settled, which means that $V_h^{\text{R}, k'}(s) = V_h^{\text{R}, k+1}(s)$ for any $k' \geq k + 1$. Thus, we call $V_h^{\text{R}, K+1}$, the reference value function after the last episode as the settled reference function. Q-EarlySettled-Advantage guarantees that

$$V_h^k(s) = \max_a Q_h^k(s, a), \pi_h^k(s) = \arg \max_a Q_h^k(s, a), \forall (h, k), \quad (9)$$

$$Q_h^{k+1} \leq Q_h^k \leq H, V_h^{k+1} \leq V_h^k \leq H, V_h^{\text{R}, k+1} \leq V_h^{\text{R}, k} \leq H, V_h^k \leq V_h^{\text{R}, k}, \forall (h, k). \quad (10)$$

Event \mathcal{E}_1 in Lemma D.2 (Lemma 2 in Li et al. (2021)) also claims that with high probability,

$$Q_h^k \geq Q_h^*, V_h^{\text{R}, k} \geq V_h^k \geq V_h^*, \forall (h, k). \quad (11)$$

Equations (9) and (11) indicate that Q-EarlySettled-Advantage is an optimism-based method that updates the policy according to an upper bound of Q_h^* .

324 Next, we introduce our surrogate reference functions $\hat{V}_h^{\text{R},k}$. They are defined as follows:
 325

$$\hat{V}_h^{\text{R},k}(s) := \max\{V_h^*(s), \min\{V_h^*(s) + \beta, V_h^{\text{R},k}(s)\}\}, \forall(s, h, k). \quad (12)$$

328 We use the word “surrogate” because the algorithm does not rely on or learn it, and $\hat{V}_h^{\text{R},k}$ differs
 329 from the actual settled reference function $V_h^{\text{R},K+1}$. $\hat{V}_h^{\text{R},k}$ is determined before episode k . In addition,
 330 Equation (11) implies that

$$V_h^* \leq \hat{V}_h^{\text{R},k} = \min\{V_h^* + \beta, V_h^{\text{R},k}\}, \forall(k, h), \quad (13)$$

333 and Lemma D.4 in Appendix D.5.2 shows that with high probability, $\hat{V}_h^{\text{R},k}(s)$ coincides with the
 334 settled reference value $V_h^{\text{R},K+1}(s)$ after the settling condition is triggered.
 335

336 Next, we discuss the usage of $\hat{V}_h^{\text{R},k}$ in our error decompositions. Our proof relies on relating the
 337 regret to multiple groups of estimation error sums that take the form $\sum_{k=1}^K \omega_{h,k}^{(i)}(Q_h^k - Q_h^*)(s_h^k, a_h^k)$.
 338 Here $\{\omega_{h,k}^{(i)}\}_k$ are nonnegative weights and i represents the group. Bounding the weighted sum
 339 via controlling each individual $Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)$ by recursion on h is a common tech-
 340 nique for model-free optimism-based algorithms, and it is also used by all of Yang et al. (2021);
 341 Zhang et al. (2020); Li et al. (2021). Yang et al. (2021) used it on gap-dependent regret anal-
 342 ysis and Zhang et al. (2020); Li et al. (2021) used it to control the reference setting errors
 343 $\sum_{k=1}^K (V_h^{\text{R},k+1}(s_h^k) - V_h^{\text{R},K+1}(s_h^k))$. However, their techniques are only limited to the Hoeffding-
 344 type update where the errors generated in the recursion take the simple form of $\tilde{O}(\sqrt{H^3/N_h^k})$ where
 345 N_h^k is short for $N_h^k(s_h^k, a_h^k)$. When analyzing the reference-advantage type update, we will face a
 346 complicated error (see Equation (15) in the proof sketch) that involves reference estimations, advan-
 347 tage estimations, and bonuses with variance estimators. See Appendix G for the details.
 348

349 Motivated by the structure of reference-advantage decomposition, we decompose our error into
 350 four parts: $\mathcal{G}_1 := \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (\mathbb{P}_{s_h^k, a_h^k, h} - \mathbb{1}_{s_{h+1}^{k^n}})(\hat{V}_{h+1}^{\text{R},k^n} - V_{h+1}^*)$, $\mathcal{G}_2 := \sum_{n=1}^{N_h^k} u_n^{N_h^k} (\mathbb{1}_{s_{h+1}^{k^n}} -$
 351 $\mathbb{P}_{s_h^k, a_h^k, h})\hat{V}_{h+1}^{\text{R},k^n}$, $\mathcal{G}_3 := \sum_{n=1}^{N_h^k} (u_n^{N_h^k} - \eta_n^{N_h^k})\mathbb{P}_{s_h^k, a_h^k, h}\hat{V}_{h+1}^{\text{R},k^n} + \sum_{n=1}^{N_h^k} u_n^{N_h^k} (V_{h+1}^{\text{R},k^n} - \hat{V}_{h+1}^{\text{R},k^n})(s_{h+1}^{k^n})$,
 352 the bonus term \mathcal{G}_4 and a negative term $\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (\hat{V}_{h+1}^{\text{R},k^n} - V_{h+1}^{\text{R},k^n})(s_{h+1}^{k^n})$. The first three terms
 353 correspond to advantage estimation error, reference estimation error, and reference settling error, re-
 354 spectively. Here, we creatively use the surrogate $\hat{V}_{h+1}^{\text{R},k}$ as it is determined before the start of episode
 355 k . Thus, $\mathcal{G}_1, \mathcal{G}_2$ are martingale sums and can be controlled by concentration inequalities. \mathcal{G}_3 corre-
 356 sponds to the reference settling error and can also be well-controlled given the settling conditions
 357 and properties of $\hat{V}_h^{\text{R},k}(s)$. \mathcal{G}_4 is controlled using the same idea of bounding $\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3$. $\hat{V}_{h+1}^{\text{R},k}$ is
 358 crucial to this process and cannot be replaced by the actual settled reference function $V_{h+1}^{\text{R},K+1}$ used
 359 in Zhang et al. (2020); Li et al. (2021). This is because $V_{h+1}^{\text{R},K+1}$ depends on the whole learning
 360 process and causes a non-martingale issue in controlling $\mathcal{G}_1, \mathcal{G}_2$. To the best of our knowledge, we
 361 are the first to introduce the novel construction of reference surrogates for reference-advantage de-
 362 composition, which is of independent interest for future research on off-policy and offline methods.
 363

3.3 GAP-DEPENDENT POLICY SWITCHING COST FOR UCB-ADVANTAGE

367 Different from Q-EarlySettled-Advantage, UCB-Advantage uses the stage design for updating the
 368 estimated Q -function. For each (s, a, h) , Zhang et al. (2020) divided the visits into consecutive
 369 stages with the stage size increasing exponentially. It updates the estimated Q -function only at the
 370 end of each stage so that the policy switches infrequently. Theorem 3.3 provides the policy switching
 371 cost for UCB-Advantage, and the proof is provided in Appendix C.

372 **Theorem 3.3.** *For UCB-Advantage (Algorithm 1 in Appendix B.1) with $\beta \in (0, H]$ and any $\delta \in$
 373 $(0, 1)$, with probability at least $1 - \delta$, N_{switch} is upper bounded by Equation (4). Here, $D_{\text{opt}} =$
 374 $\{(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H] | a = \pi_h^*(s)\}$, and $D_{\text{opt}}^c = (\mathcal{S} \times \mathcal{A} \times [H]) \setminus D_{\text{opt}}$.*

376 **Comparisons with existing works.** The first term in Equation (4) logarithmically depends on T
 377 and the second one logarithmically depends on $1/\Delta_{\min}$ and $\log T$. Next, we compare our result with
 $O(H^2 S A \log T)$ in Zhang et al. (2020), which is the best available switching cost for model-free

378 methods in the literature. For the first term in Equation (4), knowing that $|D_{\text{opt}}| < HSA$ for all
379 non-degenerated MDPs where there exists at least one state such that not all actions are optimal, the
380 coefficient is better than Equation (4). Specifically, if each state has a unique optimal action so that
381 $|D_{\text{opt}}| = SH$, Equation (4) becomes $O\left(H^2 S \log\left(\frac{T}{H^2 S} + 1\right) + H^2 SA \log\left(\frac{H^{\frac{7}{2}} S^{\frac{1}{2}} \log(\frac{SAT}{\delta})}{\beta \Delta_{\min}}\right)\right)$
382 where coefficient in the first term is Zhang et al. (2020) by a factor of A .

383 For the second term in Equation (4), when the total steps are sufficiently large such that $T = \tilde{\Omega}(\text{poly}(SAH, (\beta\Delta_{\min})^{-1}))$ for some polynomial $\text{poly}(\cdot)$, it is also better than $O(H^2 SA \log T)$.
384

385 **Key Ideas of the Proof.** The proof of Theorem 2 in Zhang et al. (2020) implies $N_{\text{switch}} \leq$
386 $\sum_{s,a,h} 4H \log\left(\frac{N_h^{K+1}(s,a)}{2H} + 1\right)$, where $N_h^{K+1}(s,a)$ is upper bounded by the total number of
387 visits to (s,a,h) . Under their worst-case MDP and noticing that $\sum_{s,a,h} N_h^{K+1}(s,a) \leq T$,
388 Zhang et al. (2020) further proved their bound $O(H^2 SA \log T)$ by applying Jensen’s inequality.
389 In our gap-dependent analysis, Equation (76) in Appendix C shows that with high probability,
390 $\sum_{(s,a,h) \in D_{\text{opt}}^c} N_h^{K+1}(s,a) \leq \tilde{O}\left(\frac{H^6 SA}{\Delta_{\min}} + \frac{H^8 S^2 A}{\beta^2}\right)$, which is much smaller than T when T is suf-
391 ficiently large. This implies the discrepancy among the number of visits to state-action-step triples
392 with optimal or suboptimal actions. Accordingly, we prove the bound in Equation (4) by using
393 Jensen’s inequality separately for triples with optimal or suboptimal actions.
394

395 4 PROOF SKETCH

400 This section provides a proof sketch to outline the key steps for proving Theorem 3.2 on the gap-
401 dependent regret of Q-EarlySettled-Advantage and explain our technical contributions. The key
402 steps for proving Theorem 3.1 are similar except for different bounds on reference settling error and
403 gap-free regret terms. For space consideration, their complete proofs are presented in the appendix.

404 **Notations.** First, we show the weights used in the algorithm. Let $\eta_n := \frac{H+1}{H+n}$. For $N \in \mathbb{N}_+$,
405 denote $\eta_0^0 := 1$ and $\eta_0^N := \prod_{i=1}^N (1 - \eta_i)$. For integers $1 \leq n \leq N$, we also denote $\eta_n^N :=$
406 $\eta_n \prod_{i=n+1}^N (1 - \eta_i)$, and $u_n^N = \sum_{i=n}^N \eta_i^N / i$. When $N > 0$, they satisfy $1 - \eta_0^N = \sum_{n=1}^N \eta_n^N =$
407 $\sum_{n=1}^N u_n^N$. For simplicity later, we use the notations $\hat{\mathbb{E}}_{h,k}^{\text{ref}} f := \sum_{n=1}^{N_h^k} u_n^{N_h^k} f(s_{h+1}^{k^n})$ and $\hat{\mathbb{E}}_{h,k}^{\text{ref}} f^{k^n} :=$
408 $\sum_{n=1}^{N_h^k} u_n^{N_h^k} f^{k^n}(s_{h+1}^{k^n})$ for any functions $f : \mathcal{S} \rightarrow \mathbb{R}$ and $f^k : \mathcal{S} \rightarrow \mathbb{R}$ with $k \in \mathbb{N}_+$, respec-
409 tively. Similarly, we denote $\hat{\mathbb{E}}_{h,k}^{\text{adv}} f := \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} f(s_{h+1}^{k^n})$ and $\hat{\mathbb{E}}_{h,k}^{\text{adv}} f^{k^n} := \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} f^{k^n}(s_{h+1}^{k^n})$.
410 We also denote $\mathbb{P}_{h,k}^{\text{ref}} f = \sum_{n=1}^{N_h^k} u_n^{N_h^k} \mathbb{P}_{s_h^k, a_h^k, h} f$, $\mathbb{P}_{h,k}^{\text{ref}} f^{k^n} = \sum_{n=1}^{N_h^k} u_n^{N_h^k} \mathbb{P}_{s_h^k, a_h^k, h} f^{k^n}$, $\mathbb{P}_{h,k}^{\text{adv}} f =$
411 $\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \mathbb{P}_{s_h^k, a_h^k, h} f$ and $\mathbb{P}_{h,k}^{\text{adv}} f^{k^n} = \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \mathbb{P}_{s_h^k, a_h^k, h} f^{k^n}$.
412

413 In what follows, we present the proof sketch of Theorem 3.2.

414 **Step 1: Bounding $Q_h^k - Q_h^*$ via decomposition and the surrogate reference function.** The update
415 of the estimated Q -function in Equations (7) and (8) guarantees that

$$416 Q_h^k(s_h^k, a_h^k) \leq \eta_0^{N_h^k} H + r_h(s_h^k, a_h^k) + \hat{\mathbb{E}}_{h,k}^{\text{adv}}(V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) + \hat{\mathbb{E}}_{h,k}^{\text{ref}} V_{h+1}^{\text{R}, k^n} + R^{h,k}. \quad (14)$$

417 Here, $R^{h,k}$ is the cumulative bonus provided in Equation (91) in Appendix D. Together with
418 $Q_h^*(s_h^k, a_h^k) \geq r_h(s_h^k, a_h^k) + (1 - \eta_0^{N_h^k}) \mathbb{P}_{s_h^k, a_h^k, h} V_{h+1}^*$ by Equation (5) and $\hat{\mathbb{E}}_{h,k}^{\text{adv}}(V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \leq$
419 $\hat{\mathbb{E}}_{h,k}^{\text{adv}}(V_{h+1}^{k^n} - \hat{V}_{h+1}^{\text{R}, k^n})$ implied by Equation (13), we have $(Q_h^k - Q_h^*)(s_h^k, a_h^k) \leq \eta_0^{N_h^k} H + R^{h,k} +$
420 $\hat{\mathbb{E}}_{h,k}^{\text{adv}}(V_{h+1}^{k^n} - \hat{V}_{h+1}^{\text{R}, k^n}) + \hat{\mathbb{E}}_{h,k}^{\text{ref}} V_{h+1}^{\text{R}, k^n} - \mathbb{P}_{h,k}^{\text{adv}} V_{h+1}^* =: G_h^k$. Denote $\hat{V}_h^{\text{adv}, k} = \hat{V}_h^{\text{R}, k} - V_h^*$, then:

$$421 G_h^k = \hat{\mathbb{E}}_{h,k}^{\text{adv}}(V_{h+1}^{k^n} - V_{h+1}^*) + (\mathbb{P}_{h,k}^{\text{adv}} - \hat{\mathbb{E}}_{h,k}^{\text{adv}}) \hat{V}_{h+1}^{\text{adv}, k^n} + (\hat{\mathbb{E}}_{h,k}^{\text{ref}} - \mathbb{P}_{h,k}^{\text{ref}}) \hat{V}_{h+1}^{\text{R}, k^n} + R^{h,k} + R_{\text{else}, 0}^{h,k}. \quad (15)$$

422 Here, $R_{\text{else}, 0}^{h,k} = H \eta_0^{N_h^k} + \hat{\mathbb{E}}_{h,k}^{\text{ref}}(V_{h+1}^{\text{R}, k^n} - \hat{V}_{h+1}^{\text{R}, k^n}) + (\mathbb{P}_{h,k}^{\text{ref}} \hat{V}_{h+1}^{\text{R}, k^n} - \mathbb{P}_{h,k}^{\text{adv}} \hat{V}_{h+1}^{\text{R}, k^n})$. Equation (88) and
423 Equation (89) in Appendix D.3.1 show that for all (k, h) simultaneously, with high probability,

$$424 (\mathbb{P}_{h,k}^{\text{adv}} - \hat{\mathbb{E}}_{h,k}^{\text{adv}}) \hat{V}_{h+1}^{\text{adv}, k^n} \leq \tilde{O}\left(\sqrt{\frac{H\beta^2}{N_h^k}}\right), \quad (\hat{\mathbb{E}}_{h,k}^{\text{ref}} - \mathbb{P}_{h,k}^{\text{ref}}) \hat{V}_{h+1}^{\text{R}, k^n} \leq \tilde{O}\left(\sqrt{\frac{\mathbb{Q}^* + \beta^2}{N_h^k}} + \frac{H}{N_h^k}\right). \quad (16)$$

432 Equation (16) corresponds to controlling $\mathcal{G}_1, \mathcal{G}_2$ discussed in Section 3.2 and holds because our
 433 surrogate reference function adapts to the learning process. To bound the bonus $R^{h,k}$, we use $\hat{V}_h^{\text{R},k}$
 434 in Appendix D.3.2. Equation (96) shows that for all (k, h) simultaneously, with high probability
 435

$$436 R^{h,k} \leq \tilde{O} \left(\sqrt{(\mathbb{Q}^* + \beta^2 H)/N_h^k} + H^2/(N_h^k)^{3/4} + \sqrt{H\Psi_h^k/N_h^k} \right). \quad (17)$$

438 where $\Psi_h^k = \sum_{n=1}^{N_h^k} (V_{h+1}^{\text{R},k^n} - \hat{V}_{h+1}^{\text{R},k^n})(s_{h+1}^{k^n})$. Equations (15) to (17) imply
 439

$$440 (Q_h^k - Q_h^*)(s_h^k, a_h^k) \leq \hat{\mathbb{E}}_{h,k}^{\text{adv}}(V_{h+1}^{k^n} - V_{h+1}^*) + \tilde{O} \left(\sqrt{(\mathbb{Q}^* + H\beta^2)/N_h^k} + H^2(N_h^k)^{-3/4} \right) + R_{\text{else}}^{h,k} \quad (18)$$

442 where $R_{\text{else}}^{h,k} = \tilde{O} \left(\eta_0^{N_h^k} H + \hat{\mathbb{E}}_{h,k}^{\text{ref}}(V_{h+1}^{\text{R},k^n} - \hat{V}_{h+1}^{\text{R},k^n}) + (\mathbb{P}_{h,k}^{\text{ref}} \hat{V}_{h+1}^{\text{R},k^n} - \mathbb{P}_{h,k}^{\text{adv}} \hat{V}_{h+1}^{\text{R},k^n}) + (\sqrt{H\Psi_h^k} + H)/N_h^k \right)$.
 443

444 *Remark 1:* We can show \mathbb{Q}^* in Equation (17) instead of its upper bound $\Theta(H^2)$ thanks to the
 445 variance estimator (line 16 of Algorithm 2 in Appendix D.1) used in Q-EarlySettled-Advantage.
 446

447 **Step 2: Bounding the Weighted Sum.** For any given h and non-negative constants $\{\omega_{h,k}\}_{h,[K]}$, we
 448 denote $\|\omega\|_{\infty,h} = \max_{k \in [K]} \omega_{h,k}$ and $\|\omega\|_{1,h} = \sum_{k \in [K]} \omega_{h,k}$. We also recursively define $\omega_{h',k}(h)$
 449 for any $h \leq h' \leq H, k \in [K]$ as follows:

$$450 \omega_{h,k}(h) := \omega_{h,k}; \omega_{h',j}(h) = \sum_{k=1}^K \sum_{n=1}^{N_h^k} \omega_{h'-1,k}(h) \eta_n^{N_h^k} \mathbb{I}[k^n = j], \forall j \in [K], h' > h. \quad (19)$$

451 Equation (19) implies the mapping from $\{\omega_{h,k}\}_{h,[K]}$ to $\{\omega_{h',k}(h)\}_{h',[K]}$ is linear. Equation (100)
 452 and Equation (101) shows that

$$453 \|\omega(h)\|_{1,h'} \leq \|\omega(h)\|_{1,h'-1}, \|\omega(h)\|_{\infty,h'} \leq (1 + 1/H) \|\omega(h)\|_{\infty,h'-1}, \forall h' > h. \quad (20)$$

454 Next, given non-negative constants $\{\omega_{h,k}\}_{h,[K]}$, we bound $\sum_{k=1}^K \omega_{h,k}(Q_h^k - Q_h^*)(s_h^k, a_h^k)$.
 455 In Equation (18) where we take summations with regard to k on both sides and
 456 apply the standard summation rearrangement technique given in Appendix D.4.1, we
 457 have $\sum_{k=1}^K \omega_{h,k}(Q_h^k - Q_h^*)(s_h^k, a_h^k) \leq \sum_{k=1}^K \omega_{h+1,k}(h)(Q_{h+1}^k - Q_{h+1}^*)(s_{h+1}^k, a_{h+1}^k) +$
 458 $\sum_{k=1}^K \omega_{h,k} \tilde{O} \left(\sqrt{(\mathbb{Q}^* + H\beta^2)/N_h^k} + H^2(N_h^k)^{-3/4} \right) + \sum_{k=1}^K \omega_{h,k} R_{\text{else}}^{h,k}$. Recurring it with regard
 459 to $h, h+1, \dots, H$, we have

$$460 \sum_{k=1}^K \omega_{h,k}(Q_h^k - Q_h^*)(s_h^k, a_h^k) \leq R_c + \sum_{k=1}^K \sum_{h'=h}^H \omega_{h',k}(h) R_{\text{else}}^{h',k}. \quad (21)$$

461 where $R_c = \sum_{k=1}^K \sum_{h'=h}^H \omega_{h',k}(h) \tilde{O} \left(\sqrt{(\mathbb{Q}^* + H\beta^2)/N_h^k} + H^2(N_h^k)^{-3/4} \right)$. Lemma D.3 in Ap-
 462 pendix D.2 implies that

$$463 R_c \leq \tilde{O} \left(H \left(\sqrt{\mathbb{Q}^* + \beta^2 H} \right) \sqrt{SA \|\omega\|_{\infty,h} \|\omega\|_{1,h}} + H^3 (SA \|\omega\|_{\infty,h})^{3/4} \|\omega\|_{1,h}^{1/4} \right). \quad (22)$$

464 **Step 3: Integrating Multiple Weighted Sums.** Next, consider multiple groups of weights related
 465 to $[\Delta_{\min}, H]$. We split it into N disjoint intervals $\mathcal{I}_i := [2^{i-1} \Delta_{\min}, 2^i \Delta_{\min})$ for $i \in [N-1]$ and
 466 $\mathcal{I}_N := [2^{N-1} \Delta_{\min}, 2^N \Delta_{\min}]$. Here, $N = \lceil \log_2(H/\Delta_{\min}) \rceil$. For any given $i \in [N]$ and $h \in [H]$, we
 467 denote $\omega_{h,k}^{(i)} = \mathbb{I}[(Q_h^k - Q_h^*)(s_h^k, a_h^k) \in \mathcal{I}_i]$. For $\{\omega_{h,k}^{(i)}\}_{h,[K]}$, we have $\|\omega^{(i)}\|_{\infty,h} \leq 1$ and

$$468 2^{i-1} \Delta_{\min} \|\omega^{(i)}\|_{1,h} \leq \sum_{k=1}^K \omega_{h,k}^{(i)} (Q_h^k - Q_h^*)(s_h^k, a_h^k) \leq 2^i \Delta_{\min} \|\omega^{(i)}\|_{1,h}. \quad (23)$$

469 Noticing that $\sum_{i=1}^N \sum_{k=1}^K \omega_{h,k}^{(i)} (Q_h^k - Q_h^*)(s_h^k, a_h^k) = \sum_{k=1}^K \text{clip}[(Q_h^k - Q_h^*)(s_h^k, a_h^k) | \Delta_{\min}]$ where
 470 $\text{clip}[x | \delta] := x \cdot \mathbb{I}[x \geq \delta]$, Equation (23) further implies

$$471 \sum_{k=1}^K \text{clip}[(Q_h^k - Q_h^*)(s_h^k, a_h^k) | \Delta_{\min}] = \Theta \left(\sum_{i=1}^N 2^i \Delta_{\min} \|\omega^{(i)}\|_{1,h} \right). \quad (24)$$

Letting $\omega_{h,k} = \omega_{h,k}^{(i)}$ in Equation (21) and applying Equations (22) and (23), we have

$$2^{i-1} \Delta_{\min} \|\omega^{(i)}\|_{1,h} \leq \tilde{O} \left(\theta_1 \sqrt{\|\omega^{(i)}\|_{1,h}} + \theta_2 \|\omega^{(i)}\|_{1,h}^{\frac{1}{4}} + \sum_{k=1}^K \sum_{h'=h}^H \omega_{h',k}(h) R_{\text{else}}^{h',k} \right). \quad (25)$$

Here, $\theta_1 = \sqrt{H^2 SA(\mathbb{Q}^* + \beta^2 H)}$, $\theta_2 = H^3 (SA)^{\frac{3}{4}}$. Thus, by Equation (106) in Appendix D.5.1,

$$\|\omega^{(i)}\|_{1,h} \leq \tilde{O} \left(\frac{(\mathbb{Q}^* + \beta^2 H) SAH^2}{4^{i-1} \Delta_{\min}^2} + \frac{H^4 SA}{(2^{i-1} \Delta_{\min})^{\frac{4}{3}}} + \frac{\sum_{k=1}^K \sum_{h'=h}^H \omega_{h',k}(h) R_{\text{else}}^{h',k}}{2^{i-1} \Delta_{\min}} \right).$$

This further implies

$$\sum_{i=1}^N 2^i \Delta_{\min} \|\omega^{(i)}\|_{1,h} \leq \tilde{O} \left(\frac{(\mathbb{Q}^* + \beta^2 H) SAH^2}{\Delta_{\min}} + \frac{H^4 SA}{\Delta_{\min}^{\frac{1}{3}}} + \sum_{k=1}^K \sum_{h'=h}^H \hat{\omega}_{h',k}(h) R_{\text{else}}^{h',k} \right). \quad (26)$$

where $\hat{\omega}_{h',k}(h) = \sum_{i=1}^N \omega_{h',k}^{(i)}(h)$. Noticing that $\hat{\omega}_{h,k}(h) = \mathbb{I}[(Q_h^k - Q_h^*)(s_h^k, a_h^k) \geq \Delta_{\min}]$, together with the linearity showed in Equation (19), Equation (20) implies $\hat{\omega}_{h',k}(h) \leq O(1)$, $\forall h \leq h' \leq H$. Thus, $\sum_{k=1}^K \sum_{h'=h}^H \hat{\omega}_{h',k}(h) R_{\text{else}}^{h',k} \leq O(\sum_{k=1}^K \sum_{h=1}^H R_{\text{else}}^{h,k})$. Appendix D.5.2 shows that with high probability,

$$\sum_{k=1}^K \sum_{h=1}^H R_{\text{else}}^{h,k} \leq \tilde{O}(H^6 SA / \beta). \quad (27)$$

Summarizing Equations (24), (26) and (27) and noticing that $H^4 SA / \Delta_{\min}^{\frac{1}{3}} = O(\beta^2 H^3 SA / \Delta_{\min} + H^4 SA / \beta + H^5 SA / \beta)$ that follows from the AM-GM inequality, we have

$$\sum_{k=1}^K \text{clip}[(Q_h^k - Q_h^*)(s_h^k, a_h^k) | \Delta_{\min}] = \tilde{O}(SAH^2(\mathbb{Q}^* + \beta^2 H) / \Delta_{\min} + H^6 SA / \beta). \quad (28)$$

Remark 2: Integrating groups of sums is first introduced in Yang et al. (2021) and also applied in Li et al. (2021). It leads to regret dependency on $1/\Delta_{\min}$ instead of $1/\Delta_{\min}^2$ that will appear if we do not use integration. We extend this method in handling $R_{\text{else}}^{h,k}$ that only appears in our proof: we apply the upper bound in Equation (27) after the integration instead of Equation (25) before the integration. This helps us remove the dependency on Δ_{\min} in the second term in Equation (28).

Remark 3: Equation (27) can be regarded as bounding the reference settling errors, which is related to $\hat{V}_h^{R,k}$ and the reference settling design in Q-EarlySettled-Advantage. UCB-Advantage and Q-EarlySettled-Advantage mainly differ on the reference settling policy, which results in different bounds for reference settling error and the gap-free regret terms in Equations (2) and (3). We show the details in Appendix D.5.2.

Step 4: Bounding the Expected Regret. By Equation (9), $Q_h^k(s_h^k, a_h^k) = V_h^k(s_h^k) \geq V_h^*(s_h^k)$. Thus,

$$\Delta_h(x_h^k, a_h^k) = \text{clip}[V_h^*(x_h^k) - Q_h^*(x_h^k, a_h^k) | \Delta_{\min}] \leq \text{clip}[(Q_h^k - Q_h^*)(x_h^k, a_h^k) | \Delta_{\min}], \forall (k, h).$$

Equation (4) of Yang et al. (2021) shows that $\mathbb{E}(\text{Regret}(K)) = \mathbb{E} \left[\sum_{k=1}^K \sum_{h=1}^H \Delta_h(x_h^k, a_h^k) \right]$, which further implies

$$\mathbb{E}(\text{Regret}(K)) \leq \mathbb{E} \left[\sum_{k=1}^K \sum_{h=1}^H \text{clip}[(Q_h^k - Q_h^*)(x_h^k, a_h^k) | \Delta_{\min}] \right]. \quad (29)$$

Using Equation (118) from Appendix D.6, which connects Equation (28) to Equation (29), we can derive the desired gap-dependent regret bound presented in Theorem 3.2.

5 CONCLUSION

In this paper, we have presented the first gap-dependent regret analysis for Q -learning using reference-advantage decomposition and also provided the first gap-dependent analysis on the policy switching cost of Q -learning, which answers two important open questions. Our novel error decomposition approach and construction of surrogate reference functions can be used in other problems using reference-advantage decomposition such as the offline Q -learning and stochastic learning.

540 REFERENCES
541

- Alekh Agarwal, Sham Kakade, and Lin F Yang. Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pp. 67–83. PMLR, 2020.
- Shipra Agrawal and Randy Jia. Optimistic posterior sampling for reinforcement learning: worst-case regret bounds. *Advances in Neural Information Processing Systems*, 30, 2017.
- Aymen Al Marjani, Aurélien Garivier, and Alexandre Proutiere. Navigating to the best policy in markov decision processes. *Advances in Neural Information Processing Systems*, 34:25852–25864, 2021.
- Amir H Ashouri, William Killian, John Cavazos, Gianluca Palermo, and Cristina Silvano. A survey on compiler autotuning using machine learning. *ACM Computing Surveys (CSUR)*, 51(5):1–42, 2018.
- Peter Auer and Ronald Ortner. Logarithmic online regret bounds for undiscounted reinforcement learning. In *Advances in Neural Information Processing Systems 19*, pp. 49–56. MIT Press, 2007.
- Peter Auer, Thomas Jaksch, and Ronald Ortner. Near-optimal regret bounds for reinforcement learning. *Advances in Neural Information Processing Systems*, 21, 2008.
- 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.
- Yu Bai, Tengyang Xie, Nan Jiang, and Yu-Xiang Wang. Provably efficient q-learning with low switching cost. *Advances in Neural Information Processing Systems*, 32, 2019.
- Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- Paul Covington, Jay Adams, and Emre Sargin. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems*, pp. 191–198, 2016.
- Christoph Dann, Lihong Li, Wei Wei, and Emma Brunskill. Policy certificates: Towards accountable reinforcement learning. In *International Conference on Machine Learning*, pp. 1507–1516. PMLR, 2019.
- Christoph Dann, Teodor Vanislavov Marinov, Mehryar Mohri, and Julian Zimmert. Beyond value-function gaps: Improved instance-dependent regret bounds for episodic reinforcement learning. *Advances in Neural Information Processing Systems*, 34:1–12, 2021.
- Simon S Du, Jianshu Chen, Lihong Li, Lin Xiao, and Dengyong Zhou. Stochastic variance reduction methods for policy evaluation. In *International Conference on Machine Learning*, pp. 1049–1058. PMLR, 2017.
- Ronan Fruit, Matteo Pirotta, Alessandro Lazaric, and Ronald Ortner. Efficient bias-span-constrained exploration-exploitation in reinforcement learning. In *International Conference on Machine Learning*, pp. 1578–1586. PMLR, 2018.
- Robert M Gower, Mark Schmidt, Francis Bach, and Peter Richtárik. Variance-reduced methods for machine learning. *Proceedings of the IEEE*, 108(11):1968–1983, 2020.
- Shixiang Gu, Ethan Holly, Timothy Lillicrap, and Sergey Levine. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3389–3396. IEEE, 2017.
- Jiafan He, Dongruo Zhou, and Quanquan Gu. Logarithmic regret for reinforcement learning with linear function approximation. In *International Conference on Machine Learning*, pp. 4171–4180. PMLR, 2021.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11:1563–1600, 2010.

- 594 Chi Jin, Zeyuan Allen-Zhu, Sébastien Bubeck, and Michael I Jordan. Is q-learning provably effi-
 595 cient? *Advances in Neural Information Processing Systems*, 31, 2018.
- 596
- 597 Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration
 598 for reinforcement learning. In *International Conference on Machine Learning*, pp. 4870–4879.
 599 PMLR, 2020.
- 600 Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance
 601 reduction. *Advances in Neural Information Processing Systems*, 26, 2013.
- 602
- 603 Anders Jonsson, Emilie Kaufmann, Pierre Ménard, Omar Darwiche Domingues, Edouard Leurent,
 604 and Michal Valko. Planning in markov decision processes with gap-dependent sample complexity.
 605 *Advances in Neural Information Processing Systems*, 33:1253–1263, 2020.
- 606 Sham Kakade, Mengdi Wang, and Lin F Yang. Variance reduction methods for sublinear reinforce-
 607 ment learning. *arXiv preprint arXiv:1802.09184*, 2018.
- 608 Koulik Khamaru, Ashwin Pananjady, Feng Ruan, Martin J Wainwright, and Michael I Jordan. Is
 609 temporal difference learning optimal? an instance-dependent analysis. *SIAM Journal on Mathe-
 610 matics of Data Science*, 3(4):1013–1040, 2021.
- 611
- 612 Jens Kober, J Andrew Bagnell, and Jan Peters. Reinforcement learning in robotics: A survey. *The
 613 International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- 614 Sanjay Krishnan, Zongheng Yang, Ken Goldberg, Joseph Hellerstein, and Ion Stoica. Learning to
 615 optimize join queries with deep reinforcement learning. *arXiv preprint arXiv:1808.03196*, 2018.
- 616
- 617 Gen Li, Yuting Wei, Yuejie Chi, Yuantao Gu, and Yuxin Chen. Sample complexity of asynchronous
 618 q-learning: Sharper analysis and variance reduction. *Advances in Neural Information Processing
 619 Systems*, 33:7031–7043, 2020.
- 620 Gen Li, Laixi Shi, Yuxin Chen, Yuantao Gu, and Yuejie Chi. Breaking the sample complexity barrier
 621 to regret-optimal model-free reinforcement learning. *Advances in Neural Information Processing
 622 Systems*, 34:17762–17776, 2021.
- 623 AA Marjani and Alexandre Proutiere. Best policy identification in discounted mdps: Problem-
 624 specific sample complexity. *arXiv preprint arXiv:2009.13405*, 2020.
- 625
- 626 Pierre Ménard, Omar Darwiche Domingues, Xuedong Shang, and Michal Valko. Ucb momen-
 627 tum q-learning: Correcting the bias without forgetting. In *International Conference on Machine
 628 Learning*, pp. 7609–7618. PMLR, 2021.
- 629 Azalia Mirhoseini, Hieu Pham, Quoc V Le, Benoit Steiner, Rasmus Larsen, Yuefeng Zhou, Naveen
 630 Kumar, Mohammad Norouzi, Samy Bengio, and Jeff Dean. Device placement optimization
 631 with reinforcement learning. In *International Conference on Machine Learning*, pp. 2430–2439.
 632 PMLR, 2017.
- 633 Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. Sarah: A novel method for machine
 634 learning problems using stochastic recursive gradient. In *International Conference on Machine
 635 Learning*, pp. 2613–2621. PMLR, 2017.
- 636
- 637 Phuoc Nguyen, Truyen Tran, Sunil Gupta, Santu Rana, Matthew Barnett, and Svetha Venkatesh.
 638 Incomplete conditional density estimation for fast materials discovery. In *Proceedings of the
 639 2019 SIAM International Conference on Data Mining*, pp. 549–557. SIAM, 2019.
- 640 Thanh Nguyen-Tang, Ming Yin, Sunil Gupta, Svetha Venkatesh, and Raman Arora. On instance-
 641 dependent bounds for offline reinforcement learning with linear function approximation. In *Pro-
 642 ceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 9310–9318, 2023.
- 643 Jungseul Ok, Alexandre Proutiere, and Damianos Tranos. Exploration in structured reinforcement
 644 learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- 645
- 646 Dan Qiao, Ming Yin, Ming Min, and Yu-Xiang Wang. Sample-efficient reinforcement learning with
 647 loglog (t) switching cost. In *International Conference on Machine Learning*, pp. 18031–18061.
 648 PMLR, 2022.

- 648 Laixi Shi, Gen Li, Yuting Wei, Yuxin Chen, and Yuejie Chi. Pessimistic q-learning for offline
 649 reinforcement learning: Towards optimal sample complexity. In *International Conference on*
 650 *Machine Learning*, pp. 19967–20025. PMLR, 2022.
- 651
- 652 Aaron Sidford, Mengdi Wang, Xian Wu, Lin Yang, and Yinyu Ye. Near-optimal time and sample
 653 complexities for solving markov decision processes with a generative model. *Advances in Neural*
 654 *Information Processing Systems*, 31, 2018.
- 655 Aaron Sidford, Mengdi Wang, Xian Wu, and Yinyu Ye. Variance reduced value iteration and faster
 656 algorithms for solving markov decision processes. *Naval Research Logistics (NRL)*, 70(5):423–
 657 442, 2023.
- 658
- 659 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
 660 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
 661 the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- 662
- 663 David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez,
 664 Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. Mastering chess and shogi
 665 by self-play with a general reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815*,
 2017.
- 666
- 667 David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez,
 668 Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. A general reinforcement
 669 learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–
 1144, 2018.
- 670
- 671 Max Simchowitz and Kevin G Jamieson. Non-asymptotic gap-dependent regret bounds for tabular
 672 mdps. *Advances in Neural Information Processing Systems*, 32, 2019.
- 673 R Sutton and A Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.
- 674
- 675 Mohammad Sadegh Talebi and Odalric-Ambrym Maillard. Variance-aware regret bounds for undis-
 676 counted reinforcement learning in mdps. In *Algorithmic Learning Theory*, pp. 770–805. PMLR,
 677 2018.
- 678 Ambuj Tewari and Peter Bartlett. Optimistic linear programming gives logarithmic regret for irre-
 679 ducible mdps. In *Advances in Neural Information Processing Systems 20: Proceedings of the*
 680 *2007 Conference*, pp. 1505–1512. Neural Information Processing Systems (NIPS) Foundation,
 681 2008.
- 682
- 683 Andrea Tirinzoni, Aymen Al Marjani, and Emilie Kaufmann. Near instance-optimal pac reinfor-
 684 cement learning for deterministic mdps. *Advances in neural information processing systems*, 35:
 685 8785–8798, 2022.
- 686
- 687 Andrea Tirinzoni, Aymen Al-Marjani, and Emilie Kaufmann. Optimistic pac reinforcement learn-
 688 ing: the instance-dependent view. In *International Conference on Algorithmic Learning Theory*,
 689 pp. 1460–1480. PMLR, 2023.
- 690
- 691 Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Juny-
 692 oung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster
 693 level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- 694
- 695 Andrew Wagenmaker and Kevin G Jamieson. Instance-dependent near-optimal policy identifica-
 696 tion in linear mdps via online experiment design. *Advances in Neural Information Processing Systems*,
 697 35:5968–5981, 2022.
- 698
- 699 Andrew J Wagenmaker, Yifang Chen, Max Simchowitz, Simon Du, and Kevin Jamieson. First-
 700 order regret in reinforcement learning with linear function approximation: A robust estimation
 701 approach. In *International Conference on Machine Learning*, pp. 22384–22429. PMLR, 2022a.
- 702
- 703 Andrew J Wagenmaker, Max Simchowitz, and Kevin Jamieson. Beyond no regret: Instance-
 704 dependent pac reinforcement learning. In *Conference on Learning Theory*, pp. 358–418. PMLR,
 705 2022b.

- 702 Hoi-To Wai, Mingyi Hong, Zhuoran Yang, Zhaoran Wang, and Kexin Tang. Variance reduced policy
 703 evaluation with smooth function approximation. *Advances in Neural Information Processing*
 704 *Systems*, 32, 2019.
- 705 Martin J Wainwright. Variance-reduced q -learning is minimax optimal. *arXiv preprint*
 706 *arXiv:1906.04697*, 2019.
- 708 Tianhao Wang, Dongruo Zhou, and Quanquan Gu. Provably efficient reinforcement learning with
 709 linear function approximation under adaptivity constraints. *Advances in Neural Information Pro-*
 710 *cessing Systems*, 34:13524–13536, 2021.
- 712 Xinqi Wang, Qiwen Cui, and Simon S Du. On gap-dependent bounds for offline reinforcement
 713 learning. *Advances in Neural Information Processing Systems*, 35:14865–14877, 2022.
- 714 Haike Xu, Tengyu Ma, and Simon Du. Fine-grained gap-dependent bounds for tabular mdps via
 715 adaptive multi-step bootstrap. In *Conference on Learning Theory*, pp. 4438–4472. PMLR, 2021.
- 717 Tengyu Xu, Zhe Wang, Yi Zhou, and Yingbin Liang. Reanalysis of variance reduced temporal
 718 difference learning. In *International Conference on Learning Representations*, 2020.
- 719 Yuling Yan, Gen Li, Yuxin Chen, and Jianqing Fan. The efficacy of pessimism in asynchronous
 720 q -learning. *IEEE Transactions on Information Theory*, 2023.
- 722 Kunhe Yang, Lin Yang, and Simon Du. Q-learning with logarithmic regret. In *International Con-*
 723 *ference on Artificial Intelligence and Statistics*, pp. 1576–1584. PMLR, 2021.
- 725 Ming Yin, Yu Bai, and Yu-Xiang Wang. Near-optimal offline reinforcement learning via double
 726 variance reduction. *Advances in Neural Information Processing Systems*, 34:7677–7688, 2021.
- 727 Ekim Yurtsever, Jacob Lambert, Alexander Carballo, and Kazuya Takeda. A survey of autonomous
 728 driving: Common practices and emerging technologies. *IEEE Access*, 8:58443–58469, 2020.
- 730 Andrea Zanette and Emma Brunskill. Tighter problem-dependent regret bounds in reinforcement
 731 learning without domain knowledge using value function bounds. In *International Conference on*
 732 *Machine Learning*, pp. 7304–7312. PMLR, 2019.
- 733 Zihan Zhang, Yuan Zhou, and Xiangyang Ji. Almost optimal model-free reinforcement learning
 734 via reference-advantage decomposition. *Advances in Neural Information Processing Systems*, 33:
 735 15198–15207, 2020.
- 737 Zihan Zhang, Xiangyang Ji, and Simon Du. Is reinforcement learning more difficult than bandits?
 738 a near-optimal algorithm escaping the curse of horizon. In *Conference on Learning Theory*, pp.
 739 4528–4531. PMLR, 2021a.
- 740 Zihan Zhang, Yuan Zhou, and Xiangyang Ji. Model-free reinforcement learning: from clipped
 741 pseudo-regret to sample complexity. In *International Conference on Machine Learning*, pp.
 742 12653–12662. PMLR, 2021b.
- 744 Zihan Zhang, Xiangyang Ji, and Simon Du. Horizon-free reinforcement learning in polynomial
 745 time: the power of stationary policies. In *Conference on Learning Theory*, pp. 3858–3904. PMLR,
 746 2022a.
- 748 Zihan Zhang, Yuhang Jiang, Yuan Zhou, and Xiangyang Ji. Near-optimal regret bounds for multi-
 749 batch reinforcement learning. *Advances in Neural Information Processing Systems*, 35:24586–
 750 24596, 2022b.
- 751 Zihan Zhang, Yuxin Chen, Jason D Lee, and Simon S Du. Settling the sample complexity of online
 752 reinforcement learning. *arXiv preprint arXiv:2307.13586*, 2023.
- 753 Zihan Zhang, Yuxin Chen, Jason D Lee, and Simon S Du. Settling the sample complexity of online
 754 reinforcement learning. In *The Thirty Seventh Annual Conference on Learning Theory*, pp. 5213–
 755 5219. PMLR, 2024.

756 Heyang Zhao, Jiafan He, Dongruo Zhou, Tong Zhang, and Quanquan Gu. Variance-dependent regret
757 bounds for linear bandits and reinforcement learning: Adaptivity and computational efficiency. In
758 *The Thirty Sixth Annual Conference on Learning Theory*, pp. 4977–5020. PMLR, 2023.

759
760 Runlong Zhou, Zhang Zihan, and Simon Shaolei Du. Sharp variance-dependent bounds in reinforce-
761 ment learning: Best of both worlds in stochastic and deterministic environments. In *International
762 Conference on Machine Learning*, pp. 42878–42914. PMLR, 2023.

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810 **A GENERAL LEMMAS**
 811

812 **Lemma A.1.** (Azuma-Hoeffding Inequality) Suppose $\{X_k\}_{k=0}^{\infty}$ is a martingale and $|X_k - X_{k-1}| \leq$
 813 $c_k, \forall k \in \mathbb{N}_+$, almost surely. Then for any positive integers N and any positive real number ϵ , it
 814 holds that:

$$815 \quad \mathbb{P}(X_N - X_0 \geq \epsilon) \leq \exp\left(-\frac{\epsilon^2}{2 \sum_{k=1}^N c_k^2}\right),$$

816 and

$$817 \quad \mathbb{P}(|X_N - X_0| \geq \epsilon) \leq 2 \exp\left(-\frac{\epsilon^2}{2 \sum_{k=1}^N c_k^2}\right).$$

818 **Lemma A.2.** (Lemma 10 in Zhang et al. (2022a)) Let X_1, X_2, \dots be a sequence of random vari-
 819 ables taking value in $[0, l]$. Define $\mathcal{F}_k = \sigma(X_1, X_2, \dots, X_{k-1})$ and $Y_k = \mathbb{E}[X_k | \mathcal{F}_k]$ for $k \geq 1$. For
 820 any $\delta > 0$, we have that

$$821 \quad \mathbb{P}\left[\exists n, \sum_{k=1}^n X_k \geq 3 \sum_{k=1}^n Y_k + l \log(1/\delta)\right] \leq \delta$$

822 and

$$823 \quad \mathbb{P}\left[\exists n, \sum_{k=1}^n Y_k \geq 3 \sum_{k=1}^n X_k + l \log(1/\delta)\right] \leq \delta.$$

824 **Lemma A.3.** (Lemma 11 in Zhang et al. (2021b)) Let $(M_n)_{n \geq 0}$ be a martingale such that $M_0 = 0$
 825 and $|M_n - M_{n-1}| \leq c$ for some $c > 0$ and any $n \geq 1$. Let

$$826 \quad \text{Var}_n = \sum_{k=1}^n \mathbb{E}[(M_k - M_{k-1})^2 | \mathcal{F}_{k-1}]$$

827 for $n \geq 0$, where $\mathcal{F}_k = \sigma(M_1, \dots, M_k)$. Then for any positive integer n and any $\epsilon, \delta > 0$, we have
 828 that

$$829 \quad \mathbb{P}\left(|M_n| \geq 2\sqrt{2\text{Var}_n \ln(1/\delta)} + 2\sqrt{\epsilon \ln(1/\delta)} + 2c \ln(1/\delta)\right) \leq 2 \left(\log_2\left(\frac{nc^2}{\epsilon}\right) + 1\right) \delta.$$

830 **B PROOF OF THEOREM 3.1**

831 **B.1 ALGORITHM DETAILS**

832 The UCB-Advantage algorithm, first introduced in Zhang et al. (2020), achieves the information-
 833 theoretic bound on regret up to logarithmic factors, using a model-free algorithm. The key innova-
 834 tion of the algorithm lies in its combination of UCB exploration (Jin et al., 2018) with a newly
 835 introduced reference-advantage decomposition for updating Q -estimates.

836 Before discussing the algorithm in detail, we will first review the special stage design used in the
 837 algorithm. For any triple (s, a, h) , we divide the samples received for the triple into consecutive
 838 stages. Define $e_1 = H$ and $e_{i+1} = \lfloor (1 + \frac{1}{H})e_i \rfloor$ for all $i \geq 1$, standing for the length of the stages.
 839 We also let $\mathcal{L} := \{\sum_{i=1}^j e_i | j = 1, 2, 3, \dots\}$ be the set of indices marking the ends of the stages.

840 We note that the definition of stages is with respect to the triple (s, a, h) . For any fixed pair of k and
 841 h , let (s_h^k, a_h^k) be the state-action pair at the h -th step during the k -th episode of the algorithm. We
 842 say that (k, h) falls in the j -th stage of (s, a, h) if and only if $(s, a) = (s_h^k, a_h^k)$ and the total visit
 843 number of (s_h^k, a_h^k) after the k -th episode is in $(\sum_{i=1}^{j-1} e_i, \sum_{i=1}^j e_i]$.

844 Now we introduce the stage-based update framework. For any (s, a, h) triple, we update $Q_h(s, a)$
 845 when the total visit number of (s, a, h) reaches the end of the current stage (in other word, the
 846 total visit number occurs in \mathcal{L}). For k -th episode at the end of a given stage, the Q -estimate
 847 $Q_h^{1,k+1}(s_h^k, a_h^k)$ learned from UCB is updated to:

$$848 \quad Q_h^{1,k+1}(s_h^k, a_h^k) = r_h^k(s_h^k, a_h^k) + \frac{1}{\check{n}_h^k} \sum_{i=1}^{\check{n}_h^k} V_{h+1}^{l_i}(s_{h+1}^{l_i}) + 2\sqrt{\frac{H^2 \iota}{\check{n}_h^k}}. \quad (30)$$

864 Here we define $\check{n}_h^k = \check{n}_h^k(s_h^k, a_h^k)$ be the number of visits to (s_h^k, a_h^k, h) during the stage immediately
 865 before the stage of k -th episode and $\check{l}_i = \check{l}_{h,k}^i$ denotes the index of the i -th episode among the \check{n}_h^k
 866 episodes. $V_h^k(s)$ is the V -estimate at the end of the episode $k - 1$ with the initial value $V_h^1(s) = H$.
 867 The term $2\sqrt{\frac{H^2\iota}{\check{n}_h^k}}$ represents the exploration bonus for \check{n}_h^k -th visit, where $\iota = \log(\frac{2}{p})$ with $p \in (0, 1)$
 868 being failure probability.
 869

870 The other estimate, denoted by $Q_h^{2,k+1}(s_h^k, a_h^k)$, uses the reference-advantage decomposition tech-
 871 nique. For k -th episode at the end of a given stage, it is updated to:
 872

$$873 \quad r_h^k(s_h^k, a_h^k) + \frac{1}{n_h^k} \sum_{i=1}^{n_h^k} V_{h+1}^{\text{ref}, k^l_i}(s_{h+1}^{l_i}) + \frac{1}{\check{n}_h^k} \sum_{i=1}^{\check{n}_h^k} \left(V_{h+1}^{k^l_i} - V_{h+1}^{\text{ref}, k^l_i} \right) (s_{h+1}^{l_i}) + b_h^{k+1}(s_h^k, a_h^k). \quad (31)$$

873 Here we define $n_h^k = n_h^k(s_h^k, a_h^k)$ be the number of visits to (s_h^k, a_h^k, h) prior to the stage of (k) -th
 874 episode and $l_i = l_{h,k}^i$ denotes the index of i -th episode among the n_h^k episodes.
 875

876 In Equation (31), $V_h^{\text{ref}, k}(s)$ is the reference function learned at the end of episode $k - 1$. We expect
 877 that for any $s \in \mathcal{S}$, sufficiently large k and some given $\beta \in (0, H]$, it holds $|V_h^{\text{ref}, k}(s) - V_h^*(s)| \leq$
 878 β . In this case, for $s_{h+1}^{k^n} \sim \mathbb{P}_h(\cdot | s_h^{k^n}, a_h^{k^n})$, the variance of the advantage term $V_{h+1}^{k^l_i}(s_{h+1}^{k^n}) -$
 879 $V_{h+1}^{\text{ref}, k^l_i}(s_{h+1}^{k^n})$, is bounded by β^2 , which can be less volatile than the stochastic term $V_{h+1}^{k^l_i}(s_{h+1}^{k^n})$,
 880 whose variance can be H^2 . Meanwhile, the reference term $\sum_{i=1}^{n_h^k} V_{h+1}^{\text{ref}, k^l_i}(s_{h+1}^{k^n})/n_h^k$ use a batch
 881 of historical visits to (s_h^k, a_h^k, h) , which can lower the variance as the increase of the sample size
 882 n_h^k . Accordingly, the exploration bonus term b_h^{k+1} is taken to be an upper confidence bound for the
 883 above-mentioned two terms combined.
 884

885 With these Q -estimates, we can update the final Q -estimate as follows:
 886

$$887 \quad Q_h^{k+1}(s_h^k, a_h^k) = \min\{Q_h^{1,k+1}(s_h^k, a_h^k), Q_h^{2,k+1}(s_h^k, a_h^k), Q_h^k(s_h^k, a_h^k)\}. \quad (32)$$

888 We also incorporate $Q_h^k(s_h^k, a_h^k)$ here to keep the monotonicity of the update. Then we can learn
 889 $V_h^{k+1}(s_h^k)$ by a greedy policy with respect to the Q -estimates $V_h^{k+1}(s_h^k) = \max_a Q_h^{k+1}(s_h^k, a)$. If
 890 the number of visits to the state-step pair (s, h) first exceeds $N_0 = O(\frac{SAH^5\iota}{\beta^2})$ at k -th episode, then
 891 we learn the final reference function $V_h^{\text{REF}}(s) = V_h^{k+1}(s)$. For the reader's convenience, we have
 892 also provided the detailed algorithm below.
 893

894 Algorithm 1 UCB-Advantage

895 1: **Initialize:** set all accumulators to 0; for all $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$, set $Q_h(s, a)$, $V_h(s, a) \leftarrow$
 896 $H - h + 1$; $V_h^{\text{ref}}(s) \leftarrow H$;
 897 2: **for** episodes $k \leftarrow 1, 2, \dots, K$ **do**
 898 3: observe s_1 ;
 899 4: **for** $h \leftarrow 1, 2, \dots, H$ **do**
 900 5: Take action $a_h \leftarrow \arg \max_a Q_h(s_h, a)$, and observe s_{h+1} .
 901 6: Update the accumulators by $n := n_h(s_h, a_h) \leftarrow^+ 1$, $\check{n} := \check{n}_h(s_h, a_h) \leftarrow^+ 1$,
 902 7: and Equation (33), Equation (34), Equation (35).
 903 8: **if** $n \in \mathcal{L}$ **then**
 904 9: $b \leftarrow 2\sqrt{\frac{\sigma_h^{\text{ref}}/n - (\mu_h^{\text{ref}}/n)^2}{n}}\iota + 2\sqrt{\frac{\check{\sigma}/\check{n} - (\check{\mu}/\check{n})^2}{\check{n}}}\iota + 5\left(\frac{H\iota}{n} + \frac{H\iota}{\check{n}} + \frac{H\iota^{3/4}}{n^{3/4}} + \frac{H\iota^{3/4}}{\check{n}^{3/4}}\right)$;
 905 10: $\bar{b} \leftarrow 2\sqrt{\frac{H^2}{\check{n}}}\iota$;
 906 11: $Q_h(s_h, a_h) \leftarrow \min\{r_h(s_h, a_h) + \frac{\check{v}}{\check{n}} + \bar{b}, r_h(s_h, a_h) + \frac{\mu^{\text{ref}}}{n} + \frac{\check{\mu}}{\check{n}} + b, Q_h(s_h, a_h)\}$;
 907 12: $V_h(s_h) \leftarrow \max_a Q_h(s_h, a)$;
 908 13: $\check{v}_h(s_h, a_h), \check{\mu}_h(s_h, a_h), \check{v}_h(s_h, a_h), \check{\sigma}_h(s_h, a_h) \leftarrow 0$;
 909 14: **end if**
 910 15: **if** $\sum_a n_h(s_h, a) = N_0$ **then** $V_h^{\text{ref}}(s_h) \leftarrow V_h(s_h)$
 911 16: **end if**
 912 17: **end for**
 913 18: **end for**

918 The accumulators in the algorithm are updated as follows.
 919

$$\check{\mu} := \check{\mu}_h(s_h, a_h) \leftarrow^+ V_{h+1}(s_{h+1}) - V_{h+1}^{\text{ref}}(s_{h+1}); \quad \check{v} := \check{v}_h(s_h, a_h) \leftarrow^+ V_{h+1}(s_{h+1}); \quad (33)$$

$$\check{\sigma} := \check{\sigma}_h(s_h, a_h) \leftarrow^+ (V_{h+1}(s_{h+1}) - V_{h+1}^{\text{ref}}(s_{h+1}))^2; \quad (34)$$

923 Meanwhile, the following two types of global accumulators are used for the samples in all stages
 924

$$\mu^{\text{ref}} := \mu_h^{\text{ref}}(s_h, a_h) \leftarrow V_{h+1}^{\text{ref}}(s_{h+1}); \quad \sigma_h^{\text{ref}} := \sigma^{\text{ref}}(s_h, a_h) \leftarrow (V_{h+1}^{\text{ref}}(s_{h+1}))^2. \quad (35)$$

927 We use $\mu_h^{\text{ref},k}$, $\sigma_h^{\text{ref},k}$, $\check{\mu}_h^k$, \check{v}_h^k , $\check{\sigma}_h^k$, b_h^k to denote respectively the values of μ^{ref} , σ^{ref} , $\check{\mu}$, \check{v} , $\check{\sigma}$, b at step h
 928 by the start of the k -th episode.

930 B.2 KEY LEMMAS

932 Before proceeding to the proof, we will first establish several key lemmas. In the algorithm, define
 933 $\iota = \log(2/p)$ with $p \in (0, 1)$ being the failure probability.

934 **Lemma B.1.** *Using $\forall(s, a, h, k)$ as the simplified notation for $\forall(s, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$.
 935 Let $N_h^k(s) = \sum_a n_h^k(s, a)$, $\lambda_h^k(s) = \mathbb{I}[N_h^k(s) < N_0]$, $\hat{V}_h^{\text{ref},k}(s) = \max\{V_h^*(s), \min\{V_h^*(s) + \beta, V_h^{\text{ref},k}(s)\}\}$. Then we have the following conclusions:*

938 (a) (Proposition 4 in Zhang et al. (2020)) *With probability at least $1 - (4H^2T^4 + 12T)p$, the
 939 following event holds:*

$$\mathcal{E}_1 = \{Q_h^*(s, a) \leq Q_h^{k+1}(s, a) \leq Q_h^k(s, a), \forall(s, a, h, k)\}.$$

942 (b) (Corollary 6 in Zhang et al. (2020)) *With probability at least $1 - (4H^2T^4 + 13T)p$, the
 943 following event holds:*

$$\mathcal{E}_2 = \{N_h^k(s) \geq N_0 \Rightarrow V_h^*(s) \leq V_h^{\text{ref},k} \leq V_h^*(s) + \beta, \forall(s, a, h, k)\}.$$

948 (c) *With probability at least $1 - p$, the following event holds:*

$$\mathcal{E}_3 = \left\{ \sum_{k=1}^K \mathbb{P}_{s_h^k, a_h^k, h} \lambda_{h+1}^k \leq 3 \sum_{k=1}^K \lambda_{h+1}^k(s_{h+1}^k) + \iota \right\}.$$

953 Especially, $\lambda_{H+1}^k(s) = 0$.

954 (d) *With probability at least $1 - SATp$, the following event holds:*

$$\mathcal{E}_4 = \left\{ \frac{\left| \sum_{i=1}^{n_h^k} \left(\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s, a, h} \right) \left(\hat{V}_{h+1}^{\text{ref}, l_i} - V_{h+1}^* \right) \right|}{n_h^k(s, a)} \leq \beta \sqrt{\frac{2\iota}{n_h^k(s, a)}}, \forall(s, a, h, k) \right\}.$$

961 (e) *With probability at least $1 - SAT^2p$, the following event holds:*

$$\mathcal{E}_5 = \left\{ \frac{\left| \sum_{i=1}^{n_h^k} \left(\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s, a, h} \right) V_{h+1}^* \right|}{n_h^k(s, a)} \leq 2 \sqrt{\frac{2\mathbb{Q}^*\iota}{n_h^k(s, a)}} + \frac{4H\iota}{n_h^k(s, a)}, \forall(s, a, h, k) \right\}.$$

967 (f) *With probability at least $1 - SAT^2p$, the following event holds:*

$$\mathcal{E}_6 = \left\{ \frac{\left| \sum_{i=1}^{\check{n}_h^k} \left(\mathbb{1}_{s_{h+1}^{\check{l}_i}} - \mathbb{P}_{s, a, h} \right) \left(\hat{V}_{h+1}^{\text{ref}, \check{l}_i} - V_{h+1}^* \right) \right|}{\check{n}_h^k(s, a)} \leq \beta \sqrt{\frac{2\iota}{\check{n}_h^k(s, a)}}, \forall(s, a, h, k) \right\}.$$

972 (g) With probability at least $1 - SATp$, the following event holds:
 973

$$974 \quad 975 \quad 976 \quad 977 \quad \mathcal{E}_7 = \left\{ \frac{\left| \sum_{i=1}^{n_h^k} (\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s,a,h}) (V_{h+1}^*)^2 \right|}{n_h^k(s,a)} \leq H^2 \sqrt{\frac{2\iota}{n_h^k(s,a)}}, \forall (s,a,h,k) \right\}.$$

978 *Proof.* We only need prove (c) to (e).
 979

980 (c) Using Lemma A.2 with $l = 1$ and $\delta = p$, we can prove this conclusion.
 981

982 (d) From the definition of $\hat{V}_h^{\text{ref},k}(s)$, we know that for any $k \in [K]$:
 983

$$984 \quad 985 \quad V_h^*(s) \leq \hat{V}_h^{\text{ref},k}(s) \leq V_h^*(s) + \beta. \quad (36)$$

986 Then the sequence $\{\sum_{i=1}^j (\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s,a,h})(\hat{V}_{h+1}^{\text{ref},l_i} - V_{h+1}^*)\}_{j \in \mathbb{N}^+}$ is a martingale sequence
 987 with $|(\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s,a,h})(\hat{V}_{h+1}^{\text{ref},l_i} - V_{h+1}^*)| \leq \beta$. Then according to Azuma-Hoeffding inequality,
 988 for any $p \in (0, 1)$, with probability at least $1 - p$, it holds for given $n_h^k(s,a) = n \in \mathbb{N}_+$
 989 that:
 990

$$991 \quad 992 \quad 993 \quad \frac{1}{n} \left| \sum_{i=1}^n (\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s,a,h}) (\hat{V}_{h+1}^{\text{ref},l_i} - V_{h+1}^*) \right| \leq \sqrt{\frac{2\beta^2\iota}{n}}.$$

994 For any all $(s,a,h,k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$, we have $n_h^k(s,a) \in [\frac{T}{H}]$. Considering all the
 995 possible combinations $(s,a,h,n) \in \mathcal{S} \times \mathcal{A} \times [H] \times [\frac{T}{H}]$, with probability at least $1 - SATp$,
 996 it holds simultaneously for all $(s,a,h,k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$ that:
 997

$$998 \quad 999 \quad 1000 \quad \frac{1}{n_h^k(s,a)} \left| \sum_{i=1}^{n_h^k} (\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s,a,h}) (\hat{V}_{h+1}^{\text{ref},l_i} - V_{h+1}^*) \right| \leq \sqrt{\frac{2\beta^2\iota}{n_h^k(s,a)}}.$$

1001 (e) The sequence $\{\sum_{i=1}^j (\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s,a,h}) V_{h+1}^*\}_{j \in \mathbb{N}^+}$ is a martingale sequence with $|(\mathbb{1}_{s_{h+1}^{l_i}} -$
 1002 $\mathbb{P}_{s,a,h}) V_{h+1}^*| \leq H$. Using Lemma A.3 with $c = H$, $\epsilon = H^2$ and $\delta = \frac{p}{2}$, for a given
 1003 $n_h^k(s,a) = n \in [\frac{T}{H}]$, with probability at least $1 - (\log_2(n) + 1)p \geq 1 - Tp$, we have:
 1004

$$1005 \quad 1006 \quad 1007 \quad \frac{1}{n} \left| \sum_{i=1}^n (\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s,a,h}) V_{h+1}^* \right| \leq 2\sqrt{\frac{2\mathbb{Q}^*\iota}{n}} + \frac{4H\iota}{n}$$

1008 Considering all the possible combinations $(s,a,h,n) \in \mathcal{S} \times \mathcal{A} \times [H] \times [\frac{T}{H}]$, with probability
 1009 at least $1 - SAT^2p$, it holds simultaneously for all $(s,a,h,k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$ that:
 1010

$$1011 \quad 1012 \quad 1013 \quad 1014 \quad \frac{1}{n_h^k(s,a)} \left| \sum_{i=1}^{n_h^k} (\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s,a,h}) V_{h+1}^* \right| \leq 2\sqrt{\frac{2\mathbb{Q}^*}{n_h^k(s,a)}} + \frac{4H\iota}{n_h^k(s,a)}.$$

1015 (f) The sequence $\{\sum_{i=1}^j (\mathbb{1}_{s_{h+1}^{\tilde{l}_i}} - \mathbb{P}_{s,a,h})(\hat{V}_{h+1}^{\text{ref},\tilde{l}_i} - V_{h+1}^*)\}_{j \in \mathbb{N}^+}$ is a martingale sequence with
 1016 $|(\mathbb{1}_{s_{h+1}^{\tilde{l}_i}} - \mathbb{P}_{s,a,h})(\hat{V}_{h+1}^{\text{ref},\tilde{l}_i} - V_{h+1}^*)| \leq \beta$. Then according to Azuma-Hoeffding inequality,
 1017 for any $p \in (0, 1)$, with probability at least $1 - p$, it holds for given $\check{n}_h^k(s,a) = \check{n} \in \mathbb{N}_+$
 1018 that:
 1019

$$1020 \quad 1021 \quad 1022 \quad \frac{1}{\check{n}} \left| \sum_{i=1}^{\check{n}} (\mathbb{1}_{s_{h+1}^{\tilde{l}_i}} - \mathbb{P}_{s,a,h}) (\hat{V}_{h+1}^{\text{ref},\tilde{l}_i} - V_{h+1}^*) \right| \leq \sqrt{\frac{2\beta^2\iota}{\check{n}}}.$$

1023 For any all $(s,a,h,k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$, we have $\check{n}_h^k(s,a) \in [\frac{T}{H}]$. Considering
 1024 all the possible combinations $(s,a,h,k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$ and $\check{n}_h^k(s,a) \in [\frac{T}{H}]$,
 1025

1026 with probability at least $1 - SAT^2/Hp \geq 1 - SAT^2p$, it holds simultaneously for all
 1027 $(s, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$ that:

$$\frac{1}{\tilde{n}_h^k(s, a)} \left| \sum_{i=1}^{\tilde{n}_h^k} \left(\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s, a, h} \right) (\hat{V}_{h+1}^{\text{ref}, l_i} - V_{h+1}^*) \right| \leq \sqrt{\frac{2\beta^2\iota}{\tilde{n}_h^k(s, a)}}.$$

1028 (g) The sequence $\{\sum_{i=1}^j (\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s, a, h})(V_{h+1}^*)^2\}_{j \in \mathbb{N}^+}$ is a martingale sequence with
 1029 $|\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s, a, h}| \leq H^2$. Then according to Azuma-Hoeffding inequality, with
 1030 probability at least $1 - p$, it holds for given $n_h^k(s, a) = n$ that:

$$\frac{1}{n} \left| \sum_{i=1}^n \left(\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s, a, h} \right) (V_{h+1}^*)^2 \right| \leq H^2 \sqrt{\frac{2\iota}{n}}$$

1031 Considering all the possible combinations $(s, a, h, n) \in \mathcal{S} \times \mathcal{A} \times [H] \times [\frac{T}{H}]$, with probability
 1032 at least $1 - SATp$, it holds simultaneously for all $(s, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$ that:

$$\frac{1}{n_h^k(s, a)} \left| \sum_{i=1}^{n_h^k} \left(\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s, a, h} \right) (V_{h+1}^*)^2 \right| \leq H^2 \sqrt{\frac{2\iota}{n_h^k(s, a)}}.$$

□

1033 From this lemma, we know that the event $\bigcap_{i=1}^7 \mathcal{E}_i$ holds with probability at least $1 - (40H^2SAT^4)p$.

1034 Next, we will discuss the relationship among the V -estimate V_h^k , the reference function $V_h^{\text{ref}, k}(s)$,
 1035 the surrogate function $\hat{V}_h^{\text{ref}, k}(s)$ and the final learned reference function $V_h^{\text{REF}}(s)$.

1036 **Lemma B.2.** *For any $(s, h, k) \in \mathcal{S} \times [H] \times [K]$, let . Under the event $\mathcal{E}_1 \cap \mathcal{E}_2$ in Lemma B.1, we
 1037 have the following conclusions:*

1038 (a) $\hat{V}_h^{\text{ref}, k}(s) = \min\{V_h^*(s) + \beta, V_h^{\text{ref}, k}(s)\}$

1039 (b) $0 \leq V_h^{\text{ref}, k}(s) - V_h^{\text{REF}}(s) \leq H\lambda_h^k(s)$.

1040 (c) $0 \leq V_h^{\text{ref}, k}(s) - \hat{V}_h^{\text{ref}, k}(s) \leq H\lambda_h^k(s)$.

1041 (d) $|\hat{V}_h^{\text{ref}, k}(s) - V_h^{\text{REF}}(s)| \leq H\lambda_h^k(s)$.

1042 *Proof.* (a) Under the event \mathcal{E}_1 in Lemma B.1, we have $V_h^{\text{ref}, k}(s) \geq V_h^k(s) \geq V_h^*(s)$. Therefore,

1043 $\min\{V_h^*(s) + \beta, V_h^{\text{ref}, k}(s)\} \geq V_h^*(s)$. According to the definition of $\hat{V}_h^{\text{ref}, k}(s)$, we have
 1044 $\hat{V}_h^{\text{ref}, k}(s) = \min\{V_h^*(s) + \beta, V_h^{\text{ref}, k}(s)\}$.

1045 (b) For any $(s, h, k) \in \mathcal{S} \times [H] \times [K]$:

1046 If $N_h^k(s) \geq N_0$, then $\lambda_h^k(s) = 0$. In this case, the reference function $V_h^{\text{ref}, k}(s)$ is updated to
 1047 its final value $V_h^{\text{REF}}(s)$ and then $V_h^{\text{ref}, k}(s) - V_h^{\text{REF}}(s) = 0 = H\lambda_h^k(s)$.

1048 If $N_h^k(s) < N_0$, then $\lambda_h^k(s) = 1$. Since the reference function is non-increasing and
 1049 $V_h^{\text{ref}, 1}(s) = H$, we have $0 \leq V_h^{\text{ref}, k}(s) - V_h^{\text{REF}}(s) \leq H = H\lambda_h^k(s)$.

1050 Combining these two cases, we can prove the conclusion (b).

1051 (c) For any $(s, h, k) \in \mathcal{S} \times [H] \times [K]$:

1052 If $N_h^k(s) \geq N_0$, then $\lambda_h^k(s) = 0$. Under the event \mathcal{E}_2 in Lemma B.1, we have $V_h^{\text{ref}, k}(s) \leq
 1053 V_h^*(s) + \beta$. Therefore, it holds that $\hat{V}_h^{\text{ref}, k}(s) = V_h^{\text{ref}, k}(s)$ by (a). In this case, $V_h^{\text{ref}, k}(s) -
 1054 \hat{V}_h^{\text{ref}, k}(s) = 0 = H\lambda_h^k(s)$.

1055 If $N_h^k(s) < N_0$, then $\lambda_h^k(s) = 1$. Since the reference function is non-increasing and
 1056 $V_h^{\text{ref}, 1}(s) = H$, we have $0 \leq V_h^{\text{ref}, k}(s) - \hat{V}_h^{\text{ref}, k}(s) \leq H$.

1057 Combining these two cases, we can prove the conclusion (c).

- 1080 (d) For any $(s, h, k) \in \mathcal{S} \times [H] \times [K]$:
1081 If $N_h^k(s) \geq N_0$, then $\lambda_h^k(s) = 0$. In this case, the reference function $V_h^{\text{ref},k}(s)$ is updated to
1082 its final value $V_h^{\text{REF}}(s)$. Under the event \mathcal{E}_2 in Lemma B.1, we have $V_h^{\text{REF}}(s) = V_h^{\text{ref},k}(s) \leq$
1083 $V_h^*(s) + \beta$. In this case, we know $\hat{V}_h^{\text{ref},k}(s) = V_h^{\text{ref},k}(s) = V_h^{\text{REF}}(s)$. Therefore, it holds
1084 that $\hat{V}_h^{\text{ref},k}(s) - V_h^{\text{REF}}(s) = 0 = H\lambda_h^k(s)$.
1085 If $N_h^k(s) < N_0$, then $\lambda_h^k(s) = 1$. Since the reference function is non-increasing and
1086 $V_h^{\text{ref},k}(s) = H$, we have $0 \leq V_h^{\text{REF}}(s) \leq V_h^{\text{ref},k}(s) \leq H$ and $0 \leq \hat{V}_h^{\text{ref},k}(s) \leq V_h^{\text{ref},k}(s) \leq$
1087 H . Therefore, it holds that $|\hat{V}_h^{\text{ref},k}(s) - V_h^{\text{REF}}(s)| \leq H = H\lambda_h^k(s)$.
1088 Combining these two cases, we can prove the conclusion (d).
1090

□

1092 **Lemma B.3.** For any $(s, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$ such that $\check{n}_h^k(s, a) \neq 0$, it holds that:
1093

$$\frac{n_h^k(s, a)}{\check{n}_h^k(s, a)} \leq 4H$$

1097 *Proof.* For $\check{n}_h^k(s, a) \neq 0$, there exists $j \in \mathbb{N}_+$ such that $\check{n}_h^k(s, a) = e_j$ and $n_h^k(s, a) = \sum_{i=1}^j e_i$. We
1098 will use the mathematical induction to prove that for any $j \in \mathbb{N}_+$, $\frac{\sum_{i=1}^j e_i}{e_j} \leq 4H$.
1099

1100 For $j = 1$, $\frac{\sum_{i=1}^j e_i}{e_j} = 1 \leq 4H$.

1101 If $\frac{\sum_{i=1}^{j-1} e_i}{e_{j-1}} \leq 4H$, then for $j \in \mathbb{N}_+$ and $j \geq 2$, because $e_j = \lfloor (1 + \frac{1}{H})e_{j-1} \rfloor \geq (1 + \frac{1}{2H})e_{j-1}$, we
1102 have:
1103

$$\frac{\sum_{i=1}^j e_i}{e_j} = 1 + \frac{\sum_{i=1}^{j-1} e_i}{e_j} \leq 1 + \frac{\sum_{i=1}^{j-1} e_i}{(1 + \frac{1}{2H})e_{j-1}} \leq 1 + \frac{4H}{1 + \frac{1}{2H}} \leq 4H.$$

1104 Therefore, we finish the proof.
1105

□

1106 **Lemma B.4.** For any non-negative weight sequence $\{\omega_{h,k}\}_{h,k}$ and $\alpha \in (0, 1)$, it holds that:
1107

$$\sum_{k=1}^K \frac{\omega_{h,k} \mathbb{I}[n_h^k(s_h^k, a_h^k) \neq 0]}{n_h^k(s_h^k, a_h^k)^\alpha} \leq \frac{2^{2-\alpha}}{1-\alpha} (SA \|\omega\|_{\infty,h})^\alpha \|\omega\|_{1,h}^{1-\alpha},$$

1108 and
1109

$$\sum_{k=1}^K \frac{\omega_{h,k} \mathbb{I}[\check{n}_h^k(s_h^k, a_h^k) \neq 0]}{\check{n}_h^k(s_h^k, a_h^k)^\alpha} \leq \frac{2^{2+\alpha} H^\alpha}{1-\alpha} (SA \|\omega\|_{\infty,h})^\alpha \|\omega\|_{1,h}^{1-\alpha}.$$

1110 Here, $\|\omega\|_{\infty,h} = \max_k \{\omega_{h,k}\}$ and $\|\omega\|_{1,h} = \sum_{k=1}^K \omega_{h,k}$.
1111

1112 For $\alpha = 1$, we have the following conclusions:
1113

$$\sum_{k=1}^K \frac{\mathbb{I}[n_h^k(s_h^k, a_h^k) \neq 0]}{n_h^k(s_h^k, a_h^k)} \leq 2SA \log(T),$$

1114 and
1115

$$\sum_{k=1}^K \frac{\mathbb{I}[\check{n}_h^k(s_h^k, a_h^k) \neq 0]}{\check{n}_h^k(s_h^k, a_h^k)} \leq 4SAH \log(T).$$

1116 *Proof.*
1117

$$\begin{aligned} 1118 \sum_{k=1}^K \frac{\omega_{h,k} \mathbb{I}[n_h^k(s_h^k, a_h^k) \neq 0]}{n_h^k(s_h^k, a_h^k)^\alpha} &= \sum_{s,a} \sum_{k=1}^K \frac{\omega_{h,k} \mathbb{I}[n_h^k(s, a) \neq 0, (s_h^k, a_h^k) = (s, a)]}{n_h^k(s, a)^\alpha} \\ 1119 &\triangleq \sum_{s,a} \sum_{k=1}^K \frac{\omega'_{h,k}(s, a)}{n_h^k(s, a)^\alpha} \end{aligned} \tag{37}$$

1134 Here we let $\omega'_{h,k}(s, a) = \omega_{h,k} \mathbb{I}[n_h^k(s, a) \neq 0, (s_h^k, a_h^k) = (s, a)]$ and $c_h(s, a) = \sum_{k=1}^K \omega'_{h,k}(s, a)$.
 1135 Then $\omega'_{h,k}(s, a) \leq \|\omega\|_{\infty, h}$ and $\sum_{s,a} c_h(s, a) \leq \sum_{k=1}^K \omega_{h,k} = \|\omega\|_{1,h}$.

1136
 1137 Because $n_h^k(s, a)$ is nondecreasing for $1 \leq k \leq K$, given the term $\sum_{k=1}^K \frac{\omega'_{h,k}}{n_h^k(s, a)^\alpha}$, when the weights
 1138 $\omega'_{h,k}(s, a)$ concentrates on former terms, we can obtain the largest value. For a given state-action
 1139 pair (s, a) and $j \in \mathbb{N}_+$, according to the stage design, the set $\{k : n_h^k(s, a) = \sum_{i=1}^j e_i\}$ has
 1140 at most $e_{j+1} \leq (1 + \frac{1}{H})e_j$ elements. Thus, the upper bound for the sum of the coefficients of
 1141 $n_h^k(s, a) = \sum_{i=1}^j e_i$ in Equation (37) is given by $(1 + \frac{1}{H})e_j \|\omega\|_{\infty, h}$.
 1142

1143 Let:

$$1144 \quad k_0 = \max \left\{ k : \sum_{j=1}^{k-1} \left(1 + \frac{1}{H}\right) e_j \|\omega\|_{\infty, h} < c_h(s, a), k \in \mathbb{N}_+ \right\}.$$

1145 Because $e_{j+1} \leq (1 + \frac{1}{H})e_j$ for any $j \in \mathbb{N}_+$, we have

$$1146 \quad \sum_{j=2}^{k_0} e_j \|\omega\|_{\infty, h} < c_h(s, a),$$

1147 and then k_0 satisfies

$$1148 \quad \sum_{j=1}^{k_0} e_j \|\omega\|_{\infty, h} \leq \sum_{j=1}^{k_0-1} \left(1 + \frac{1}{H}\right) e_j \|\omega\|_{\infty, h} + \sum_{j=2}^{k_0} e_j \|\omega\|_{\infty, h} < 2c_h(s, a). \quad (38)$$

1149 Therefore, back to Equation (37), concentrating the weight to the terms with $n_h^k(s, a) = \sum_{i=1}^j e_i$,
 1150 $j \in \{1, 2, \dots, k_0\}$, for any given state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, we have:

$$1151 \quad \sum_{k=1}^K \frac{\omega'_{h,k}}{n_h^k(s, a)^\alpha} \leq \sum_{j=1}^{k_0} \frac{(1 + \frac{1}{H})e_j \|\omega\|_{\infty, h}}{\left(\sum_{i=1}^j e_i\right)^\alpha} = (1 + \frac{1}{H}) \|\omega\|_{\infty, h} \left(\sum_{j=1}^{k_0} \frac{e_j}{\left(\sum_{i=1}^j e_i\right)^\alpha} \right). \quad (39)$$

1152 For any $0 < y < x$ and $\alpha \in (0, 1)$, we have:

$$1153 \quad \frac{x-y}{x^\alpha} \leq \frac{1}{1-\alpha} (x^{1-\alpha} - y^{1-\alpha}).$$

1154 For any $j \in \mathbb{N}_+$, let $x = \sum_{i=1}^j e_i$ and $y = \sum_{i=1}^{j-1} e_i$, then we have:

$$1155 \quad \frac{e_j}{\left(\sum_{i=1}^j e_i\right)^\alpha} \leq \frac{1}{1-\alpha} \left(\left(\sum_{i=1}^j e_i \right)^{1-\alpha} - \left(\sum_{i=1}^{j-1} e_i \right)^{1-\alpha} \right).$$

1156 Sum the above inequality from 1 to k_0 , then it holds that:

$$1157 \quad \sum_{j=1}^{k_0} \frac{e_j}{\left(\sum_{i=1}^j e_i\right)^\alpha} \leq \frac{1}{1-\alpha} \left(\sum_{i=1}^{k_0} e_i \right)^{1-\alpha} < \frac{1}{1-\alpha} \left(\frac{2c_h(s, a)}{\|\omega\|_{\infty, h}} \right)^{1-\alpha}.$$

1158 The last inequality is because of Equation (38). Applying this inequality to Equation (39), we have:
 1159

$$1160 \quad \sum_{k=1}^K \frac{\omega'_{h,k}}{n_h^k(s, a)^\alpha} \leq \frac{2^{2-\alpha}}{1-\alpha} \|\omega\|_{\infty, h}^\alpha c_h(s, a)^{1-\alpha}.$$

1161 Using this inequality in Equation (37), we have:

$$1162 \quad \sum_{k=1}^K \frac{\omega_{h,k} \mathbb{I}[n_h^k(s, a) \neq 0]}{\sqrt{n_h^k(s, a)^\alpha}} \leq \frac{2^{2-\alpha}}{1-\alpha} \|\omega\|_{\infty, h}^\alpha \sum_{s,a} c_h(s, a)^{1-\alpha} \leq \frac{2^{2-\alpha}}{1-\alpha} (SA \|\omega\|_{\infty, h})^\alpha \|\omega\|_{1,h}^{1-\alpha}.$$

1188 The last inequality holds due to Hölder's inequality, as $\sum_{s,a} c_h(s,a)^{1-\alpha} \leq (SA)^\alpha \|\omega\|_{1,h}^{1-\alpha}$.
 1189

1190 By Lemma B.3, it is easy to prove the second conclusion:

1191
 1192
$$\sum_{k=1}^K \frac{\omega_{h,k} \mathbb{I}[\check{n}_h^k(s_h^k, a_h^k) \neq 0]}{\check{n}_h^k(s_h^k, a_h^k)^\alpha} \leq \frac{2^{2+\alpha} H^\alpha}{1-\alpha} (SA \|\omega\|_{\infty,h})^\alpha \|\omega\|_{1,h}^{1-\alpha}.$$

 1193

1194 The case of $\alpha = 1$ is proved in Lemma 11 of Zhang et al. (2020). \square
 1195

1196 **Lemma B.5.** For any non-negative sequence $\{X_h^k\}_{k,h}$, we have that
 1197

1198
 1199
$$\sum_{k=1}^K \frac{\mathbb{I}[n_h^k(s_h^k, a_h^k) \neq 0]}{n_h^k(s_h^k, a_h^k)} \sum_{i=1}^{n_h^k} X_h^{l_i} \leq 3 \log(T) \sum_{k=1}^K X_h^k,$$

 1200

1201
 1202
$$\sum_{k=1}^K \frac{\mathbb{I}[\check{n}_h^k(s_h^k, a_h^k) \neq 0]}{\check{n}_h^k(s_h^k, a_h^k)} \sum_{i=1}^{\check{n}_h^k} X_h^{\check{l}_i} \leq \left(1 + \frac{1}{H}\right) \sum_{k=1}^K X_h^k.$$

 1203

1204
 1205 *Proof.* For the first conclusion,
 1206

1207
 1208
$$\sum_{k=1}^K \frac{\mathbb{I}[n_h^k(s_h^k, a_h^k) \neq 0]}{n_h^k(s_h^k, a_h^k)} \sum_{i=1}^{n_h^k} X_h^{l_i} = \sum_{k=1}^K \frac{\sum_{i=1}^{n_h^k} X_h^{l_i}}{n_h^k(s_h^k, a_h^k)} \cdot \sum_{j=1}^K \mathbb{I}[l_i = j, n_h^k(s_h^k, a_h^k) \neq 0]$$

 1209
 1210
 1211
$$= \sum_{k=1}^K \sum_{i=1}^{n_h^k} \sum_{j=1}^K \frac{X_h^j}{n_h^k(s_h^k, a_h^k)} \cdot \mathbb{I}[l_i = j, n_h^k(s_h^k, a_h^k) \neq 0]$$

 1212
 1213
 1214
$$= \sum_{j=1}^K \left(\sum_{k=1}^K \frac{\sum_{i=1}^{n_h^k} \mathbb{I}[l_i = j, n_h^k(s_h^k, a_h^k) \neq 0]}{n_h^k(s_h^k, a_h^k)} \right) X_h^j. \quad (40)$$

 1215
 1216

1217 For a given episode k , according to the definition of l_i , $\sum_{i=1}^{n_h^k} \mathbb{I}[l_i = j, n_h^k(s_h^k, a_h^k) \neq 0] = 1$ if and
 1218 only if $(s_h^k, a_h^k) = (s_h^j, a_h^j)$ and (j, h) falls in the stage before that (k, h) falls in. As a result, for
 1219 $n_h^k(s_h^k, a_h^k) = \sum_{i=1}^{j-1} e_i$, the set $\{k : \sum_{i=1}^{n_h^k} \mathbb{I}[l_i = j, n_h^k(s_h^k, a_h^k) \neq 0] = 1\}$ has at most e_j elements.
 1220 Then it holds that:
 1221

1222
 1223
$$\sum_{k=1}^K \frac{\sum_{i=1}^{n_h^k} \mathbb{I}[l_i = j, n_h^k(s_h^k, a_h^k) \neq 0]}{n_h^k(s_h^k, a_h^k)} \leq \sum_{j \in A} \frac{e_j}{\sum_{i=1}^{j-1} e_i} \leq \sum_{j \in A} \sum_{p=1}^{e_j} \frac{3}{\sum_{i=1}^{j-1} e_i + p} \leq 3 \log(T) \quad (41)$$

 1224
 1225

1226 Here, $A = \{j : H \leq \sum_{i=1}^{j-1} e_i \leq T, j \in \mathbb{N}_+\}$. The second inequality is because $e_j \leq (1 + \frac{1}{H})e_{j-1}$
 1227 and then for any $p \in [e_j]$, $\sum_{i=1}^{j-1} e_i + p \leq 3 \sum_{i=1}^{j-1} e_i$. Then we finish the proof of the first conclusion.
 1228 For the second conclusion,
 1229

1230
 1231
$$\sum_{k=1}^K \frac{\mathbb{I}[\check{n}_h^k(s_h^k, a_h^k) \neq 0]}{\check{n}_h^k(s_h^k, a_h^k)} \sum_{i=1}^{\check{n}_h^k} X_h^{\check{l}_i} = \sum_{k=1}^K \frac{\sum_{i=1}^{\check{n}_h^k} X_h^{\check{l}_i}}{\check{n}_h^k(s_h^k, a_h^k)} \cdot \sum_{j=1}^K \mathbb{I}[\check{l}_i = j, \check{n}_h^k(s_h^k, a_h^k) \neq 0]$$

 1232
 1233
 1234
$$= \sum_{k=1}^K \sum_{i=1}^{\check{n}_h^k} \sum_{j=1}^K \frac{X_h^j}{\check{n}_h^k(s_h^k, a_h^k)} \cdot \mathbb{I}[\check{l}_i = j, \check{n}_h^k(s_h^k, a_h^k) \neq 0]$$

 1235
 1236
 1237
$$= \sum_{j=1}^K \left(\sum_{k=1}^K \frac{\sum_{i=1}^{\check{n}_h^k} \mathbb{I}[\check{l}_i = j, \check{n}_h^k(s_h^k, a_h^k) \neq 0]}{\check{n}_h^k(s_h^k, a_h^k)} \right) X_h^j. \quad (42)$$

 1238
 1239

1240 For a given episode k , according to the definition of \check{l}_i , $\sum_{i=1}^{\check{n}_h^k} \mathbb{I}[\check{l}_i = j, \check{n}_h^k(s_h^k, a_h^k) \neq 0] = 1$ if
 1241 only if $(s_h^k, a_h^k) = (s_h^j, a_h^j)$ and (j, h) falls in the previous stage of that (k, h) falls in.

As a result, in the stage of (j, h) , the number of visits to (s_h^k, a_h^k, h) is $\check{n}_h^k(s_h^k, a_h^k)$, and the set $\{k : \sum_{i=1}^{n_h^k} \mathbb{I}[\check{l}_i = j, \check{n}_h^k \neq 0] = 1\}$ has at most $(1 + \frac{1}{H})\check{n}_h^k(s_h^k, a_h^k)$ elements. Then it holds that:

$$\sum_{k=1}^K \frac{\sum_{i=1}^{\check{n}_h^k} \mathbb{I}[\check{l}_i = j, \check{n}_h^k \neq 0]}{\check{n}_h^k(s_h^k, a_h^k)} \leq 1 + \frac{1}{H} \quad (43)$$

Therefore, we prove the second conclusion. \square

B.3 PROOF SKETCH OF THEOREM 3.1

Next, we will begin to prove Theorem 3.1 under $\bigcap_{i=1}^7 \mathcal{E}_i$.

Step 1: Bounding the term $Q_h^k - Q_h^*$. By Equation (31) and Bellman Optimality Equation (5), it holds that:

$$\begin{aligned} & Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \\ & \leq \mathbb{I}[n_h^k \neq 0] \left(\frac{\sum_{i=1}^{n_h^k} V_{h+1}^{\text{ref}, l_i}(s_{h+1}^{l_i})}{n_h^k(s_h^k, a_h^k)} + \frac{\sum_{i=1}^{\check{n}_h^k} (V_{h+1}^{\check{l}_i} - V_{h+1}^{\text{ref}, \check{l}_i})(s_{h+1}^{\check{l}_i})}{\check{n}_h^k(s_h^k, a_h^k)} + b_h^k(s_h^k, a_h^k) \right) \\ & \quad + \mathbb{I}[n_h^k = 0] H - \mathbb{P}_{s_h^k, a_h^k, h} V_{h+1}^* \\ & \leq \mathbb{I}[n_h^k \neq 0] \left(\frac{\sum_{i=1}^{n_h^k} V_{h+1}^{\text{ref}, l_i}(s_{h+1}^{l_i})}{n_h^k(s_h^k, a_h^k)} + \frac{\sum_{i=1}^{\check{n}_h^k} (V_{h+1}^{\check{l}_i} - V_{h+1}^{\text{REF}})(s_{h+1}^{\check{l}_i})}{\check{n}_h^k(s_h^k, a_h^k)} + b_h^k(s_h^k, a_h^k) \right) \\ & \quad + \mathbb{I}[n_h^k = 0] H - \mathbb{I}[\check{n}_h^k \neq 0] \mathbb{P}_{s_h^k, a_h^k, h} V_{h+1}^* \\ & = \mathbb{I}[n_h^k = 0] H + \mathbb{I}[n_h^k \neq 0] (G_1 + b_h^k(s_h^k, a_h^k)) + \mathbb{I}[\check{n}_h^k \neq 0] (G_2 + G_3) \end{aligned}$$

The second inequality is because $V_{h+1}^{\text{ref}, \check{l}_i}(s_{h+1}^{\check{l}_i}) \geq V_{h+1}^{\text{REF}}(s_{h+1}^{\check{l}_i})$. In the last equality we use $\mathbb{I}[n_h^k(s_h^k, a_h^k) = 0] = \mathbb{I}[\check{n}_h^k(s_h^k, a_h^k) = 0]$. Here

$$\begin{aligned} G_1 &= \frac{\sum_{i=1}^{n_h^k} \left(V_{h+1}^{\text{ref}, l_i}(s_{h+1}^{l_i}) - \mathbb{P}_{s_h^k, a_h^k, h} V_{h+1}^{\text{REF}} \right)}{n_h^k(s_h^k, a_h^k)}, \\ G_2 &= \frac{\sum_{i=1}^{\check{n}_h^k} \left(\mathbb{P}_{s_h^k, a_h^k, h} - \mathbb{I}_{s_{h+1}^{l_i}} \right) (V_{h+1}^{\text{REF}} - V_{h+1}^*)}{\check{n}_h^k(s_h^k, a_h^k)}, \\ G_3 &= \frac{\sum_{i=1}^{\check{n}_h^k} \left(V_{h+1}^{\check{l}_i}(s_{h+1}^{\check{l}_i}) - V_{h+1}^*(s_{h+1}^{\check{l}_i}) \right)}{\check{n}_h^k(s_h^k, a_h^k)}. \end{aligned}$$

The upper bounds of G_1 , G_2 and b_h^k is given in Appendix B.4. Combining the three upper bounds Equation (54), Equation (58) and Equation (63), the following inequality holds:

$$(Q_h^k - Q_h^*)(s_h^k, a_h^k) \lesssim \mathbb{I}[\check{n}_h^k \neq 0] \left(G_3 + \frac{H\iota^{\frac{3}{4}}}{\check{n}_h^k(s_h^k, a_h^k)^{\frac{3}{4}}} \right) + \mathbb{I}[n_h^k \neq 0] \sqrt{\frac{(\mathbb{Q}^* + \beta^2 H)\iota}{n_h^k(s_h^k, a_h^k)}} + Y_h^k. \quad (44)$$

Here, for any $h' \in [H]$ and $k \in [K]$, $Y_{h'}^k$ is defined as:

$$\begin{aligned} Y_{h'}^k &= H\mathbb{I}[n_{h'}^k = 0] + \frac{\mathbb{I}[n_{h'}^k \neq 0]}{n_{h'}^k(s_{h'}^k, a_{h'}^k)} \left(\sum_{i=1}^{n_{h'}^k} H \left(\mathbb{I}_{s_{h'+1}^{l_i}} + \mathbb{P}_{s_{h'}^k, a_{h'}^k, h'} \right) \lambda_{h'+1}^{l_i} + \sqrt{H\Gamma_{h'}^k(s_{h'}^k, a_{h'}^k)\iota} \right) \\ & \quad + \frac{\mathbb{I}[\check{n}_{h'}^k \neq 0]}{\check{n}_{h'}^k(s_{h'}^k, a_{h'}^k)} \left(\sum_{i=1}^{\check{n}_{h'}^k} H \left(\mathbb{P}_{s_{h'}^k, a_{h'}^k, h'} + \mathbb{I}_{s_{h'+1}^{l_i}} \right) \lambda_{h'+1}^{l_i} + H\sqrt{\check{\Gamma}_{h'}^k(s_{h'}^k, a_{h'}^k)\iota} + H\iota \right). \end{aligned}$$

1296 Then in Equation (44), for the given h , $W_h = \sum_{k=1}^K \omega_{h,k} Y_h^k$.
 1297

1298 **Step 2: Bounding the weighted sum.** For any given h and non-negative constants $\{\omega_{h,k}\}_{h,[K]}$, we
 1299 denote $\|\omega\|_{\infty,h} = \max_{k \in [K]} \omega_{h,k}$ and $\|\omega\|_{1,h} = \sum_{k \in [K]} \omega_{h,k}$. We also recursively define $\omega_{h',k}(h)$
 1300 for any $h \leq h' \leq H$, $k \in [K]$ as follows:

$$1301 \quad \omega_{h,k}(h) := \omega_{h,k}; \quad \omega_{h'+1,j}(h) = \sum_{k=1}^K \omega_{h',k}(h) \frac{\sum_{i=1}^{\check{n}_{h'}^k} \mathbb{I}[\check{l}_i = j, \check{n}_{h'}^k \neq 0]}{\check{n}_{h'}^k(s_h^k, a_h^k)}. \quad (45)$$

1304 By Equation (43), it is easy to show that
 1305

$$1306 \quad \|\omega(h)\|_{1,h'+1} \leq \|\omega(h)\|_{1,h'+1}, \quad \|\omega(h)\|_{\infty,h'} \leq (1 + 1/H) \|\omega(h)\|_{\infty,h'}, \forall h' > h. \quad (46)$$

1308 Given the weight $\{\omega_{h,k}\}$, we will bound $\sum_{k=1}^K \omega_{h,k}(Q_h^k - Q_h^*)(s_h^k, a_h^k)$. With Equation (44), we
 1309 have:

$$\begin{aligned} 1310 \quad & \sum_{k=1}^K \omega_{h,k}(Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)) \\ 1311 \quad & \leq \sum_{k=1}^K \omega_{h,k} \mathbb{I}[\check{n}_h^k \neq 0] G_3 + \sum_{k=1}^K \omega_{h,k} \left(\mathbb{I}[n_h^k \neq 0] \sqrt{\frac{(\mathbb{Q}^* + \beta^2 H) \iota}{n_h^k(s_h^k, a_h^k)}} + \mathbb{I}[\check{n}_h^k \neq 0] \frac{H \iota^{3/4}}{\check{n}_h^k(s_h^k, a_h^k)^{3/4}} \right) \\ 1312 \quad & + \sum_{k=1}^K \omega_{h,k} Y_h^k. \\ 1313 \quad & \lesssim \sum_{j=1}^K \omega_{h+1,j}(h) (Q_{h+1}^j - Q_{h+1}^*) (s_{h+1}^j, a_{h+1}^j) + \sqrt{(\mathbb{Q}^* + \beta^2 H) S A \|\omega\|_{\infty,h} \|\omega\|_{1,h} \iota} \\ 1314 \quad & + H^{11/4} (S A \|\omega\|_{\infty,h} \iota)^{3/4} \|\omega\|_{1,h}^{1/4} + \sum_{k=1}^K \omega_{h,k} Y_h^k. \end{aligned} \quad (47)$$

1325 In this inequality, the upper bound of $\sum_{k=1}^K \omega_{h,k} \mathbb{I}[\check{n}_h^k \neq 0] G_3$ is given in Appendix B.5. The upper
 1326 bounds of middle two terms is given by Lemma B.4 with $\alpha = \frac{1}{2}$ and $\alpha = \frac{3}{4}$.
 1327

1328 Recurring Equation (47) for $h, h+1, \dots, H$, since $Q_{H+1}^k(s, a) = Q_{H+1}^*(s, a) = 0$ and the weight
 1329 relationship Equation (45) and Equation (46), we have:

$$\begin{aligned} 1330 \quad & \sum_{k=1}^K \omega_{h,k}(Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)) \\ 1331 \quad & \lesssim H \sqrt{(\mathbb{Q}^* + \beta^2 H) S A \|\omega\|_{\infty,h} \|\omega\|_{1,h} \iota} + H^{11/4} (S A \|\omega\|_{\infty,h} \iota)^{3/4} \|\omega\|_{1,h}^{1/4} + \sum_{h'=h}^H \sum_{k=1}^K \omega_{h',k}(h) Y_{h'}^k. \end{aligned} \quad (48)$$

1338 **Step 3: Integrating multiple weighted sums.** For any $N = \lceil \log_2(H/\Delta_{\min}) \rceil$, $n \in [N]$, $k \in [K]$
 1339 and the given $h \in [H]$, let:

$$1340 \quad \omega_{h,k}^{(n)} = \mathbb{I}[Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \in [2^{n-1} \Delta_{\min}, 2^n \Delta_{\min})],$$

1342 and

$$1343 \quad \|\omega^{(n)}\|_{\infty,h} = \max_k \omega_{h,k}^{(n)} \leq 1, \quad \|\omega^{(n)}\|_{1,h} = \sum_{k=1}^K \omega_{h,k}^{(n)}.$$

1346 For $h \leq h' \leq H$ and any $n \in [N]$ the weight $\{\omega_{h',k}^{(n)}\}_k$ can be defined recursively by Equation (45):
 1347

$$1348 \quad \omega_{h,j}^{(n)}(h) = \omega_{h,j}^{(n)}; \quad \omega_{h'+1,j}^{(n)}(h) = \sum_{k=1}^K \omega_{h',k}^{(n)}(h) \frac{\sum_{i=1}^{\check{n}_{h'}^k} \mathbb{I}[\check{l}_i = j, \check{n}_{h'}^k \neq 0]}{\check{n}_{h'}^k(s_h^k, a_h^k)}.$$

Therefore, for any $j \in [K]$, it holds that:

$$\sum_{n=1}^N \omega_{h'+1,j}^{(n)}(h) = \sum_{k=1}^K \left(\sum_{n=1}^N \omega_{h',k}^{(n)}(h) \right) \frac{\sum_{i=1}^{\check{n}_{h'}^k} \mathbb{I}[\check{l}_i = j, \check{n}_{h'}^k \neq 0]}{\check{n}_{h'}^k(s_{h'}^k, a_{h'}^k)}.$$

Then by mathematical induction on $h' \in [h, H]$, it is straightforward to prove that for any $j \in [K]$,

$$\sum_{n=1}^N \omega_{h',j}^{(n)}(h) \leq \left(1 + \frac{1}{H} \right)^{h'-h} < 3, \quad (49)$$

given that for any $j \in [K]$

$$\sum_{k=1}^K \frac{\sum_{i=1}^{\check{n}_{h'}^k} \mathbb{I}[\check{l}_i = j, \check{n}_{h'}^k \neq 0]}{\check{n}_{h'}^k(s_{h'}^k, a_{h'}^k)} \leq 1 + \frac{1}{H}$$

by Equation (43) and $\sum_{n=1}^N \omega_{h,j}^{(n)}(h) = \sum_{n=1}^N \omega_{h,j}^{(n)} \leq 1$.

Applying the weight $\{\omega_{h,k}^{(n)}\}_k$ to Equation (48), for any $n \in [N]$, it holds that:

$$\begin{aligned} & \sum_{k=1}^K \omega_{h,k}^{(n)}(Q_h^k(s_h^k, a_h^k) - Q_h^\star(s_h^k, a_h^k)) \\ & \lesssim H \sqrt{(\mathbb{Q}^\star + \beta^2 H) S A \|\omega^{(n)}\|_{1,h} \iota} + H^{\frac{11}{4}} (S A \iota)^{\frac{3}{4}} \|\omega^{(n)}\|_{1,h}^{\frac{1}{4}} + \sum_{h'=h}^H \sum_{k=1}^K \omega_{h',k}^{(n)}(h) Y_{h'}^k. \end{aligned}$$

On the other hand, according to the definition of $\omega_{h,k}^{(n)}$,

$$\sum_{k=1}^K \omega_{h,k}^{(n)}(Q_h^k(s_h^k, a_h^k) - Q_h^\star(s_h^k, a_h^k)) \geq 2^{n-1} \Delta_{\min} \|\omega\|_{1,h}^{(n)}.$$

Therefore, we obtain the following inequality:

$$\begin{aligned} & 2^{n-1} \Delta_{\min} \|\omega\|_{1,h}^{(n)} \\ & \lesssim H \sqrt{(\mathbb{Q}^\star + \beta^2 H) S A \|\omega^{(n)}\|_{1,h} \iota} + H^{\frac{11}{4}} (S A \iota)^{\frac{3}{4}} (\|\omega\|_{1,h}^{(n)})^{\frac{1}{4}}, \end{aligned} \quad (50)$$

Then at least one of the following three inequalities holds:

$$\begin{aligned} & 2^{n-1} \Delta_{\min} \|\omega\|_{1,h}^{(n)} \lesssim H \sqrt{(\mathbb{Q}^\star + \beta^2 H) S A \|\omega^{(n)}\|_{1,h} \iota}, \\ & 2^{n-1} \Delta_{\min} \|\omega^{(n)}\|_{1,h} \lesssim H^{\frac{11}{4}} (S A \iota)^{\frac{3}{4}} (\|\omega\|_{1,h}^{(n)})^{\frac{1}{4}}, \\ & 2^{n-1} \Delta_{\min} \|\omega\|_{1,h}^{(n)} \lesssim \sum_{h'=h}^H \sum_{k=1}^K \omega_{h',k}^{(n)}(h) Y_{h'}^k. \end{aligned}$$

Solving this three inequalities, we know that:

$$\begin{aligned} \|\omega\|_{1,h}^{(n)} & \leq O \left(\max \left\{ \frac{(\mathbb{Q}^\star + \beta^2 H) S A H^2 \iota}{4^{n-2} \Delta_{\min}^2}, \frac{H^{\frac{11}{3}} S A \iota}{(2^{n-1} \Delta_{\min})^{\frac{4}{3}}}, \frac{\sum_{h'=h}^H \sum_{k=1}^K \omega_{h',k}^{(n)}(h) Y_{h'}^k}{2^{n-1} \Delta_{\min}} \right\} \right) \\ & \leq O \left(\frac{(\mathbb{Q}^\star + \beta^2 H) S A H^2 \iota}{4^{n-2} \Delta_{\min}^2} + \frac{H^{\frac{11}{3}} S A \iota}{(2^{n-1} \Delta_{\min})^{\frac{4}{3}}} + \frac{\sum_{h'=h}^H \sum_{k=1}^K \omega_{h',k}^{(n)}(h) Y_{h'}^k}{2^{n-1} \Delta_{\min}} \right). \end{aligned}$$

By Equation (49), we have:

$$\sum_{n=1}^N \sum_{h'=h}^H \sum_{k=1}^K \omega_{h',k}^{(n)}(h) Y_{h'}^k = \sum_{h'=h}^H \sum_{k=1}^K \left(\sum_{n=1}^N \omega_{h',k}^{(n)}(h) \right) Y_{h'}^k \leq 3 \sum_{h'=1}^H \sum_{k=1}^K Y_{h'}^k.$$

1404 Therefore,

$$1406 \sum_{n=1}^N 2^n \Delta_{\min} \|\omega\|_{1,h}^{(n)} \leq O \left(\frac{(\mathbb{Q}^* + \beta^2 H) SAH^2 \iota}{\Delta_{\min}} + \frac{H^{\frac{11}{3}} SA \iota}{(\Delta_{\min})^{\frac{1}{3}}} + \sum_{h'=1}^H \sum_{k=1}^K Y_{h'}^k \right). \quad (51)$$

1409 From Appendix B.6, we know the upper bound for $\sum_{h'=1}^H \sum_{k=1}^K Y_{h'}^k$ is $O(\frac{H^7 S^2 A \iota \log(T)}{\beta^2})$. Therefore, back to Equation (51), it holds that:

$$\begin{aligned} 1412 \sum_{n=1}^N 2^n \Delta_{\min} \|\omega\|_{1,h}^{(n)} &\leq O \left(\frac{(\mathbb{Q}^* + \beta^2 H) SAH^2 \iota}{\Delta_{\min}} + \frac{H^{\frac{11}{3}} SA \iota}{(\Delta_{\min})^{\frac{1}{3}}} + \frac{H^7 S^2 A \iota \log(T)}{\beta^2} \right) \\ 1415 &\leq O \left(\frac{(\mathbb{Q}^* + \beta^2 H) SAH^2 \iota}{\Delta_{\min}} + \frac{H^7 S^2 A \iota \log(T)}{\beta^2} \right) \end{aligned} \quad (52)$$

1418 The last inequality is because:

$$1420 \frac{H^{\frac{11}{3}} SA \iota}{(\Delta_{\min})^{\frac{1}{3}}} \lesssim \frac{\beta^2 H^3 SA \iota}{\Delta_{\min}} + \frac{H^4 SA \iota}{\beta} + \frac{H^4 SA \iota}{\beta} \lesssim \frac{(\mathbb{Q}^* + \beta^2 H) SAH^2 \iota}{\Delta_{\min}} + \frac{H^7 SA \iota \log(T)}{\beta^2}.$$

1422 **Step 4: Bounding the expected gap-dependent regret.** Let $p = (40SAH^2T^5)^{-1}$, then $\mathcal{E} = \bigcap_{i=1}^7 \mathcal{E}_i$ holds with probability at least $1 - \frac{1}{T}$ and $\iota \lesssim \log(SAT)$. Therefore, by Equation (29), we have:

$$\begin{aligned} 1426 \mathbb{E}(\text{Regret}(K)) &\leq \mathbb{E} \left[\sum_{k=1}^K \sum_{h=1}^H \text{clip}[(Q_h^k - Q_h^*)(s_h^k, a_h^k) \mid \Delta_{\min}] \right] \\ 1429 &= \mathbb{E} \left[\sum_{k=1}^K \sum_{h=1}^H \text{clip}[(Q_h^k - Q_h^*)(x_h^k, a_h^k) \mid \Delta_{\min}] \middle| \mathcal{E} \right] \mathbb{P}(\mathcal{E}) \\ 1432 &\quad + \mathbb{E} \left[\sum_{k=1}^K \sum_{h=1}^H \text{clip}[(Q_h^k - Q_h^*)(x_h^k, a_h^k) \mid \Delta_{\min}] \middle| \mathcal{E}^c \right] \mathbb{P}(\mathcal{E}^c) \\ 1435 &\leq \sum_{h=1}^H \sum_{n=1}^N 2^n \Delta_{\min} \|\omega\|_{1,h}^{(n)} + \frac{1}{T} \cdot TH \\ 1438 &\leq O \left(\frac{(\mathbb{Q}^* + \beta^2 H) H^3 SA \log(SAT)}{\Delta_{\min}} + \frac{H^8 S^2 A \log(SAT) \log(T)}{\beta^2} \right). \end{aligned}$$

1441 The third inequality is because

$$\begin{aligned} 1443 \sum_{k=1}^K \text{clip}[(Q_h^k - Q_h^*)(s_h^k, a_h^k) \mid \Delta_{\min}] &= \sum_{k=1}^K \sum_{n=1}^N \omega_{h,k}^{(n)} (Q_h^k - Q_h^*)(s_h^k, a_h^k) \\ 1446 &\leq \sum_{n=1}^N 2^n \Delta_{\min} \sum_{k=1}^K \omega_{h,k}^{(n)} = \sum_{n=1}^N 2^n \Delta_{\min} \|\omega\|_{1,h}^{(n)}. \end{aligned}$$

1448 The last inequality is by Equation (52).

1450 B.4 BOUNDING THE TERM $Q_h^k - Q_h^*$

1452 B.4.1 BOUNDING THE TERM G_1

1454 We can split G_1 into four terms:

$$1456 \frac{\sum_{i=1}^{n_h^k} \left(V_{h+1}^{\text{ref}, l_i}(s_{h+1}^{l_i}) - \mathbb{P}_{s_h^k, a_h^k, h} V_{h+1}^{\text{REF}} \right)}{n_h^k(s_h^k, a_h^k)} = G_{1,1} + G_{1,2} + G_{1,3} + G_{1,4}, \quad (53)$$

1458 where

$$G_{1,1} = \frac{\sum_{i=1}^{n_h^k} \left(\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s_h^k, a_h^k, h} \right) \left(V_{h+1}^{\text{ref}, l_i} - \hat{V}_{h+1}^{\text{ref}, l_i} \right)}{n_h^k(s_h^k, a_h^k)},$$

$$G_{1,2} = \frac{\sum_{i=1}^{n_h^k} \left(\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s_h^k, a_h^k, h} \right) \left(\hat{V}_{h+1}^{\text{ref}, l_i} - V_{h+1}^* \right)}{n_h^k(s_h^k, a_h^k)},$$

$$G_{1,3} = \frac{\sum_{i=1}^{n_h^k} \left(\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s_h^k, a_h^k, h} \right) V_{h+1}^*}{n_h^k(s_h^k, a_h^k)}$$

1468 and

$$G_{1,4} = \frac{\sum_{i=1}^{n_h^k} \mathbb{P}_{s_h^k, a_h^k, h} \left(V_{h+1}^{\text{ref}, l_i} - V_{h+1}^{\text{REF}} \right)}{n_h^k(s_h^k, a_h^k)}.$$

1471 According to (c) in Lemma B.2, we have:

$$G_{1,1} \leq \frac{\sum_{i=1}^{n_h^k} H \left(\mathbb{1}_{s_{h+1}^{l_i}} + \mathbb{P}_{s_h^k, a_h^k, h} \right) \lambda_{h+1}^{l_i}}{n_h^k(s_h^k, a_h^k)}.$$

1476 Under the event \mathcal{E}_4 in Lemma B.1, we can bound $G_{1,2}$:

$$G_{1,2} \leq \beta \sqrt{\frac{2\iota}{n_h^k(s_h^k, a_h^k)}}.$$

1479 Under the event \mathcal{E}_5 in Lemma B.1, we can bound $G_{1,3}$:

$$G_{1,3} \leq 2 \sqrt{\frac{2Q^* \iota}{n_h^k(s_h^k, a_h^k)}} + \frac{4H\iota}{n_h^k(s_h^k, a_h^k)}.$$

1484 The upper bound of $G_{1,4}$ is given by (b) in Lemma B.2:

$$G_{1,4} \leq \frac{\sum_{i=1}^{n_h^k} H \mathbb{P}_{s_h^k, a_h^k, h} \lambda_{h+1}^{l_i}}{n_h^k(s_h^k, a_h^k)}.$$

1487 Combining these four upper bounds together, we can bound G_1 :

$$G_1 \lesssim \frac{\sum_{i=1}^{n_h^k} H \left(\mathbb{1}_{s_{h+1}^{l_i}} + 2\mathbb{P}_{s_h^k, a_h^k, h} \right) \lambda_{h+1}^{l_i}}{n_h^k(s_h^k, a_h^k)} + \sqrt{\frac{(Q^* + \beta^2)\iota}{n_h^k(s_h^k, a_h^k)}} + \frac{H\iota}{n_h^k(s_h^k, a_h^k)}. \quad (54)$$

1492 BOUNDING THE TERM G_2

1493 We can split the term of G_2 into two terms:

$$G_2 = \frac{\sum_{i=1}^{\check{n}_h^k} \left(\mathbb{P}_{s_h^k, a_h^k, h} - \mathbb{1}_{s_{h+1}^{l_i}} \right) \left[\left(V_{h+1}^{\text{REF}} - \hat{V}_{h+1}^{\text{ref}, l_i} \right) + \left(\hat{V}_{h+1}^{\text{ref}, l_i} - V_{h+1}^* \right) \right]}{\check{n}_h^k(s_h^k, a_h^k)}. \quad (55)$$

1498 According to (d) in Lemma B.2, we can bound the first term in Equation (55):

$$\frac{\sum_{i=1}^{\check{n}_h^k} \left(\mathbb{P}_{s_h^k, a_h^k, h} - \mathbb{1}_{s_{h+1}^{l_i}} \right) \left(V_{h+1}^{\text{REF}} - \hat{V}_{h+1}^{\text{ref}, l_i} \right)}{\check{n}_h^k(s_h^k, a_h^k)} \leq \frac{\sum_{i=1}^{\check{n}_h^k} H \left(\mathbb{P}_{s_h^k, a_h^k, h} + \mathbb{1}_{s_{h+1}^{l_i}} \right) \lambda_{h+1}^{l_i}}{\check{n}_h^k(s_h^k, a_h^k)}. \quad (56)$$

1503 The upper bound for the second term in Equation (55) is given by the event \mathcal{E}_6 in Lemma B.1:

$$\frac{\sum_{i=1}^{\check{n}_h^k} \left(\mathbb{P}_{s, a, h} - \mathbb{1}_{s_{h+1}^{l_i}} \right) \left(\hat{V}_{h+1}^{\text{ref}, l_i} - V_{h+1}^* \right)}{\check{n}_h^k(s, a)} \leq \sqrt{\frac{2\beta^2\iota}{\check{n}_h^k(s, a)}} \lesssim \sqrt{\frac{\beta^2 H\iota}{n_h^k(s, a)}}. \quad (57)$$

1508 The last inequality is because of Lemma B.3. Applying Equation (56) and Equation (57) to Equation (55), we have:

$$G_2 \lesssim \frac{\sum_{i=1}^{\check{n}_h^k} H \left(\mathbb{P}_{s_h^k, a_h^k, h} + \mathbb{1}_{s_{h+1}^{l_i}} \right) \lambda_{h+1}^{l_i}}{\check{n}_h^k(s_h^k, a_h^k)} + \sqrt{\frac{\beta^2 H\iota}{n_h^k(s, a)}}. \quad (58)$$

1512 B.4.3 BOUNDING THE TERM $b_h^k(s_h^k, a_h^k)$
 1513

1514 According to the definition of $b_h^k(s_h^k, a_h^k)$ in the algorithm, we have
 1515

$$1516 \quad b_h^k(s_h^k, a_h^k) = 2\sqrt{\frac{\nu_h^{\text{ref},k}\iota}{n_h^k}} + 2\sqrt{\frac{\check{\nu}_h^k\iota}{\check{n}_h^k}} + 5\left(\frac{H\iota}{n_h^k} + \frac{H\iota}{\check{n}_h^k} + \frac{H\iota^{\frac{3}{4}}}{(n_h^k)^{\frac{3}{4}}} + \frac{H\iota^{\frac{3}{4}}}{(\check{n}_h^k)^{\frac{3}{4}}}\right), \quad (59)$$

1519 where $\nu_h^{\text{ref},k} = \sigma_h^{\text{ref},k}/n_h^k - (\mu_h^{\text{ref},k}/n_h^k)^2$ and $\check{\nu}_h^k = \check{\sigma}_h^k/\check{n}_h^k - (\check{\mu}_h^k/\check{n}_h^k)^2$.
 1520

1521 Since $V_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i}) \geq \hat{V}_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i})$, it holds that
 1522

$$1523 \quad \sqrt{\frac{\nu_h^{\text{ref},k}\iota}{n_h^k}} = \sqrt{\frac{\frac{\sigma_h^{\text{ref},k}(s_h^k, a_h^k)}{n_h^k(s_h^k, a_h^k)} - \left(\frac{\mu_h^{\text{ref},k}(s_h^k, a_h^k)}{n_h^k(s_h^k, a_h^k)}\right)^2}{n_h^k(s_h^k, a_h^k)}}\iota \leq \sqrt{\frac{I_1^{h,k} + I_2^{h,k}}{n_h^k(s_h^k, a_h^k)}}\iota,$$

1526 where:

$$1527 \quad I_1^{h,k} = \frac{\sum_{i=1}^{n_h^k} \left(\left(V_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i})\right)^2 - \left(\hat{V}_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i})\right)^2 \right)}{n_h^k(s_h^k, a_h^k)},$$

1530 and

$$1531 \quad I_2^{h,k} = \frac{\sum_{i=1}^{n_h^k} \left(\hat{V}_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i}) \right)^2}{n_h^k(s_h^k, a_h^k)} - \left(\frac{\sum_{i=1}^{n_h^k} \hat{V}_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i})}{n_h^k(s_h^k, a_h^k)} \right)^2.$$

1535 Next we want to bound both $I_1^{h,k}$ and $I_2^{h,k}$.

$$1536 \quad I_1^{h,k} = \frac{\sum_{i=1}^{n_h^k} \left(V_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i}) + \hat{V}_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i}) \right) \left(V_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i}) - \hat{V}_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i}) \right)}{n_h^k(s_h^k, a_h^k)}$$

$$1539 \quad \leq \frac{\sum_{i=1}^{n_h^k} 2H \left(V_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i}) - \hat{V}_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i}) \right)}{n_h^k(s_h^k, a_h^k)} \triangleq \frac{2H\Gamma_h^k(s_h^k, a_h^k)}{n_h^k(s_h^k, a_h^k)}, \quad (60)$$

1542 where

$$1544 \quad \Gamma_h^k(s_h^k, a_h^k) = \sum_{i=1}^{n_h^k} \left(V_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i}) - \hat{V}_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i}) \right).$$

1547 For the second term $I_2^{h,k}$, because of Cauchy's Inequality, we have:

$$1548 \quad I_2^{h,k} = \frac{\sum_{n=1}^{n_h^k} \left(\hat{V}_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i}) - \frac{\sum_{i=1}^{n_h^k} \hat{V}_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i})}{n_h^k(s_h^k, a_h^k)} \right)^2}{n_h^k(s_h^k, a_h^k)} \leq 2 \left(I_{2,1}^{h,k} + I_{2,2}^{h,k} \right),$$

1552 where:

$$1554 \quad I_{2,1}^{h,k} = \frac{\sum_{i=1}^{n_h^k} \left(\hat{V}_{h+1}^{\text{ref},l_i}(s_{h+1}^{l_i}) - V_{h+1}^*(s_{h+1}^{l_i}) + \frac{\sum_{n=1}^{n_h^k} V_{h+1}^*(s_{h+1}^{l_n})}{n_h^k(s_h^k, a_h^k)} - \frac{\sum_{n=1}^{n_h^k} \hat{V}_{h+1}^{\text{ref},l_n}(s_{h+1}^{l_n})}{n_h^k(s_h^k, a_h^k)} \right)^2}{n_h^k(s_h^k, a_h^k)},$$

1558 and

$$1559 \quad I_{2,2}^{h,k} = \frac{\sum_{i=1}^{n_h^k} \left(V_{h+1}^*(s_{h+1}^{l_i}) - \frac{\sum_{n=1}^{n_h^k} V_{h+1}^*(s_{h+1}^{l_n})}{n_h^k(s_h^k, a_h^k)} \right)^2}{n_h^k(s_h^k, a_h^k)}$$

$$1564 \quad = \frac{\sum_{i=1}^{n_h^k} \left(V_{h+1}^*(s_{h+1}^{l_i}) \right)^2}{n_h^k(s_h^k, a_h^k)} - \left(\frac{\sum_{i=1}^{n_h^k} V_{h+1}^*(s_{h+1}^{l_i})}{n_h^k(s_h^k, a_h^k)} \right)^2.$$

1566 Since $V_{h+1}^*(s_{h+1}^{l_i}) \leq \hat{V}_{h+1}^{\text{ref}, l_i}(s_{h+1}^{l_i}) \leq V_{h+1}^*(s_{h+1}^{l_i}) + \beta$, it holds that:
 1567
 1568
$$\left| \hat{V}_{h+1}^{\text{ref}, l_i}(s_{h+1}^{l_i}) - V_{h+1}^*(s_{h+1}^{l_i}) + \frac{\sum_{n=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{l_n})}{n_h^k(s_h^k, a_h^k)} - \frac{\sum_{n=1}^{N_h^k} \hat{V}_{h+1}^{\text{ref}, l_n}(s_{h+1}^{l_n})}{n_h^k(s_h^k, a_h^k)} \right|$$

 1569
 1570
 1571
$$\leq \left| \hat{V}_{h+1}^{\text{ref}, l_i}(s_{h+1}^{l_i}) - V_{h+1}^*(s_{h+1}^{l_i}) \right| + \left| \frac{\sum_{n=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{l_n})}{n_h^k(s_h^k, a_h^k)} - \frac{\sum_{n=1}^{N_h^k} \hat{V}_{h+1}^{\text{ref}, l_n}(s_{h+1}^{l_n})}{n_h^k(s_h^k, a_h^k)} \right| \leq 2\beta.$$

 1572
 1573
 1574

Therefore, applying this inequality to $I_{2,1}^{h,k}$, we have $I_{2,1}^{h,k} \leq 4\beta^2$.

Moreover, according to the definition of \mathbb{Q}^* , it holds that

1575
 1576
$$I_{2,2}^{h,k} - \mathbb{Q}^* \leq I_{2,2}^{h,k} - \left(\mathbb{P}_{s_h^k, a_h^k, h} (V_{h+1}^*)^2 - \left(\mathbb{P}_{s_h^k, a_h^k, h} V_{h+1}^* \right)^2 \right)$$

 1577
 1578
$$= - \left(\frac{\sum_{i=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{l_i})}{n_h^k(s_h^k, a_h^k)} + \mathbb{P}_{s_h^k, a_h^k, h} V_{h+1}^* \right) \left(\frac{\sum_{i=1}^{N_h^k} (\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s_h^k, a_h^k, h}) V_{h+1}^*}{n_h^k(s_h^k, a_h^k)} \right)$$

 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588

$$+ \frac{\sum_{i=1}^{N_h^k} (\mathbb{1}_{s_{h+1}^{l_i}} - \mathbb{P}_{s_h^k, a_h^k, h}) (V_{h+1}^*)^2}{n_h^k(s_h^k, a_h^k)}$$

$$\lesssim H^2 \sqrt{\frac{\iota}{n_h^k(s_h^k, a_h^k)}}.$$

The last inequality is because of $\mathbb{Q}^* \leq H^2$, the event \mathcal{E}_5 and the event \mathcal{E}_7 in Lemma B.1. Therefore,

$$I_{2,2}^{h,k} \lesssim \mathbb{Q}^* + H^2 \sqrt{\frac{\iota}{n_h^k(s_h^k, a_h^k)}}.$$

Combining the upper bounds of $I_1^{h,k}$ Equation (60), $I_{2,1}^{h,k}$ and $I_{2,2}^{h,k}$, we have:

$$\sqrt{\frac{\nu_h^{\text{ref}, k} \iota}{n_h^k}} \lesssim \frac{\sqrt{H \Gamma_h^k(s_h^k, a_h^k) \iota}}{n_h^k(s_h^k, a_h^k)} + \sqrt{\frac{(\mathbb{Q}^* + \beta^2) \iota}{n_h^k(s_h^k, a_h^k)}} + \frac{H \iota^{\frac{3}{4}}}{n_h^k(s_h^k, a_h^k)^{\frac{3}{4}}}. \quad (61)$$

Using the inequality (80) of Zhang et al. (2020), we have:

$$\sqrt{\frac{\check{\nu}_h^k \iota}{\check{n}_h^k}} \leq \sqrt{\frac{\beta^2 \iota}{\check{n}_h^k}} + \sqrt{\frac{H^2 \check{\Gamma}_h^k(s_h^k, a_h^k) \iota}{\check{n}_h^k}} \lesssim \sqrt{\frac{\beta^2 H \iota}{n_h^k}} + \sqrt{\frac{H^2 \check{\Gamma}_h^k(s_h^k, a_h^k) \iota}{\check{n}_h^k}}. \quad (62)$$

where $\check{\Gamma}_h^k(s_h^k, a_h^k) = \sum_{i=1}^{\check{n}_h^k} (V_{h+1}^{\text{ref}, \check{l}_i}(s_{h+1}^{\check{l}_i}) - \hat{V}_{h+1}^{\text{ref}, \check{l}_i}(s_{h+1}^{\check{l}_i}))$. The last inequality is by Lemma B.3.

Applying Equation (61) and Equation (62) to Equation (59), we have:

$$b_h^k(s_h^k, a_h^k) \lesssim \frac{\sqrt{H \Gamma_h^k(s_h^k, a_h^k) \iota}}{n_h^k(s_h^k, a_h^k)} + \sqrt{\frac{(\mathbb{Q}^* + \beta^2 H) \iota}{n_h^k(s_h^k, a_h^k)}} + \frac{H \iota^{\frac{3}{4}}}{\check{n}_h^k(s_h^k, a_h^k)^{\frac{3}{4}}} + \frac{H \sqrt{\check{\Gamma}_h^k(s_h^k, a_h^k) \iota} + H \iota}{\check{n}_h^k(s_h^k, a_h^k)}. \quad (63)$$

B.5 REARRANGE THE WEIGHTED SUM OF G_3

Similar to Equation (42), it holds that:

1611
 1612
 1613
 1614
 1615
$$\sum_{k=1}^K \omega_{h,k} \mathbb{1}[\check{n}_h^k \neq 0] G_3 = \sum_{k=1}^K \omega_{h,k} \mathbb{1}[\check{n}_h^k \neq 0] \frac{\sum_{i=1}^{\check{n}_h^k} (V_{h+1}^{\check{l}_i}(s_{h+1}^{\check{l}_i}) - V_{h+1}^*(s_{h+1}^{\check{l}_i}))}{\check{n}_h^k(s_h^k, a_h^k)}$$

 1616
 1617
 1618
 1619

$$= \sum_{j=1}^K \left(\sum_{k=1}^K \omega_{h,k} \frac{\sum_{i=1}^{\check{n}_h^k} \mathbb{1}[\check{l}_i = j, \check{n}_h^k \neq 0]}{\check{n}_h^k(s_h^k, a_h^k)} \right) (V_{h+1}^j(s_{h+1}^j) - V_{h+1}^*(s_{h+1}^j))$$

$$\begin{aligned} & \leq \sum_{j=1}^K \left(\sum_{k=1}^K \omega_{h,k} \frac{\sum_{i=1}^{\check{n}_h^k} \mathbb{I}[\check{l}_i = j, \check{n}_h^k \neq 0]}{\check{n}_h^k(s_h^k, a_h^k)} \right) (Q_{h+1}^j - Q_{h+1}^*) (s_{h+1}^j, a_{h+1}^j) \\ & \end{aligned} \quad (64)$$

$$= \sum_{j=1}^K \omega_{h+1,j}(h) \left(Q_{h+1}^j(s_{h+1}^j, a_{h+1}^j) - Q_{h+1}^*(s_{h+1}^j, a_{h+1}^j) \right). \quad (65)$$

B.6 BOUNDING THE TERM $\sum_{h'=1}^H \sum_{k=1}^K Y_{h'}^k$

$$\begin{aligned} & \sum_{k=1}^K Y_{h'}^k = \sum_{k=1}^K \mathbb{I}[n_{h'}^k = 0] H \\ & + \sum_{k=1}^K \frac{\mathbb{I}[n_{h'}^k \neq 0]}{n_{h'}^k(s_{h'}^k, a_{h'}^k)} \left(\sum_{i=1}^{\check{n}_{h'}^k} H \left(\mathbb{1}_{s_{h'+1}^{l_i}} + \mathbb{P}_{s_{h'}^k, a_{h'}^k, h'} \right) \lambda_{h'+1}^{l_i} + \sqrt{H\Gamma_{h'}^k(s_{h'}^k, a_{h'}^k)\iota} \right) \\ & + \sum_{k=1}^K \frac{\mathbb{I}[\check{n}_{h'}^k \neq 0]}{\check{n}_{h'}^k(s_{h'}^k, a_{h'}^k)} \left(\sum_{i=1}^{\check{n}_{h'}^k} H \left(\mathbb{P}_{s_{h'}^k, a_{h'}^k, h'} + \mathbb{1}_{s_{h'+1}^{l_i}} \right) \lambda_{h'+1}^{l_i} + H\sqrt{\check{\Gamma}_{h'}^k(s_{h'}^k, a_{h'}^k)\iota} + H\iota \right). \end{aligned} \quad (66)$$

In this equation,

$$\sum_{h'=1}^H \sum_{k=1}^K \mathbb{I}[n_{h'}^k(s_{h'}^k, a_{h'}^k) = 0] H = \sum_{h'=1}^H \sum_{s,a} H \sum_{k=1}^K \mathbb{I}[n_{h'}^k(s, a) = 0, (s_{h'}^k, a_{h'}^k) = (s, a)] \leq H^3 S A. \quad (67)$$

By Lemma B.4, we have the following inequalities:

$$\sum_{h'=1}^H \sum_{k=1}^K \frac{\mathbb{I}[\check{n}_{h'}^k \neq 0]}{\check{n}_{h'}^k(s_{h'}^k, a_{h'}^k)} H\iota \lesssim H^3 S A \iota \log(T). \quad (68)$$

According to Lemma B.5, we have:

$$\sum_{k=1}^K \frac{\mathbb{I}[n_{h'}^k \neq 0]}{n_{h'}^k(s_{h'}^k, a_{h'}^k)} \sum_{i=1}^{\check{n}_{h'}^k} H \left(\mathbb{1}_{s_{h'+1}^{l_i}} + \mathbb{P}_{s_{h'}^k, a_{h'}^k, h'} \right) \lambda_{h'+1}^{l_i} \leq 3H \log T \sum_{k=1}^K \left(\mathbb{1}_{s_{h'+1}^k} + \mathbb{P}_{s_{h'}^k, a_{h'}^k, h'} \right) \lambda_{h'+1}^k. \quad (69)$$

and

$$\sum_{k=1}^K \frac{\mathbb{I}[\check{n}_{h'}^k \neq 0]}{\check{n}_{h'}^k(s_{h'}^k, a_{h'}^k)} \sum_{i=1}^{\check{n}_{h'}^k} H \left(\mathbb{P}_{s_{h'}^k, a_{h'}^k, h'} + \mathbb{1}_{s_{h'+1}^{l_i}} \right) \lambda_{h'+1}^{l_i} \leq 2H \sum_{k=1}^K \left(\mathbb{P}_{s_{h'}^k, a_{h'}^k, h'} + \mathbb{1}_{s_{h'+1}^k} \right) \lambda_{h'+1}^k. \quad (70)$$

Note that

$$\sum_{k=1}^K \lambda_{h'+1}^k = \sum_{k=1}^K \mathbb{I}[N_{h'+1}^k(s_{h'+1}^k) < N_0] = \sum_s \sum_{k=1}^K \mathbb{I}[N_{h'+1}^k(s) < N_0, s_{h'+1}^k = s] \leq S N_0.$$

Then under the event \mathcal{E}_3 in Lemma B.1, it holds that:

$$\sum_{k=1}^K \mathbb{P}_{s_{h'}^k, a_{h'}^k, h'} \lambda_{h'+1}^k \leq 3 \sum_{k=1}^K \lambda_{h'+1}^k + \iota \leq 4S N_0.$$

Applying these two inequalities to Equation (69) and Equation (70), then the following two inequalities holds:

$$\sum_{h'=1}^H \sum_{k=1}^K \frac{\mathbb{I}[n_{h'}^k \neq 0]}{n_{h'}^k(s_{h'}^k, a_{h'}^k)} \sum_{i=1}^{\check{n}_{h'}^k} H \left(\mathbb{1}_{s_{h'+1}^{l_i}} + \mathbb{P}_{s_{h'}^k, a_{h'}^k, h'} \right) \lambda_{h'+1}^{l_i} \lesssim H^2 S N_0 \log(T), \quad (71)$$

1674

and

1675

1676

1677

1678

1679

Meanwhile, according to Lemma B.2 we have:

1680

1681

1682

1683

1684

Then it holds that:

1685

1686

1687

1688

1689

1690

1691

1692

1693

$$\begin{aligned} \Gamma_{h'}^k(s_{h'}^k, a_{h'}^k) &= \sum_{i=1}^{n_{h'}^k} \left(V_{h'+1}^{\text{ref}, l_i}(s_{h'+1}^{k^n}) - \hat{V}_{h'+1}^{\text{ref}, l_i}(s_{h'+1}^{l_i}) \right) \leq H \sum_{i=1}^{n_{h'}^k} \lambda_{h'+1}^{l_i}(s_{h'+1}^{l_i}) \triangleq \Theta_{h'}^k(s_{h'}^k, a_{h'}^k). \\ &\sum_{k=1}^K \frac{\sqrt{\Gamma_{h'}^k(s_{h'}^k, a_{h'}^k)}}{n_{h'}^k(s_{h'}^k, a_{h'}^k)} \leq \sum_{k=1}^K \frac{\sqrt{\Theta_{h'}^k(s_{h'}^k, a_{h'}^k)}}{n_{h'}^k(s_{h'}^k, a_{h'}^k)} \leq \sum_{s,a} \left(\sum_{j \in C} \frac{e_j}{\sum_{i=1}^{j-1} e_i} \right) \sqrt{\Theta_{h'}^K(s, a)} \\ &\lesssim \log T \sum_{s,a} \sqrt{\Theta_{h'}^K(s, a)} \leq \log T \sqrt{SA \sum_{s,a} \Theta_{h'}^K(s, a)} \end{aligned} \quad (73)$$

Here, $C = \{j : H \leq \sum_{i=1}^{j-1} e_i \leq T\}$. The second inequality is by Equation (41) and the monotonicity of $\Theta_{h'}^n(s, a)$. The last inequality is by Cauchy's inequality. To continue, note that:

1694

1695

1696

1697

$$\sum_{h'=1}^H \sqrt{\sum_{s,a} \Theta_{h'}^K(s, a)} \leq \sum_{h'=1}^H \sqrt{H \sum_{k=1}^K \lambda_{h'+1}^k(s_{h'+1}^k)} \leq H \sqrt{HSN_0}$$

1698

Together with Equation (73), it holds:

1699

1700

1701

1702

1703

Since $\check{\Gamma}_{h'}^k(s_{h'}^k, a_{h'}^k) \leq \Gamma_{h'}^k(s_{h'}^k, a_{h'}^k)$ and $4H\check{n}_{h'}^k(s_{h'}^k, a_{h'}^k) \geq n_{h'}^k(s_{h'}^k, a_{h'}^k)$ by Lemma B.3, it holds:

1704

1705

1706

1707

1708

$$\sum_{h'=1}^H \sum_{k=1}^K \frac{H \sqrt{\check{\Gamma}_{h'}^k(s_{h'}^k, a_{h'}^k) \iota}}{\check{n}_{h'}^k(s_{h'}^k, a_{h'}^k)} \lesssim \log T \sum_{h'=1}^H \sqrt{HSA \sum_{s,a} \Theta_{h'}^K(s, a) \iota} \lesssim H^2 S \log(T) \sqrt{AN_0 \iota}. \quad (74)$$

1709

1710

1711

Applying the inequalities Equation (67), Equation (68), Equation (71), Equation (72), Equation (74) and Equation (75) to Equation (66), since $N_0 = O(\frac{SAH^5\iota}{\beta^2})$, we have:

1712

1713

1714

1715

$$\sum_{h'=1}^H \sum_{k=1}^K Y_{h'}^k \leq O\left(\frac{H^7 S^2 A \iota \log(T)}{\beta^2}\right).$$

1716

C PROOF OF THEOREM 3.3

1717

In this section, we will prove Theorem 3.3.

1718

1719

1720

1721

1722

1723

Proof. For $\delta \in (0, 1)$, let $p \leftarrow \frac{\delta}{40SAH^2T^4}$, then $\iota = \log(\frac{2}{p}) = O(\frac{SAT}{\delta})$. Now with probability at least $1 - \delta$, $\bigcap_{i=1}^7 \mathcal{E}_i$ holds. Next, we will prove the upper bound for policy switching cost under the event $\bigcap_{i=1}^7 \mathcal{E}_i$.

1724

From the proof of Theorem 2 in Zhang et al. (2020), we have:

1725

1726

1727

$$N_{\text{switch}} \leq \sum_{s,a,h} 4H \log \left(\frac{N_h^{K+1}(s, a)}{2H} + 1 \right).$$

Next for any $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$, we will bound the term $N_h^{K+1}(s, a)$. For $a \neq \pi_h^*(s)$. In this case, we have $\Delta_h(s, a) > 0$ and then $\Delta_h(s, a) \geq \Delta_{\min}$. For any $h \in [H]$, let set D_h be all triples of (s, a, h) such that $a \neq \pi_h^*(s)$, that is:

$$D_h = \{(s, a, h) | a \neq \pi_h^*(s)\}.$$

We also let the set $D = \bigcup_{h=1}^H D_h$ and the set $D_{\text{opt}} = \{(s, a, h) | a = \pi_h^*(s)\}$. Then we have $|D| + |D_{\text{opt}}| = SAH$. Since for every state-step pair (s, h) , there exists at least one optimal action. Therefore we know $|D_{\text{opt}}| \geq SH$ and then $0 \leq |D| \leq SA(H - 1)$.

If for given $(h, k) \in [H] \times [k]$, $(s_h^k, a_h^k) \in D_h$, we have $\Delta_h(s_h^k, a_h^k) \geq \Delta_{\min}$. Then it holds that:

$$Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) = V_h^k(s_h^k) - Q_h^*(s_h^k, a_h^k) \geq \Delta_h(s_h^k, a_h^k) \geq \Delta_{\min}.$$

The first inequality is because $V_h^k(s) \geq V_h^*(s)$. Therefore, we have

$$\begin{aligned} \sum_{(s, a, h) \in D_h} \mathbb{I}[(s_h^k, a_h^k) = (s, a)] &= \mathbb{I}[(s_h^k, a_h^k, h) \in D_h] \\ &\leq \mathbb{I}[Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \geq \Delta_{\min}] \\ &= \sum_{n=1}^N \omega_{h,k}^{(n)}. \end{aligned}$$

and then

$$\begin{aligned} \sum_{(s, a, h) \in D} N_h^{K+1}(s, a) &= \sum_{h=1}^H \sum_{(s, a, h) \in D_h} N_h^{K+1}(s, a) = \sum_{h=1}^H \sum_{(s, a, h) \in D_h} \sum_{k=1}^K \mathbb{I}[(s_h^k, a_h^k) = (s, a)] \\ &\leq \sum_{h=1}^H \sum_{k=1}^K \sum_{n=1}^N \omega_{h,k}^{(n)} = \sum_{h=1}^H \sum_{n=1}^N \|\omega\|_{1,h}^{(n)}. \end{aligned}$$

By Equation (52), we know:

$$\sum_{(s, a, h) \in D_{\text{opt}}^c} N_h^{K+1}(s, a) \leq \sum_{h=1}^H \sum_{n=1}^N \|\omega\|_{1,h}^{(n)} \leq O\left(\frac{(\mathbb{Q}^* + \beta^2 H) SAH^3 \ell}{\Delta_{\min}} + \frac{H^8 S^2 A \ell \log(T)}{\beta^2}\right) \quad (76)$$

Therefore we have:

$$\begin{aligned} N_{\text{switch}} &\leq \sum_{s, a, h} 4H \log\left(\frac{N_h^{K+1}(s, a)}{2H} + 1\right) \\ &= \sum_{(s, a, h) \in D_{\text{opt}}^c} 4H \log\left(\frac{N_h^{K+1}(s, a)}{2H} + 1\right) + \sum_{(s, a, h) \notin D_{\text{opt}}} 4H \log\left(\frac{N_h^{K+1}(s, a)}{2H} + 1\right) \quad (77) \\ &\leq 4H(SAH - |D_{\text{opt}}|) \log\left(1 + \frac{\sum_{h=1}^H \sum_{n=1}^N \|\omega\|_{1,h}^{(n)}}{2H(SAH - |D_{\text{opt}}|)}\right) + 4H|D_{\text{opt}}| \log\left(\frac{T}{2H|D_{\text{opt}}|} + 1\right) \\ &\leq O\left(H(SAH - |D_{\text{opt}}|) \log\left(\frac{(\mathbb{Q}^* + \beta^2 H) H^2 S A \ell}{(SAH - |D_{\text{opt}}|) \Delta_{\min}^2} + \frac{H^7 S^2 A \ell \log(T)}{\beta^2 (SAH - |D_{\text{opt}}|) \Delta_{\min}}\right)\right. \\ &\quad \left.+ H|D_{\text{opt}}| \log\left(\frac{K}{|D_{\text{opt}}|} + 1\right)\right). \quad (78) \end{aligned}$$

The first inequality is because of Jensen's Inequality. The last inequality is by Equation (52). Since $\mathbb{Q}^* \leq H^2$ and $\beta \leq H$, then we have:

$$\frac{(\mathbb{Q}^* + \beta^2 H) H^2 S A \ell}{(SAH - |D_{\text{opt}}|) \Delta_{\min}^2} \leq \frac{H^7 S A \ell}{\beta^2 (SAH - |D_{\text{opt}}|) \Delta_{\min}^2}.$$

1782 By $\Delta_{\min} \leq H$, we also have:
 1783

$$\frac{H^7 S^2 A \iota \log(T)}{\beta^2(SAH - |D_{\text{opt}}|) \Delta_{\min}} \leq \frac{H^8 S^2 A \iota \log(T)}{\beta^2(SAH - |D_{\text{opt}}|) \Delta_{\min}^2}.$$

1784
 1785
 1786 For $\delta \in (0, 1)$, let $p \leftarrow \frac{\delta}{60SAH^2T^5}$, then $\iota = \log(\frac{2}{p}) \leq O(\log(\frac{SAT}{\delta}))$. Applying the above two
 1787 inequalities to Equation (78), with probability at least $1 - \delta$, we have it holds that:
 1788

$$\begin{aligned} N_{\text{switch}} &\leq O\left(H(SAH - |D_{\text{opt}}|) \log\left(\frac{H^8 S^2 A \iota \log(T)}{\beta^2(SAH - |D_{\text{opt}}|) \Delta_{\min}^2}\right) + H|D_{\text{opt}}| \log\left(\frac{K}{|D_{\text{opt}}|} + 1\right)\right) \\ &= O\left(H(SAH - |D_{\text{opt}}|) \log\left(\frac{H^4 S A^{\frac{1}{2}} \iota}{\beta \sqrt{(SAH - |D_{\text{opt}}|) \Delta_{\min}}}\right) + H|D_{\text{opt}}| \log\left(\frac{K}{|D_{\text{opt}}|} + 1\right)\right) \\ &= O\left(H|D_{\text{opt}}^c| \log\left(\frac{H^4 S A^{\frac{1}{2}} \log(\frac{SAT}{\delta})}{\beta \sqrt{|D_{\text{opt}}^c| \Delta_{\min}}}\right) + H|D_{\text{opt}}| \log\left(\frac{K}{|D_{\text{opt}}|} + 1\right)\right). \end{aligned}$$

1805
 1806
 1807 Especially, if the optimal policy is deterministic and unique, which means $|D_{\text{opt}}| = SA$, then the
 1808 policy switching cost is upper bounded by:
 1809

$$O\left(H^2 S A \log\left(\frac{H^{\frac{7}{2}} S^{\frac{1}{2}} \log(\frac{SAT}{\delta})}{\beta \Delta_{\min}}\right) + H^2 S \log\left(\frac{K}{HS} + 1\right)\right).$$

1810
 1811 \square
 1812

D PROOF OF THEOREM 3.2

D.1 ALGORITHM DETAILS

1825 Before continuing, let us briefly introduce the refined algorithm, which is similar to the original
 1826 version in Li et al. (2021). Before diving into the algorithm itself, we will first discuss the key
 1827 auxiliary functions used for estimating the Q -value functions. For any $\delta \in [0, 1]$, let $\iota = \log(\frac{SAT}{\delta})$.

1828 In the algorithm, μ_h^{ref} and σ_h^{ref} are updated to represent the current mean and second moment of
 1829 the reference function. μ_h^{adv} and σ_h^{adv} are updated to be the current weighted mean and weighted
 1830 second moment of the reference function with weight to be the learning rate $\eta_n = \frac{H+1}{H+n}$. b_h^R is the
 1831 exploration bonus for Q-EarlySettled-Advantage. With these update functions, we can then discuss
 1832 the Q-EarlySettled-Advantage algorithm.
 1833

1836 **Algorithm 2** Auxiliary functions

1837 1: **function** UPDATE-UCB-Q

1838 2: $Q_h^{\text{UCB}}(s_h, a_h) \leftarrow (1 - \eta_n)Q_h^{\text{UCB}}(s_h, a_h) + \eta_n \left(r_h(s_h, a_h) + V_{h+1}(s_{h+1}) + c_b \sqrt{\frac{H^3 \iota}{n}} \right)$.

1839 3: **function** UPDATE-LCB-Q

1840 4: $Q_h^{\text{LCB}}(s_h, a_h) \leftarrow (1 - \eta_n)Q_h^{\text{LCB}}(s_h, a_h) + \eta_n \left(r_h(s_h, a_h) + V_{h+1}^{\text{LCB}}(s_{h+1}) - c_b \sqrt{\frac{H^3 \iota}{n}} \right)$.

1841 5: **function** UPDATE-UCB-UCB-ADVANTAGE

1842 6: $[\mu_h^{\text{ref}}, \sigma_h^{\text{ref}}, \mu_h^{\text{adv}}, \sigma_h^{\text{adv}}](s_h, a_h) \leftarrow \text{UPDATE-MOMENTS}()$;

1843 7: $[\delta_h^{\text{R}}, B_h^{\text{R}}](s_h, a_h) \leftarrow \text{UPDATE-BONUS}()$;

1844 8: $b_h^{\text{R}} \leftarrow B_h^{\text{R}}(s_h, a_h) + (1 - \eta_n) \frac{\delta_h^{\text{R}}(s_h, a_h)}{\eta_n} + c_b \frac{H^2 \iota}{n^{3/4}}$;

1845 9: $Q_h^{\text{R}}(s_h, a_h) \leftarrow (1 - \eta_n)Q_h^{\text{R}}(s_h, a_h)$
 $+ \eta_n (r_h(s_h, a_h) + V_{h+1}(s_{h+1}) - V_{h+1}^{\text{R}}(s_{h+1}) + \mu_h^{\text{ref}}(s_h, a_h) + b_h^{\text{R}})$.

1846 10: **function** UPDATE-MOMENTS

1847 11: $\mu_h^{\text{ref}}(s_h, a_h) \leftarrow (1 - \frac{1}{n}) \mu_h^{\text{ref}}(s_h, a_h) + \frac{1}{n} V_{h+1}^{\text{R}}(s_{h+1})$;

1848 12: $\sigma_h^{\text{ref}}(s_h, a_h) \leftarrow (1 - \frac{1}{n}) \sigma_h^{\text{ref}}(s_h, a_h) + \frac{1}{n} (V_{h+1}^{\text{R}}(s_{h+1}))^2$;

1849 13: $\mu_h^{\text{adv}}(s_h, a_h) \leftarrow (1 - \eta_n) \mu_h^{\text{adv}}(s_h, a_h) + \eta_n (V_{h+1}(s_{h+1}) - V_{h+1}^{\text{R}}(s_{h+1}))$;

1850 14: $\sigma_h^{\text{adv}}(s_h, a_h) \leftarrow (1 - \eta_n) \sigma_h^{\text{adv}}(s_h, a_h) + \eta_n (V_{h+1}(s_{h+1}) - V_{h+1}^{\text{R}}(s_{h+1}))^2$;

1851 15: **function** UPDATE-BONUS

1852 16: $B_h^{\text{next}}(s_h, a_h) \leftarrow$
 $c_b \sqrt{\frac{\iota}{n}} \left(\sqrt{\sigma_h^{\text{ref}}(s_h, a_h) - (\mu_h^{\text{ref}}(s_h, a_h))^2} + \sqrt{H} \sqrt{\sigma_h^{\text{adv}}(s_h, a_h) - (\mu_h^{\text{adv}}(s_h, a_h))^2} \right)$;

1853 17: $\delta_h^{\text{R}}(s_h, a_h) = B_h^{\text{next}}(s_h, a_h) - B_h^{\text{R}}(s_h, a_h)$;

1854 18: $B_h^{\text{R}}(s_h, a_h) \leftarrow B_h^{\text{next}}(s_h, a_h) - \delta_h^{\text{R}}(s_h, a_h)$.

1863

1864

1865

1866

1867 **Algorithm 3** Refined Q-EarlySettled-Advantage

1868 1: **Parameters:** Some universal constant $c_b > 0$ and probability of failure $\delta \in (0, 1)$;

1869 2: **Initialize** $Q_h^1(s, a), Q_h^{\text{UCB}, 1}(s, a), Q_h^{\text{R}, 1}(s, a) \leftarrow H; Q_h^{\text{LCB}, 1}(s, a) \leftarrow 0$;
 $V_h^1(s) \leftarrow H, N_h^1(s, a), \mu_h^{\text{ref}}(s, a), \sigma_h^{\text{ref}}(s, a), \mu_h^{\text{adv}}(s, a), \sigma_h^{\text{adv}}(s, a), \delta_h^{\text{R}}(s, a), B_h^{\text{R}}(s, a) \leftarrow 0$;
 and $u_h^k(s, a, h) \leftarrow \text{True}$, for all $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$.

1870 3: **for** Episode $k = 1$ to K **do**

1871 4: Set initial state $s_1^k \leftarrow s_1^k$;

1872 5: **for** Step $h = 1$ to H **do**

1873 6: Take action $a_h^k = \pi_h^k(s_h^k) = \arg \max_a Q_h^k(s_h^k, a)$, and draw $s_{h+1}^k \sim P_h(\cdot | s_h^k, a_h^k)$;

1874 7: $N_h^k(s_h^k, a_h^k) \leftarrow N_h^{k-1}(s_h^k, a_h^k) + 1; n \leftarrow N_h^k(s_h^k, a_h^k)$;

1875 8: $\eta_n \leftarrow \frac{H+1}{H+n}$;

1876 9: $Q_h^{\text{UCB}, k+1}(s_h^k, a_h^k) \leftarrow \text{UPDATE-UCB-Q}()$.

1877 10: $Q_h^{\text{LCB}, k+1}(s_h^k, a_h^k) \leftarrow \text{UPDATE-LCB-Q}()$.

1878 11: $Q_h^{\text{R}, k+1}(s_h^k, a_h^k) \leftarrow \text{UPDATE-UCB-UCB-Advantage}()$.

1879 12: $Q_h^{k+1}(s_h^k, a_h^k) \leftarrow \min\{Q_h^{\text{R}, k+1}(s_h^k, a_h^k), Q_h^{\text{UCB}, k+1}(s_h^k, a_h^k), Q_h^k(s_h^k, a_h^k)\}$;

1880 13: $V_h^{k+1}(s_h^k) \leftarrow \max_a Q_h^{k+1}(s_h^k, a)$;

1881 14: $V_h^{\text{LCB}, k+1}(s_h^k) \leftarrow \max \left\{ \max_a Q_h^{\text{LCB}, k+1}(s_h^k, a), V_h^{\text{LCB}, k}(s_h^k) \right\}$;

1882 15: **if** $V_h^{k+1}(s_h^k) - V_h^{\text{LCB}, k+1}(s_h^k) > \beta$ **then**

1883 16: $V_h^{\text{R}, k+1}(s_h^k) \leftarrow V_h^{k+1}(s_h^k)$;

1884 17: **else if** $u_h^k(s_h^k) = \text{True}$ **then**

1885 18: $V_h^{\text{R}, k+1}(s_h^k) \leftarrow V_h^{k+1}(s_h^k); u_h^{k+1}(s_h^k) = \text{False}$.

At the beginning of the k -th episode, we can obtain V -estimate $V_h^k(s)$, the reference function $V_h^{R,k}(s)$ and the policy π^k from the previous episode $k - 1$ and select a initial state s_1^k (For the first episode, we randomly choose a policy π^1 and $V_h^1(s) = V_h^{R,1} = H$). At step $h \in [H]$, we can process the trajectory with $a_h^k = \pi_h^k(s_h^k)$ and $s_{h+1}^k \sim \mathbb{P}_h(\cdot | s_h^k, a_h^k)$. Now we need to update the estimates of both Q -value and V -value functions at the end of k -th episode. In the algorithm, the estimate learned from the UCB by the end of k -th episode, denoted as $Q_h^{\text{UCB},k+1}$, is updated to:

$$Q_h^{\text{UCB},k+1} = r_h^k(s_h^k, a_h^k) + \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(V_{h+1}^{k^n}(s_{h+1}^{k^n}) + c_b \sqrt{\frac{H^3 \iota}{n}} \right) \quad (79)$$

Here we define $N_h^k = N_h^k(s_h^k, a_h^k)$ as the number of times that the state-action pair (s_h^k, a_h^k) has been visited at step h at the beginning of the k -th episode and $k^n = k_h^n(s_h^k, a_h^k)$ denotes the index of the episode in which the state-action pair (s_h^k, a_h^k) is visited for the n -th time at step h . The term $c_b \sqrt{\frac{H^3 \iota}{n}}$ represents the exploration bonus for n -th visit, where $c_b > 0$ is a sufficiently large constant and $\iota = \log(\frac{SAT}{\delta})$ with $\delta \in (0, 1)$ being failure probability.

Another Q -estimate obtained from LCB at the end of k -th episode, denoted as $Q_h^{\text{LCB},k+1}$, is updated similarly to $Q_h^{\text{UCB},k+1}$, but with the exploration bonus subtracted instead.

The last estimate of Q -value function, denoted as $Q_h^{R,k+1}$, uses reference-advantage decomposition techniques. At the end of k -th episode, $Q_h^{R,k+1}$ is updated to:

$$Q_h^{R,k+1} = r_h^k(s_h^k, a_h^k) + \sum_{n=1}^{N_h^{k+1}} \eta_n^{N_h^{k+1}} \left(V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) + \frac{\sum_{i=1}^n V_{h+1}^{R,k^i}(s_{h+1}^{k^i})}{n} + b_h^{R,k^n+1} \right). \quad (80)$$

In Equation (80), $V_h^{R,k}(s)$ is the reference function learned at the end of episode $k - 1$. The key idea of the reference-advantage decomposition is that we expect to maintain a collection of reference values $\{V_h^{R,k}(s)\}_{s,k,h}$, which form reasonable estimates of $\{V_h^*(s)\}_{s,h}$ and become increasingly more accurate as the algorithm progresses. It means for any $s \in \mathcal{S}$, sufficiently large k and some given $\beta \in (0, H]$, it holds $|V_h^{R,k}(s) - V_h^*(s)| \leq \beta$. In this case, for $s_{h+1}^{k^n} \sim \mathbb{P}_h(\cdot | s_h^{k^n}, a_h^{k^n})$, the variance of the advantage term $V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{R,k^n}(s_{h+1}^{k^n})$, is bounded by β^2 , which can be less volatile than the stochastic term $V_{h+1}^{k^n}(s_{h+1}^{k^n})$. Meanwhile, the reference term $\sum_{i=1}^n V_{h+1}^{R,k^i}(s_{h+1}^{k^i})/n$ use a batch of historical visits to (s_h^k, a_h^k, h) , which can lower the variance as the increase of the sample size n . Accordingly, the exploration bonus term b_h^{R,k^n+1} is taken to be an upper confidence bound for the above-mentioned two terms combined. Given that the uncertainty of Equation (80) largely stems from the advantage and the reference terms (which can both be much smaller than the variability in Equation (79)), the incorporation of the reference function helps accelerate convergence and lower the regret upper bound.

With two additional Q-estimates in hand — $Q_h^{\text{UCB},k+1}$ learned from UCB and $Q_h^{R,k+1}$ obtained from the reference-advantage decomposition, it is natural to combine them as follows to further reduce the bias without violating the optimism principle:

$$Q_h^{k+1}(s_h^k, a_h^k) = \min\{Q_h^{\text{UCB},k+1}(s_h^k, a_h^k), Q_h^{R,k+1}(s_h^k, a_h^k), Q_h^k(s_h^k, a_h^k)\}. \quad (81)$$

We also incorporate $Q_h^k(s_h^k, a_h^k)$ here to keep the monotonicity of the update. Then we can learn $V_h^{k+1}(s_h^k, a_h^k)$ and $V_h^{\text{LCB},k+1}(s_h^k, a_h^k)$ by a greedy policy with respect to these Q -estimates:

$$V_h^{k+1}(s_h^k) = \max_a Q_h^{k+1}(s_h^k, a), V_h^{\text{LCB},k+1}(s_h^k) = \max \left\{ \max_a Q_h^{\text{LCB},k+1}(s_h^k, a), V_h^{\text{LCB},k}(s_h^k) \right\}.$$

In the algorithm, $V_h^{\text{LCB},k}(s)$ is used as lower bound estimates of $V_h^*(s)$. We learn the final value $V_h^R(s)$ of the reference function for the state-step pair (s, h) when it first meets the condition $V_h^k(s) - V_h^{\text{LCB},k}(s) \leq \beta$.

1944 D.2 AUXILIARY LEMMAS
 1945

1946 As can be easily verified, we have

$$1948 \quad \sum_{n=1}^N \eta_n^N = \begin{cases} 1, & \text{if } N > 0, \\ 0, & \text{if } N = 0. \end{cases} \quad (82)$$

1950 **Lemma D.1.** For any integer $N > 0$, the following properties hold:

$$1952 \quad \frac{1}{N^a} \leq \sum_{n=1}^N \frac{\eta_n^N}{n^a} \leq \frac{2}{N^a}, \quad \text{for all } \frac{1}{2} \leq a \leq 1, \quad (83)$$

$$1956 \quad \max_{1 \leq n \leq N} \eta_n^N \leq \frac{2H}{N}, \quad \sum_{n=1}^N (\eta_n^N)^2 \leq \frac{2H}{N}, \quad \sum_{n=n}^{\infty} \eta_n^N \leq 1 + \frac{1}{H}. \quad (84)$$

1959 *Proof.* It is proved in Appendix B of Li et al. (2021). \square
 1960

1961 Let $u_i^N = \sum_{n=i}^N \frac{\eta_n^N}{n}$. Then according to Equation (83), we know $u_i^N \leq \frac{2}{N}$ for any $i \leq N \in \mathbb{N}_+$.

1963 **Lemma D.2.** Consider any $\delta \in (0, 1)$ and $\iota = \log(\frac{SAT}{\delta})$. Using $\forall(s, a, h, k)$ as the simplified notation for $\forall(s, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$ and $\forall(s, h, k)$ as the simplified notation for $\forall(s, a, h, k) \in \mathcal{S} \times [H] \times [K]$. Let $\hat{V}_h^{R,k}(s) = \max\{V_h^*(s), \min\{V_h^*(s) + \beta, V_h^{R,k}(s)\}\}$. Then we have the following conclusions:

1967 (a) (Lemma 2 of Li et al. (2021)) With probability at least $1 - \delta$, the following event holds:
 1968

$$1969 \quad \mathcal{E}_1 = \left\{ Q_h^*(s, a) \leq Q_h^{k+1}(s, a) \leq Q_h^k(s, a), V_h^*(s) \leq V_h^k(s) \leq V_h^{R,k}(s), \forall(s, a, h, k) \right\}.$$

1972 (b) (Lemma 3 of Li et al. (2021)) With probability at least $1 - \delta$, the following event holds:
 1973

$$1974 \quad \mathcal{E}_2 = \left\{ Q_h^{\text{LCB},k}(s, a) \leq Q_h^*(s, a), V_h^{\text{LCB},k}(s) \leq V_h^*(s), \forall(s, a, h, k) \text{ and} \right. \\ 1975 \quad \left. \sum_{h=1}^H \sum_{k=1}^K \mathbb{1} \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB},k}(s_h^k, a_h^k) > \varepsilon \right) \lesssim \frac{SAH^6\iota}{\varepsilon^2}, \text{ for any } \varepsilon \in (0, H] \right\}.$$

1980 (c) (Paraphrased from Lemma 4 of Li et al. (2021)) With probability at least $1 - \delta$, the following event holds:
 1981

$$1982 \quad \mathcal{E}_3 = \left\{ \left| V_h^k(s) - V_h^{R,k}(s) \right| \leq 2\beta \text{ and} \right. \\ 1983 \quad \left. \sum_{h=1}^H \sum_{k=1}^K \left(V_h^{R,k}(s_h^k) - V_h^{R,k+1}(s_h^k) \right) \leq \frac{2H^6SAT}{\beta}, \forall(s, h, k) \right\}.$$

1988 (d) With probability at least $1 - \delta$, the following event holds:
 1989

$$1990 \quad \mathcal{E}_4 = \left\{ \sum_{i=1}^{N_h^k} u_i^{N_h^k} \left(\mathbb{1}_{s_{h+1}^{k^i}} - \mathbb{P}_{s,a,h} \right) (\hat{V}_{h+1}^{R,k^i} - V_{h+1}^*) \leq 2 \sqrt{\frac{2\beta^2\iota}{N_h^k(s, a)}}, \forall(s, a, h, k) \right\}.$$

1994 (e) With probability at least $1 - \delta$, the following event holds:
 1995

$$1996 \quad \mathcal{E}_5 = \left\{ \sum_{i=1}^{N_h^k} u_i^{N_h^k} \left(\mathbb{1}_{s_{h+1}^{k^i}} - \mathbb{P}_{s,a,h} \right) V_{h+1}^* \leq 8 \sqrt{\frac{\mathbb{Q}^*\iota}{N_h^k(s, a)}} + 16 \frac{H\iota}{N_h^k(s, a)}, \forall(s, a, h, k) \right\}.$$

1998 (f) With probability at least $1 - \delta$, the following event holds:
 1999

$$2000 \quad \mathcal{E}_6 = \left\{ \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(\mathbb{P}_{s,a,h} - \mathbb{1}_{s_{h+1}^{k^n}} \right) \left(\hat{V}_{h+1}^{R,k^i} - V_{h+1}^* \right) \leq 2 \sqrt{\frac{\beta^2 H \iota}{N_h^k(s,a)}}, \forall (s,a,h,k) \right\}. \\ 2001 \\ 2002$$

2003 (g) With probability at least $1 - \delta$, the following event holds:
 2004

$$2005 \quad \mathcal{E}_7 = \left\{ \sum_{h=1}^H \sum_{k=1}^K \mathbb{P}_{s_h^k, a_h^k, h} \left| V_h^{R,K+1}(s_h^k) - \hat{V}_h^{R,k}(s_h^k) \right| \right. \\ 2006 \quad \left. \leq 3 \sum_{h=1}^H \sum_{k=1}^K \left| V_h^{R,K+1}(s_h^k) - \hat{V}_h^{R,k}(s_h^k) \right| + H \iota \forall (s,h,k) \right\}. \\ 2007 \\ 2008 \\ 2009 \\ 2010$$

2011 *Proof.* (d) From the definition of $\hat{V}_h^{R,k}(s)$, we know that for any $k \in [K]$:
 2012

$$2013 \quad V_h^*(s) \leq \hat{V}_h^{R,k}(s) \leq V_h^*(s) + \beta. \quad (85)$$

2014 Then the sequence

$$2015 \quad \left\{ \sum_{i=1}^j u_i^N \left(\mathbb{1}_{s_{h+1}^{k^i}} - \mathbb{P}_{s,a,h} \right) \left(\hat{V}_{h+1}^{R,k^i} - V_{h+1}^* \right) \right\}_{j \in \mathbb{N}^+}$$

2016 is a martingale sequence with

$$2017 \quad \left| u_i^N \left(\mathbb{1}_{s_{h+1}^{k^i}} - \mathbb{P}_{s,a,h} \right) \left(\hat{V}_{h+1}^{R,k^i} - V_{h+1}^* \right) \right| \leq \frac{2\beta}{N}.$$

2018 Then according to Azuma-Hoeffding inequality, for any $\delta \in (0, 1)$, with probability at least
 2019 $1 - \frac{\delta}{SAT}$, it holds for given $N_H^k(s, a) = N \in \mathbb{N}_+$ that:

$$2020 \quad \sum_{i=1}^N \left(\mathbb{1}_{s_{h+1}^{k^i}} - \mathbb{P}_{s,a,h} \right) \left(\hat{V}_{h+1}^{R,k^i} - V_{h+1}^* \right) \leq 2 \sqrt{\frac{2\beta^2 \iota}{N}}.$$

2021 For any all $(s, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$, we have $N_h^k(s, a) \in [\frac{T}{H}]$. Considering all the
 2022 possible combinations $(s, a, h, N) \in \mathcal{S} \times \mathcal{A} \times [H] \times [\frac{T}{H}]$, with probability at least $1 - \delta$,
 2023 it holds simultaneously for all $(s, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$ that:

$$2024 \quad \sum_{i=1}^{N_h^k} u_i^{N_h^k} \left(\mathbb{1}_{s_{h+1}^{k^i}} - \mathbb{P}_{s,a,h} \right) \left(\hat{V}_{h+1}^{R,k^i} - V_{h+1}^* \right) \leq 2 \sqrt{\frac{2\beta^2 \iota}{N_h^k(s,a)}}.$$

2025 (e) The sequence

$$2026 \quad \left\{ \sum_{i=1}^j u_i^N \left(\mathbb{1}_{s_{h+1}^{k^i}} - \mathbb{P}_{s,a,h} \right) V_{h+1}^* \right\}_{j \in \mathbb{N}_+}$$

2027 is a martingale sequence with

$$2028 \quad \left| u_i^N \left(\mathbb{1}_{s_{h+1}^{k^i}} - \mathbb{P}_{s,a,h} \right) V_{h+1}^* \right| \leq \frac{2H}{N}.$$

2029 Using Lemma A.3 with $c = \frac{2H}{N}$, $\epsilon = c^2$ and δ being $\frac{\delta}{SAT^2}$, for any given $N_h^k(s, a) = N \in \mathbb{N}_+$, with probability at least $1 - (\log_2(n) + 1) \frac{\delta}{SAT^2} \geq 1 - \frac{\delta}{SAT}$, we have:

$$2030 \quad \sum_{i=1}^N u_i^N \left(\mathbb{1}_{s_{h+1}^{k^i}} - \mathbb{P}_{s,a,h} \right) V_{h+1}^* \leq 8 \sqrt{\frac{\mathbb{Q}^* \iota}{N}} + 16 \frac{H \iota}{N}.$$

2031 For any all $(s, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$, we have $N_h^k(s, a) \in [\frac{T}{H}]$. Considering all the
 2032 possible combinations $(s, a, h, N) \in \mathcal{S} \times \mathcal{A} \times [H] \times [\frac{T}{H}]$, with probability at least $1 - \delta$,
 2033 it holds simultaneously for all $(s, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$:

$$2034 \quad \sum_{i=1}^{N_h^k} u_i^{N_h^k} \left(\mathbb{1}_{s_{h+1}^{k^i}} - \mathbb{P}_{s,a,h} \right) V_{h+1}^* \leq 8 \sqrt{\frac{\mathbb{Q}^* \iota}{N_h^k(s,a)}} + 16 \frac{H \iota}{N_h^k(s,a)}.$$

2052 (f) The sequence
 2053
 2054
 2055

$$\left\{ \sum_{n=1}^j \eta_n^N \left(\mathbb{P}_{s,a,h} - \mathbb{1}_{s_{h+1}^{k^n}} \right) \left(\hat{V}_{h+1}^{\mathbf{R}, k^n} - V_{h+1}^* \right) \right\}_{j \in \mathbb{N}^+}$$

2056 is a martingale sequence with
 2057
 2058

$$\eta_n^N \left(\mathbb{P}_{s,a,h} - \mathbb{1}_{s_{h+1}^{k^n}} \right) \left(\hat{V}_{h+1}^{\mathbf{R}, k^n} - V_{h+1}^* \right) \leq \eta_n^N \beta.$$

2060 Then according to Azuma-Hoeffding inequality and Equation (84), for any $\delta \in (0, 1)$, with
 2061 probability at least $1 - \frac{\delta}{SAT}$, it holds for given $N_h^k(s, a) = N \in \mathbb{N}_+$ that:
 2062

$$\sum_{n=1}^N \eta_n^N \left(\mathbb{P}_{s,a,h} - \mathbb{1}_{s_{h+1}^{k^n}} \right) \left(\hat{V}_{h+1}^{\mathbf{R}, k^n} - V_{h+1}^* \right) \leq 2 \sqrt{\frac{\beta^2 H \iota}{N}}.$$

2066 For any all $(s, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$, we have $N_h^k(s, a) \in [\frac{T}{H}]$. Considering all the
 2067 possible combinations $(s, a, h, N) \in \mathcal{S} \times \mathcal{A} \times [H] \times [\frac{T}{H}]$, with probability at least $1 - \delta$,
 2068 it holds simultaneously for all $(s, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$ that:
 2069

$$\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(\mathbb{P}_{s,a,h} - \mathbb{1}_{s_{h+1}^{k^n}} \right) \left(\hat{V}_{h+1}^{\mathbf{R}, k^n} - V_{h+1}^* \right) \leq 2 \sqrt{\frac{\beta^2 H \iota}{N_h^k(s, a)}}$$

2073 (g) This conclusion is directly proved by Lemma A.2 with $l = H$.
 2074
 2075 \square

2076 **Lemma D.3.** For any non-negative weight sequence $\{\omega_{h,k}\}_{h,k}$ and $\alpha \in (0, 1)$, it holds that:
 2077

$$\sum_{k=1}^K \frac{\omega_{h,k}}{N_h^k(s_h^k, a_h^k)^\alpha} \leq \frac{1}{1-\alpha} (SA \|\omega\|_{\infty, h})^\alpha \|\omega\|_{1,h}^{1-\alpha},$$

2081 Here, $\|\omega\|_{\infty, h} = \max_k \{\omega_{h,k}\}$ and $\|\omega\|_{1,h} = \sum_{k=1}^K \omega_{h,k}$.
 2082 For $\alpha = 1$, we have the following conclusions:
 2083

$$\sum_{k=1}^K \frac{1}{N_h^k(s_h^k, a_h^k)} \lesssim SA \log(T),$$

2088 *Proof.*

$$\sum_{k=1}^K \frac{\omega_{h,k}}{N_h^k(s_h^k, a_h^k)^\alpha} = \sum_{s,a} \sum_{i=1}^{N_h^K(s,a)} \frac{\omega_{h,k^i(s,a)}}{i^\alpha} \quad (86)$$

2093 Here $k^i(s, a)$ is the episode index of the i -th visits to (s, a, h) . Let $c_h(s, a) = \sum_{i=1}^{N_h^K(s,a)} \omega_{h,k^i(s,a)}$
 2094 and then we have $\sum_{s,a} c_h(s, a) = \sum_{k=1}^K \omega_{h,k} = \|\omega\|_{1,h}$. Given the term $\sum_{k=1}^K \frac{\omega_{h,k^i(s,a)}}{i^\alpha}$, when
 2095 the weights $\omega_{h,k^i(s,a)}$ concentrates on former terms, we can obtain the largest value. Let
 2096

$$k_{s,a,h} = \left\lceil \frac{c_h(s, a)}{\|\omega\|_{\infty, h}} \right\rceil \text{ and } d_{s,a,h} = c_h(s, a) - (k_{s,a,h} - 1) \|\omega\|_{\infty, h}.$$

2099 Then we have:
 2100

$$\begin{aligned} \sum_{k=1}^K \frac{\omega_{h,k}}{N_h^k(s_h^k, a_h^k)^\alpha} &\leq \sum_{s,a} \sum_{i=1}^{k_{s,a,h}-1} \frac{\|\omega\|_{\infty, h}}{i^\alpha} + \frac{d_{s,a,h}}{k_{s,a,h}^\alpha} \\ &\leq \sum_{s,a} \|\omega\|_{\infty, h} \sum_{i=1}^{k_{s,a,h}-1} \frac{i^{1-\alpha} - (i-1)^{1-\alpha}}{1-\alpha} + \frac{d_{s,a,h}}{k_{s,a,h}^\alpha} \end{aligned}$$

$$\begin{aligned}
&= \sum_{s,a} \frac{\|\omega\|_{\infty,h} (k_{s,a,h} - 1)^{1-\alpha}}{1-\alpha} + \frac{d_{s,a,h}}{k_{s,a,h}^\alpha} \\
&= \sum_{s,a} \|\omega\|_{\infty,h}^\alpha \left(\frac{[(k_{s,a,h} - 1)\|\omega\|_{\infty,h}]^{1-\alpha}}{1-\alpha} + \frac{d_{s,a,h}}{(k_{s,a,h}\|\omega\|_{\infty,h})^\alpha} \right) \\
&\leq \sum_{s,a} \|\omega\|_{\infty,h}^\alpha \left(\frac{[(k_{s,a,h} - 1)\|\omega\|_{\infty,h}]^{1-\alpha}}{1-\alpha} + \frac{d_{s,a,h}}{c_h(s,a)^\alpha} \right). \tag{87}
\end{aligned}$$

Here the last inequality is because $k_{s,a,h}\|\omega\|_{\infty,h} \geq c_h(s,a)$. The second inequality is because for any $0 < y < x$ and $\alpha \in (0, 1)$, we have:

$$\frac{x-y}{x^\alpha} \leq \frac{1}{1-\alpha}(x^{1-\alpha} - y^{1-\alpha}).$$

Then, let $x = i$ and $y = i - 1$, it holds that:

$$\frac{1}{i^\alpha} \leq \frac{1}{1-\alpha}(i^{1-\alpha} - (i-1)^{1-\alpha}).$$

Also let $x = c_h(s,a)$ and $y = (k_{s,a,h} - 1)\|\omega\|_{\infty,h}$, we have:

$$\frac{d_{s,a,h}}{c_h(s,a)^\alpha} + \frac{[(k_{s,a,h} - 1)\|\omega\|_{\infty,h}]^{1-\alpha}}{1-\alpha} \leq \frac{c_h(s,a)^{1-\alpha}}{1-\alpha}.$$

Applying this inequality to Equation (87), we have:

$$\sum_{k=1}^K \frac{\omega_{h,k}}{N_h^k(s_h^k, a_h^k)^\alpha} \leq \sum_{s,a} \|\omega\|_{\infty,h}^\alpha \frac{c_h(s,a)^{1-\alpha}}{1-\alpha} \leq \frac{1}{1-\alpha} (SA\|\omega\|_{\infty,h})^\alpha \|\omega\|_{1,h}^{1-\alpha}$$

The last inequality is by Hölder's inequality, as $\sum_{s,a} c_h(s,a)^{1-\alpha} \leq (SA)^\alpha \|\omega\|_{1,h}^{1-\alpha}$.
For $\alpha = 1$, it holds that:

$$\sum_{k=1}^K \frac{1}{N_h^k(s_h^k, a_h^k)} = \sum_{s,a} \sum_{i=1}^{N_h^K(s,a)} \frac{1}{i} \leq \sum_{s,a} (\log(N_h^K(s,a)) + 1) \lesssim SA \log T.$$

□

D.3 STEP 1: BOUNDING $Q_h^k - Q_h^*$

D.3.1 BOUNDING THE EMPIRICAL ESTIMATION ERRORS

By \mathcal{E}_6 in Lemma D.2 we have:

$$\left(\hat{\mathbb{P}}_{h,k}^{\text{adv}} - \hat{\mathbb{E}}_{h,k}^{\text{adv}}\right) \hat{V}_{h+1}^{\text{adv},k^n} = \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(\mathbb{P}_{s_h^k, a_h^k, h} - \mathbb{1}_{s_{h+1}^{k^n}}\right) \left(\hat{V}_{h+1}^{\mathbb{R},k^i} - V_{h+1}^*\right) \leq 2\sqrt{\frac{\beta^2 H \iota}{N_h^k(s_h^k, a_h^k)}}. \tag{88}$$

By \mathcal{E}_4 in Lemma D.2, it holds that:

$$\left(\hat{\mathbb{E}}_{h,k}^{\text{ref}} - \mathbb{P}_{h,k}^{\text{ref}}\right) \left(\hat{V}_{h+1}^{\mathbb{R},k^n} - V_{h+1}^*\right) = \sum_{i=1}^{N_h^k} u_i^{N_h^k} \left(\mathbb{1}_{s_{h+1}^{k^i}} - \mathbb{P}_{s_h^k, a_h^k, h}\right) \left(\hat{V}_{h+1}^{\mathbb{R},k^i} - V_{h+1}^*\right) \leq 2\sqrt{\frac{2\beta^2 \iota}{N_h^k(s_h^k, a_h^k)}}.$$

By \mathcal{E}_5 in Lemma D.2, it holds that:

$$\left(\hat{\mathbb{E}}_{h,k}^{\text{ref}} - \mathbb{P}_{h,k}^{\text{ref}}\right) V_{h+1}^* = \sum_{i=1}^{N_h^k} u_i^{N_h^k} \left(\mathbb{1}_{s_{h+1}^{k^i}} - \mathbb{P}_{s_h^k, a_h^k, h}\right) V_{h+1}^* \leq 8\sqrt{\frac{\mathbb{Q}^* \iota}{N_h^k(s_h^k, a_h^k)}} + 16\frac{H \iota}{N_h^k(s_h^k, a_h^k)}.$$

Therefore, combining these two inequalities, we have:

$$\left(\hat{\mathbb{E}}_{h,k}^{\text{ref}} - \mathbb{P}_{h,k}^{\text{ref}}\right) \hat{V}_{h+1}^{\text{R},k^n} \lesssim \sqrt{\frac{\mathbb{Q}^* + \beta^2}{N_h^k(s_h^k, a_h^k)} \iota} + \frac{H \iota}{N_h^k(s_h^k, a_h^k)}. \tag{89}$$

2160 D.3.2 BOUNDING THE BONUS
 2161

2162 Since the term ι^2 in the last inequality of Lemma 7 in Li et al. (2021) can be easily improved to ι ,
 2163 we can paraphrase the equation (87) and equation (88) of Li et al. (2021) to the following form:

2164 2165 $b_h^{R,k^n+1} = \left(1 - \frac{1}{\eta_n}\right) B_h^{R,k^n}(s_h^k, a_h^k) + \frac{1}{\eta_n} B_h^{R,k^n+1}(s_h^k, a_h^k) + \frac{c_b}{n^{3/4}} H^2 \iota.$ (90)
 2166

2167 This taken collectively with the definition of η_n^N allows us to expand
 2168

2169 2170 $R^{h,k} = \sum_{n=1}^{N_h^k} \eta_n^N b_h^{R,k^n+1}$
 2171
 2172 2173 $= \sum_{n=1}^{N_h^k} \eta_n \prod_{i=n+1}^{N_h^k} \left(1 - \frac{1}{\eta_i}\right) \left(\left(1 - \frac{1}{\eta_n}\right) B_h^{R,k^n}(s_h^k, a_h^k) + \frac{1}{\eta_n} B_h^{R,k^n+1}(s_h^k, a_h^k)\right) + c_b \sum_{n=1}^{N_h^k} \frac{\eta_n^{N_h^k}}{n^{3/4}} H^2 \iota$
 2174
 2175 2176 $= B_h^{R,k^{N_h^k}+1} + c_b \sum_{n=1}^{N_h^k} \frac{\eta_n^{N_h^k}}{n^{3/4}} H^2 \iota.$ (91)
 2177
 2178

2179 Then with $B_h^{R,k^{N_h^k}+1} = B_h^{R,k}$ and Equation (83) in Lemma D.1, it holds that
 2180

2181 2182 $R^{h,k} \lesssim B_h^{R,k} + \frac{H^2 \iota}{N_h^k(s_h^k, a_h^k)^{\frac{3}{4}}}.$ (92)
 2183

2184 Similar to equation (158) of Li et al. (2021), we have:

2185 2186 $\sqrt{\frac{\sigma_h^{\text{adv},k}(s_h^k, a_h^k) - (\mu_h^{\text{adv},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}} \leq \sqrt{\frac{\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{R,k^n}(s_{h+1}^{k^n}))^2}{N_h^k(s_h^k, a_h^k)}} \leq 2\beta$
 2187
 2188
 2189 2190 2191 2192 2193

Equation (93) is because $|V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{R,k^n}(s_{h+1}^{k^n})| \leq 2\beta$ by \mathcal{E}_3 in Lemma D.2 and $\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \leq 1$. Meanwhile, since $V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) \geq \hat{V}_{h+1}^{R,k^n}(s_{h+1}^{k^n})$, it also holds that

2194 2195 2196 2197 $\sqrt{\frac{\sigma_h^{\text{ref},k}(s_h^k, a_h^k) - (\mu_h^{\text{ref},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}} \leq \sqrt{\frac{J_1^{h,k} + J_2^{h,k}}{N_h^k(s_h^k, a_h^k)}},$

2198 where:

2199 2200 2201 $J_1^{h,k} = \frac{\sum_{n=1}^{N_h^k} \left((V_{h+1}^{R,k^n}(s_{h+1}^{k^n}))^2 - (\hat{V}_{h+1}^{R,k^n}(s_{h+1}^{k^n}))^2 \right)}{N_h^k(s_h^k, a_h^k)},$

2202 and

2203 2204 2205 $J_2^{h,k} = \frac{\sum_{n=1}^{N_h^k} (\hat{V}_{h+1}^{R,k^n}(s_{h+1}^{k^n}))^2}{N_h^k(s_h^k, a_h^k)} - \left(\frac{\sum_{n=1}^{N_h^k} \hat{V}_{h+1}^{R,k^n}(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} \right)^2.$

2206 Next we want to bound both $J_1^{h,k}$ and $J_2^{h,k}$.

2207 2208 2209 2210 2211 2212 2213 $J_1^{h,k} = \frac{\sum_{n=1}^{N_h^k} (V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) + \hat{V}_{h+1}^{R,k^n}(s_{h+1}^{k^n})) (V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) - \hat{V}_{h+1}^{R,k^n}(s_{h+1}^{k^n}))}{N_h^k(s_h^k, a_h^k)}$
 $\leq \frac{\sum_{n=1}^{N_h^k} 2H (V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) - \hat{V}_{h+1}^{R,k^n}(s_{h+1}^{k^n}))}{N_h^k(s_h^k, a_h^k)}.$

2214 Therefore, we have
 2215
 2216
 2217

$$J_1^{h,k} \leq \frac{2H\Psi_h^k(s_h^k, a_h^k)}{N_h^k(s_h^k, a_h^k)}, \quad (94)$$

2218 where
 2219
 2220

$$\Psi_h^k(s_h^k, a_h^k) = \sum_{n=1}^{N_h^k} \left(V_{h+1}^{\mathbb{R}, k^n}(s_{h+1}^{k^n}) - \hat{V}_{h+1}^{\mathbb{R}, k^n}(s_{h+1}^{k^n}) \right).$$

2221 For the second term $J_2^{h,k}$, because of Cauchy's Inequality, we have:
 2222

$$J_2^{h,k} = \frac{\sum_{n=1}^{N_h^k} \left(\hat{V}_{h+1}^{\mathbb{R}, k^n}(s_{h+1}^{k^n}) - \frac{\sum_{i=1}^{N_h^k} \hat{V}_{h+1}^{\mathbb{R}, k^n}(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} \right)^2}{N_h^k(s_h^k, a_h^k)} \leq 2(J_{2,1}^{h,k} + J_{2,2}^{h,k}),$$

2223 where:
 2224
 2225

$$J_{2,1}^{h,k} = \frac{\sum_{n=1}^{N_h^k} \left(\hat{V}_{h+1}^{\mathbb{R}, k^n}(s_{h+1}^{k^n}) - V_{h+1}^*(s_{h+1}^{k^n}) + \frac{\sum_{i=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} - \frac{\sum_{i=1}^{N_h^k} \hat{V}_{h+1}^{\mathbb{R}, k^n}(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} \right)^2}{N_h^k(s_h^k, a_h^k)},$$

2226 and
 2227

$$\begin{aligned} J_{2,2}^{h,k} &= \frac{\sum_{n=1}^{N_h^k} \left(V_{h+1}^*(s_{h+1}^{k^n}) - \frac{\sum_{i=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} \right)^2}{N_h^k(s_h^k, a_h^k)} \\ &= \frac{\sum_{n=1}^{N_h^k} (V_{h+1}^*(s_{h+1}^{k^n}))^2}{N_h^k(s_h^k, a_h^k)} - \left(\frac{\sum_{n=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} \right)^2. \end{aligned}$$

2228 Since $V_{h+1}^*(s_{h+1}^{k^n}) \leq \hat{V}_{h+1}^{\mathbb{R}, k^n}(s_{h+1}^{k^n}) \leq V_{h+1}^*(s_{h+1}^{k^n}) + \beta$, it holds that:
 2229

$$\begin{aligned} &\left| \hat{V}_{h+1}^{\mathbb{R}, k^n}(s_{h+1}^{k^n}) - V_{h+1}^*(s_{h+1}^{k^n}) + \frac{\sum_{i=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} - \frac{\sum_{i=1}^{N_h^k} \hat{V}_{h+1}^{\mathbb{R}, k^n}(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} \right| \\ &\leq \left| \hat{V}_{h+1}^{\mathbb{R}, k^n}(s_{h+1}^{k^n}) - V_{h+1}^*(s_{h+1}^{k^n}) \right| + \left| \frac{\sum_{i=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} - \frac{\sum_{i=1}^{N_h^k} \hat{V}_{h+1}^{\mathbb{R}, k^n}(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} \right| \leq 2\beta. \end{aligned}$$

2230 Therefore, applying this inequality to $J_{2,1}^{h,k}$, we have $J_{2,1}^{h,k} \leq 4\beta^2$. Moreover, according to equation
 2231 (165) of Li et al. (2021), the following inequality holds:
 2232

$$J_{2,2}^{h,k} \lesssim \mathbb{Q}^* + H^2 \sqrt{\frac{\ell}{N_h^k(s_h^k, a_h^k)}}.$$

2233 Combining the upper bounds of $J_1^{h,k}$ Equation (94), $J_{2,1}^{h,k}$ and $J_{2,2}^{h,k}$, we have:
 2234

$$\sqrt{\frac{\sigma_h^{\text{ref},k}(s_h^k, a_h^k) - (\mu_h^{\text{ref},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}} \lesssim \frac{\sqrt{H\Psi_h^k(s_h^k, a_h^k)}}{N_h^k(s_h^k, a_h^k)} + \sqrt{\frac{\mathbb{Q}^* + \beta^2}{N_h^k(s_h^k, a_h^k)}} + \frac{H\ell^{\frac{1}{4}}}{N_h^k(s_h^k, a_h^k)^{\frac{3}{4}}}. \quad (95)$$

2235 Back to the definition of $B_h^{\mathbb{R},k}$ in Algorithm 2, combining Equation (93) and Equation (95), it holds
 2236 that:
 2237

$$B_h^{\mathbb{R},k} \leq c_b \sqrt{\ell} \sqrt{\frac{\sigma_h^{\text{ref},k}(s_h^k, a_h^k) - (\mu_h^{\text{ref},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}} + c_b \sqrt{H\ell} \sqrt{\frac{\sigma_h^{\text{adv},k}(s_h^k, a_h^k) - (\mu_h^{\text{adv},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}}$$

$$\begin{aligned} & \lesssim \frac{\sqrt{H\Psi_h^k(s_h^k, a_h^k)\iota}}{N_h^k(s_h^k, a_h^k)} + \sqrt{\frac{(\mathbb{Q}^\star + \beta^2 H)\iota}{N_h^k(s_h^k, a_h^k)}} + \frac{H\iota^{\frac{3}{4}}}{N_h^k(s_h^k, a_h^k)^{\frac{3}{4}}}. \end{aligned}$$

Then by Equation (92), we have

$$R^{h,k} \lesssim \frac{\sqrt{H\Psi_h^k(s_h^k, a_h^k)\iota}}{N_h^k(s_h^k, a_h^k)} + \sqrt{\frac{(\mathbb{Q}^\star + \beta^2 H)\iota}{N_h^k(s_h^k, a_h^k)}} + \frac{H^2\iota}{N_h^k(s_h^k, a_h^k)^{\frac{3}{4}}}. \quad (96)$$

Applying Equation (88), Equation (89), Equation (96) to Equation (15), it holds that:

$$(Q_h^k - Q_h^\star)(s_h^k, a_h^k) \leq \hat{\mathbb{E}}_{h,k}^{\text{adv}}(V_{h+1}^{k^n} - V_{h+1}^\star) + \sqrt{\frac{(\mathbb{Q}^\star + \beta^2 H)\iota}{N_h^k(s_h^k, a_h^k)}} + \frac{H^2\iota}{N_h^k(s_h^k, a_h^k)^{\frac{3}{4}}} + R_{\text{else}}^{h,k}. \quad (97)$$

Here

$$R_{\text{else}}^{h,k} = \eta_0^{N_h^k} H + \hat{\mathbb{E}}_{h,k}^{\text{ref}}(V_h^{\text{R}, k^n} - \hat{V}_h^{\text{R}, k^n}) + \frac{\sqrt{H\Psi_h^k(s_h^k, a_h^k)\iota}}{N_h^k(s_h^k, a_h^k)} + \frac{H\iota}{N_h^k(s_h^k, a_h^k)}.$$

D.4 STEP 2: BOUNDING THE WEIGHTED SUM

D.4.1 REARRANGING THE SUMMATION

$$\begin{aligned} \sum_{k=1}^K \omega_{h,k} \hat{\mathbb{E}}_{h,k}^{\text{adv}}(V_{h+1}^{k^n} - V_{h+1}^\star) &= \sum_{k=1}^K \sum_{n=1}^{N_h^k} \omega_{h,k} \eta_n^{N_h^k} \left(V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^\star(s_{h+1}^{k^n}) \right) \\ &= \sum_{j=1}^K \left(\sum_{k=1}^K \sum_{n=1}^{N_h^k} \omega_{h,k} \eta_n^{N_h^k} \mathbb{I}[k^n = j] \right) \left(V_{h+1}^j(s_{h+1}^j) - V_{h+1}^\star(s_{h+1}^j) \right) \\ &\leq \sum_{j=1}^K \left(\sum_{k=1}^K \sum_{n=1}^{N_h^k} \omega_{h,k} \eta_n^{N_h^k} \mathbb{I}[k^n = j] \right) \left(Q_{h+1}^j - Q_{h+1}^\star \right) (s_{h+1}^j, a_{h+1}^j) \\ &\triangleq \sum_{j=1}^K \omega_{h+1,j}(h) \left(Q_{h+1}^j(s_{h+1}^j, a_{h+1}^j) - Q_{h+1}^\star(s_{h+1}^j, a_{h+1}^j) \right). \end{aligned} \quad (98)$$

Here, for any $j \in [K]$

$$\omega_{h+1,j}(h) = \sum_{k=1}^K \sum_{n=1}^{N_h^k} \omega_{h,k} \eta_n^{N_h^k} \mathbb{I}[k^n = j].$$

The inequality is because $Q_{h+1}^j(s_{h+1}^j, a_{h+1}^j) = V_{h+1}^j(s_{h+1}^j)$, $Q_{h+1}^\star(s_{h+1}^j, a_{h+1}^j) \leq V_{h+1}^\star(s_{h+1}^j)$.

D.4.2 PROOF OF EQUATION (22)

For any given h and non-negative constants $\{\omega_{h,k}\}_{h,[K]}$, we denote $\|\omega\|_{\infty,h} = \max_{k \in [K]} \omega_{h,k}$ and $\|\omega\|_{1,h} = \sum_{k \in [K]} \omega_{h,k}$. We also recursively define $\omega_{h',k}(h)$ for any $h \leq h' \leq H+1, k \in [K]$ as follows:

$$\omega_{h,k}(h) := \omega_{h,k}; \quad \omega_{h',j}(h) = \sum_{k=1}^K \sum_{n=1}^{N_h^k} \omega_{h'-1,k}(h) \eta_n^{N_h^k} \mathbb{I}[k^n = j], \quad \forall j \in [K], h < h' \leq H+1.$$

According to the definition of k^n , $\mathbb{I}[k^n = j] = 1$ if and only if $(s_h^j, a_h^j) = (s_h^k, a_h^k)$, $j \leq k-1$ and $n = N_h^{j+1}(s_h^j, a_h^j)$. Then by Equation (84) in Lemma D.1, we have:

$$\sum_{k=1}^K \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \mathbb{I}[k^n = j] = \sum_{k=j+1}^K \eta_{N_h^{j+1}}^{N_h^k} \mathbb{I}[(s_h^j, a_h^j) = (s_h^k, a_h^k)] \leq \sum_{t=N_h^{j+1}}^{\infty} \eta_{N_h^{j+1}}^t \leq 1 + \frac{1}{H}. \quad (99)$$

Therefore, for $h < h' \leq H + 1$, it holds that:

$$\omega_{h',j}(h) \leq \|\omega(h)\|_{\infty,h'-1} \sum_{k=1}^K \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \mathbb{I}[k^n = j] \leq (1 + \frac{1}{H}) \|\omega(h)\|_{\infty,h'-1}. \quad (100)$$

It also holds that:

$$\sum_{j=1}^K \omega_{h',j}(h) = \sum_{k=1}^K \omega_{h,k} \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \leq \|\omega(h)\|_{1,h'-1}. \quad (101)$$

Combining Equation (97) with Equation (98), the weighted sum $\sum_{k=1}^K \omega_{h,k}(Q_h^k - Q_h^*)(s_h^k, a_h^k)$ can be bounded by

$$\begin{aligned} & \sum_{k=1}^K \omega_{h,k}(Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)) \\ & \lesssim \sum_{k=1}^K \omega_{h+1,k}(h)(Q_{h+1}^k - Q_{h+1}^*)(s_{h+1}^k, a_{h+1}^k) + \sum_{k=1}^K \omega_{h,k} \left(\sqrt{\frac{(\mathbb{Q}^* + \beta^2 H)\iota}{N_h^k(s_h^k, a_h^k)}} + \frac{H^2\iota}{(N_h^k)^{\frac{3}{4}}} + R_{\text{else}}^{h,k} \right) \\ & \leq \sum_{k=1}^K \omega_{h+1,k}(h)(Q_{h+1}^k - Q_{h+1}^*)(s_{h+1}^k, a_{h+1}^k) + \sqrt{(\mathbb{Q}^* + \beta^2)SA\|\omega\|_{\infty,h}\|\omega\|_{1,h}\iota} \\ & \quad + H^2\iota(SA\|\omega\|_{\infty,h})^{\frac{3}{4}}\|\omega\|_{1,h}^{\frac{1}{4}} + \sum_{k=1}^K \omega_{h+1,k}(h)R_{\text{else}}^{h,k}. \end{aligned} \quad (102)$$

The last inequality is by Lemma D.3 with $\alpha = \frac{1}{2}$ and $\frac{3}{4}$. Recurring Equation (102) with regard to $h, h+1, \dots, H$, since $Q_{H+1}^k(s, a) = Q_{H+1}^*(s, a) = 0$ and the weight recursions relationship Equation (100) and Equation (101), we have

$$\begin{aligned} & \sum_{k=1}^K \omega_{h,k}(Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)) \\ & \lesssim H \sqrt{(\mathbb{Q}^* + \beta^2 H)SA\|\omega\|_{\infty,h}\|\omega\|_{1,h}\iota} + H^3\iota(SA\|\omega\|_{\infty,h})^{\frac{3}{4}}\|\omega\|_{1,h}^{\frac{1}{4}} + \sum_{k=1}^K \sum_{h'=h}^H \omega_{h',k}(h)R_{\text{else}}^{h',k}. \end{aligned} \quad (103)$$

D.5 STEP 3: INTEGRATING MULTIPLE WEIGHTED SUMS

D.5.1 PROOF OF EQUATION (26)

For any $N = \lceil \log_2(H) \rceil$, $t \in [N]$, $k \in [K]$ and the given $h \in [H]$, let:

$$\omega_{h,k}^{(i)} = \mathbb{I}[Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \in [2^{i-1}\Delta_{\min}, 2^i\Delta_{\min})],$$

and then

$$\|\omega\|_{\infty,h}^{(i)} = \max_k \omega_{h,k}^{(i)} \leq 1, \quad \|\omega\|_{1,h}^{(i)} = \sum_{k=1}^K \omega_{h,k}^{(i)}.$$

For any given $i \in [N]$ and $h \leq h' \leq H$ and the weight $\{\omega_{h',k}^{(i)}\}_k$ can be defined recursively by Equation (19). Therefore, for any $j \in [K]$, it holds that:

$$\sum_{i=1}^N \omega_{h'+1,j}^{(i)}(h) = \sum_{k=1}^K \sum_{n=1}^{N_h^k} \left(\sum_{i=1}^N \omega_{h',k}^{(i)}(h) \right) \eta_n^{N_h^k} \mathbb{I}[k^n = j].$$

Here for any $i \in [N]$, $\omega_{h,k}^{(i)}(h) = \omega_{h,k}^{(i)}$. Then by mathematical induction on $h' \in [h, H]$, it is straightforward to prove that for any $j \in [K]$,

$$\sum_{i=1}^N \omega_{h',j}^{(i)}(h) \leq \left(1 + \frac{1}{H}\right)^{h'-h} < 3, \quad (104)$$

given that for any $j \in [K]$

$$\sum_{k=1}^K \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \mathbb{I}[k^n = j] \leq 1 + \frac{1}{H}$$

by Equation (99) and $\sum_{i=1}^N \omega_{h,j}^{(i)}(h) = \sum_{i=1}^N \omega_{h,j}^{(i)} \leq 1$.

Applying the weight $\{\omega_{h,k}^{(i)}\}_k$ to Equation (103), since $\|\omega\|_{\infty,h} \leq 1$, then for any $i \in [N]$, it holds that:

$$\begin{aligned} \sum_{k=1}^K \omega_{h,k}^{(i)} (Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)) &\lesssim H \sqrt{(\mathbb{Q}^* + \beta^2 H) S A \|\omega^{(i)}\|_{\infty,h} \|\omega^{(i)}\|_{1,h} \iota} \\ &+ H^3 \iota (S A \|\omega^{(i)}\|_{\infty,h})^{\frac{3}{4}} (\|\omega^{(i)}\|_{1,h})^{\frac{1}{4}} + \sum_{k=1}^K \sum_{h'=h}^H \omega_{h',k}^{(i)}(h) R_{\text{else}}^{h',k}. \end{aligned}$$

On the other hand, according to the definition of $\omega_{h,k}^{(i)}$,

$$\sum_{k=1}^K \omega_{h,k}^{(i)} (Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)) \geq 2^{i-1} \Delta_{\min} \|\omega\|_{1,h}^{(i)}.$$

Therefore, since $\|\omega^{(i)}\|_{\infty,h} \leq 1$, we obtain the following inequality for any $i \in [N]$:

$$\begin{aligned} 2^{i-1} \Delta_{\min} \|\omega^{(i)}\|_{1,h} &\lesssim H \sqrt{(\mathbb{Q}^* + \beta^2 h) S A \|\omega^{(i)}\|_{1,h} \iota} + H^3 \iota (S A)^{\frac{3}{4}} (\|\omega^{(i)}\|_{1,h})^{\frac{1}{4}} \\ &+ \sum_{k=1}^K \sum_{h'=h}^H \omega_{h',k}^{(i)}(h) R_{\text{else}}^{h',k}. \end{aligned} \quad (105)$$

Then at least one of the following three inequalities holds:

$$\begin{aligned} 2^{i-1} \Delta_{\min} \|\omega^{(i)}\|_{1,h} &\lesssim H \sqrt{(\mathbb{Q}^* + \beta^2 H) S A \|\omega^{(i)}\|_{1,h} \iota} \\ 2^{i-1} \Delta_{\min} \|\omega^{(i)}\|_{1,h} &\lesssim H^3 \iota (S A)^{\frac{3}{4}} (\|\omega^{(i)}\|_{1,h})^{\frac{1}{4}}, \\ 2^{i-1} \Delta_{\min} \|\omega^{(i)}\|_{1,h} &\lesssim \sum_{k=1}^K \sum_{h'=h}^H \omega_{h',k}^{(i)}(h) R_{\text{else}}^{h',k}. \end{aligned}$$

Solving this three inequalities, we know that:

$$\begin{aligned} \|\omega^{(i)}\|_{1,h} &\leq O \left(\max \left\{ \frac{(\mathbb{Q}^* + \beta^2 H) S A H^2 \iota}{4^{i-1} \Delta_{\min}^2}, \frac{H^4 S A \iota^{\frac{4}{3}}}{(2^{i-1} \Delta_{\min})^{\frac{4}{3}}}, \frac{\sum_{k=1}^K \sum_{h'=h}^H \omega_{h',k}^{(i)}(h) R_{\text{else}}^{h',k}}{2^{i-1} \Delta_{\min}} \right\} \right) \\ &\leq O \left(\frac{(\mathbb{Q}^* + \beta^2 H) S A H^2 \iota}{4^{i-1} \Delta_{\min}^2} + \frac{H^4 S A \iota^{\frac{4}{3}}}{(2^{i-1} \Delta_{\min})^{\frac{4}{3}}} + \frac{\sum_{k=1}^K \sum_{h'=h}^H \omega_{h',k}^{(i)}(h) R_{\text{else}}^{h',k}}{2^{i-1} \Delta_{\min}} \right). \end{aligned} \quad (106)$$

By Equation (104), we have:

$$\sum_{i=1}^N \sum_{k=1}^K \sum_{h'=h}^H \omega_{h',k}^{(i)}(h) R_{\text{else}}^{h',k} = \sum_{h'=h}^H \sum_{k=1}^K \left(\sum_{i=1}^N \omega_{h',k}^{(i)}(h) \right) R_{\text{else}}^{h',k} \leq 3 \sum_{h'=1}^H \sum_{k=1}^K R_{\text{else}}^{h',k}.$$

Using this inequality, we have

$$\sum_{i=1}^N 2^i \Delta_{\min} \|\omega^{(i)}\|_{1,h} \leq O \left(\frac{(\mathbb{Q}^* + \beta^2 H) S A H^2 \iota}{\Delta_{\min}} + \frac{H^4 S A \iota^{\frac{4}{3}}}{(\Delta_{\min})^{\frac{1}{3}}} + \sum_{h'=1}^H \sum_{k=1}^K R_{\text{else}}^{h',k} \right). \quad (107)$$

2430 D.5.2 PROOF OF EQUATION (27) AND EQUATION (28)

2431
2432 Next we will bound the term $\sum_{h'=1}^H \sum_{k=1}^K R_{\text{else}}^{h',k}$, where

2433
2434
2435
$$R_{\text{else}}^{h,k} = \eta_0^{N_h^k} H + \hat{\mathbb{E}}_{h,k}^{\text{ref}} \left(V_{h+1}^{\text{R},k^n} - \hat{V}_{h+1}^{\text{R},k^n} \right) + \left(\mathbb{P}_{h,k}^{\text{ref}} \hat{V}_{h+1}^{\text{R},k^n} - \mathbb{P}_{h,k}^{\text{adv}} \hat{V}_{h+1}^{\text{R},k^n} \right) + \frac{\sqrt{H\Psi_h^k \iota}}{N_h^k} + \frac{H\iota}{N_h^k}.$$

2436 According to equation (149) of Li et al. (2021), we have:

2437
2438
2439
$$\sum_{h'=1}^H \sum_{k=1}^K \eta_0^{N_h^k} H \leq H^2 SA \leq \frac{H^6 SA \log(T)\iota}{\beta}. \quad (108)$$

2440
2441 Since for any $k \in [K]$, $V_{h'+1}^{\text{R},k}(s) - \hat{V}_{h'+1}^{\text{R},k}(s) \geq 0$, we have

2442
2443
2444
$$\begin{aligned} \sum_{h'=1}^H \sum_{k=1}^K \hat{\mathbb{E}}_{h',k}^{\text{ref}} \left(V_{h'+1}^{\text{R},k^n} - \hat{V}_{h'+1}^{\text{R},k^n} \right) &\leq \sum_{h'=1}^H \sum_{k=1}^K \sum_{n=1}^{N_h^k} u_n^{N_h^k} \left(V_{h'+1}^{\text{R},k^n} - \hat{V}_{h'+1}^{\text{R},k^n} \right) (s_{h'+1}^{k^n}) \left(\sum_{j=1}^K \mathbb{I}[k^n = j] \right) \\ &= \sum_{h'=1}^H \sum_{j=1}^K \left(\sum_{k=1}^K \sum_{n=1}^{N_h^k} u_i^{N_h^k} \mathbb{I}[k^n = j] \right) (V_{h'+1}^{\text{R},j} - \hat{V}_{h'+1}^{\text{R},j})(s_{h'+1}^{k^n}). \end{aligned}$$

2445
2446 Here $\mathbb{I}[k^n = j] = 1$ if and only if $(s_{h'}^j, a_{h'}^j) = (s_{h'}^k, a_{h'}^k)$, $j \leq k-1$ and $n = N_{h'}^{j+1}(s_{h'}^j, a_{h'}^j) > 0$.
2447 Then we have:

2448
2449
2450
$$\sum_{k=1}^K \sum_{n=1}^{N_h^k} u_i^{N_h^k} \mathbb{I}[k^n = j] = \sum_{k=j+1}^K u_{N_{h'}^{j+1}}^{N_{h'}^k} \mathbb{I}\left[(s_{h'}^j, a_{h'}^j) = (s_{h'}^k, a_{h'}^k)\right] \leq \sum_{t=N_{h'}^{j+1}}^{N_{h'}^K} u_{N_{h'}^{j+1}}^t \lesssim \log T. \quad (109)$$

2451
2452 The last inequality is because for any $N \in \mathbb{N}_+$ and $i \in [N]$, $u_i^N \leq \frac{2}{N}$. Therefore it holds that:

2453
2454
2455
$$\sum_{h'=1}^H \sum_{k=1}^K \hat{\mathbb{E}}_{h',k}^{\text{ref}} \left(V_{h'}^{\text{R},k^n} - \hat{V}_{h'}^{\text{R},k^n} \right) \lesssim \log T \sum_{h'=1}^H \sum_{j=1}^K \left(V_{h'+1}^{\text{R},j} - \hat{V}_{h'+1}^{\text{R},j} \right) (s_{h'+1}^{k^i}). \quad (110)$$

2456 To continue, we will first prove a lemma

2457 **Lemma D.4.** For any $h' \in [H]$ and $k \in [K]$,

- 2458
2459 • If $V_{h'+1}^k(s) - V_{h'+1}^{\text{LCB},k}(s) \leq \beta$, then $V_{h'+1}^{\text{R},K+1}(s) = V_{h'+1}^{\text{R},k}(s) = \hat{V}_{h'+1}^{\text{R},k}(s)$.
- 2460
2461 • If $V_{h'+1}^k(s) - V_{h'+1}^{\text{LCB},k}(s) > \beta$, then we have:

2462
2463
$$0 \leq V_{h'+1}^{\text{R},k}(s) - \hat{V}_{h'+1}^{\text{R},k}(s) \leq V_{h'+1}^k(s) - V_{h'+1}^{\text{LCB},k}(s),$$

2464 and

2465
$$|\hat{V}_{h'+1}^{\text{R},k}(s) - V_{h'+1}^{\text{R},K+1}(s)| \leq V_{h'+1}^k(s) - V_{h'+1}^{\text{LCB},k}(s).$$

2466 *Proof.* • If for given $k \in [K]$, $V_{h'+1}^k(s) - V_{h'+1}^{\text{LCB},k}(s) \leq \beta$, then there exists $k_1 \in [K]$ such
2467 that:

2468
2469
$$k_1 = \min \left\{ k : V_{h'+1}^k(s) - V_{h'+1}^{\text{LCB},k}(s) \leq \beta \right\}.$$

2470 Then according the algorithm, we have $u_{\text{ref}}^{k_1}(s) = \text{TRUE}$, or it is contradictory to the mini-
2471 mality of k_1 . Therefore, in this case we have:

2472
2473
$$V_{h'+1}^{\text{R},K+1}(s) = V_{h'+1}^{\text{R},k}(s) = V_{h'+1}^{\text{R},k_1}(s) = V_{h'+1}^{k_1}(s) \leq V_{h'+1}^{\text{LCB},k_1}(s) + \beta \leq V_{h'+1}^{\star}(s) + \beta,$$

2474 and

2475
$$V_{h'+1}^{\text{R},k}(s) = V_{h'+1}^{\text{R},k_1}(s) = V_{h'+1}^{k_1}(s) \geq V_{h'+1}^{\star}(s).$$

2476 According to the definition of $\hat{V}_{h'+1}^{\text{R},k}(s)$, we have $\hat{V}_{h'+1}^{\text{R},k}(s) = V_{h'+1}^{\text{R},k}(s) = V_{h'+1}^{\text{R},K+1}(s)$.

2477 Thus $V_{h'+1}^k(s) - V_{h'+1}^{\text{LCB},k}(s) \leq \beta$ is the sufficient condition of $V_{h'+1}^{\text{R},k}(s) = \hat{V}_{h'+1}^{\text{R},k}(s) =$
2478 $V_{h'+1}^{\text{R},K+1}(s)$.

- 2484 • Moreover, if $V_{h'+1}^k(s) - V_{h'+1}^{\text{LCB},k}(s) > \beta$, according to the algorithm, we have $V_{h'+1}^{\text{R},k}(s) =$
 2485 $V_{h'+1}^k(s)$ and then $0 \leq V_{h'+1}^{\text{R},k}(s) - \hat{V}_{h'+1}^{\text{R},k}(s) \leq V_{h'+1}^k(s) - V_{h'+1}^{\text{LCB},k}(s)$.
 2486

2487 In this case, we also have $V_{h'+1}^{\text{LCB},k}(s) \leq V_{h'+1}^{\text{R},K+1}(s) \leq V_{h'+1}^k(s)$ and then $V_{h'+1}^{\text{LCB},k}(s) \leq$
 2488 $\hat{V}_{h'+1}^{\text{R},k}(s) \leq V_{h'+1}^{\text{R},k}(s) = V_{h'+1}^k(s)$. These two inequalities imply that $|\hat{V}_{h'+1}^{\text{R},j}(s) -$
 2489 $V_{h'+1}^{\text{R},K+1}(s)| \leq V_{h'+1}^k(s) - V_{h'+1}^{\text{LCB},k}(s)$.
 2490

□

2491
 2492 According to this lemma, the following inequality holds:
 2493

2494
 2495
$$\sum_{h'=1}^H \sum_{j=1}^K \left(V_{h'+1}^{\text{R},j}(s_{h'+1}^j) - \hat{V}_{h'+1}^{\text{R},j}(s_{h'+1}^j) \right)$$

 2496
 2497
$$\leq \sum_{h'=1}^H \sum_{j=1}^K \left(V_{h'+1}^k(s_{h'+1}^j) - V_{h'+1}^{\text{LCB},k}(s_{h'+1}^j) \right) \mathbb{I} [V_{h'+1}^j(s_{h'+1}^j) - V_{h'+1}^{\text{LCB},j}(s_{h'+1}^j) > \beta] \lesssim \frac{H^6 S A \iota}{\beta}.$$

 2498
 2499
 2500

2501 Applying this inequality to Equation (110), it holds that:
 2502

2503
$$\sum_{h'=1}^H \sum_{k=1}^K \hat{\mathbb{E}}_{h',k}^{\text{ref}} \left(V_{h'}^{\text{R},k^n} - \hat{V}_{h'}^{\text{R},k^n} \right) \lesssim \frac{H^6 S A \log(T) \iota}{\beta} \quad (111)$$

 2504
 2505

2506 For the third term in $R_{\text{else}}^{h',k}$, because $\sum_{n=1}^{N_{h'}^k} u_n^{N_{h'}^k} = \sum_{n=1}^{N_{h'}^k} \eta_n^{N_{h'}^k}$, then
 2507

2508
$$\begin{aligned} & \mathbb{P}_{h,k}^{\text{ref}} \hat{V}_{h'+1}^{\text{R},k^n} - \mathbb{P}_{h,k}^{\text{adv}} \hat{V}_{h'+1}^{\text{R},k^n} \\ &= \sum_{n=1}^{N_{h'}^k} u_n^{N_{h'}^k} \mathbb{P}_{s_{h'}, a_{h'}, h'} (\hat{V}_{h'+1}^{\text{R},k^n} - V_{h'+1}^{\text{R},K+1}) - \sum_{n=1}^{N_{h'}^k} \eta_n^{N_{h'}^k} \mathbb{P}_{s_{h'}, a_{h'}, h'} (\hat{V}_{h'+1}^{\text{R},k^n} - V_{h'+1}^{\text{R},K+1}) \\ &\leq \sum_{n=1}^{N_{h'}^k} u_n^{N_{h'}^k} \mathbb{P}_{s_{h'}, a_{h'}, h'} |\hat{V}_{h'+1}^{\text{R},k^n} - V_{h'+1}^{\text{R},K+1}| + \sum_{n=1}^{N_{h'}^k} \eta_n^{N_{h'}^k} \mathbb{P}_{s_{h'}, a_{h'}, h'} |\hat{V}_{h'+1}^{\text{R},k^n} - V_{h'+1}^{\text{R},K+1}| \end{aligned}$$

2516 Similar to Equation (110), we have:
 2517

2518
$$\sum_{h'=1}^H \sum_{k=1}^K \sum_{n=1}^{N_{h'}^k} u_n^{N_{h'}^k} \mathbb{P}_{s_{h'}, a_{h'}, h'} |\hat{V}_{h'+1}^{\text{R},k^n} - V_{h'+1}^{\text{R},K+1}| \lesssim \log(T) \sum_{h'=1}^H \sum_{j=1}^K \mathbb{P}_{s_{h'}, a_{h'}, h'} |\hat{V}_{h'+1}^{\text{R},j} - V_{h'+1}^{\text{R},K+1}|.$$

 2519
 2520

2521 and
 2522

2523
$$\sum_{h'=1}^H \sum_{k=1}^K \sum_{n=1}^{N_{h'}^k} \eta_n^{N_{h'}^k} \mathbb{P}_{s_{h'}, a_{h'}, h'} |\hat{V}_{h'+1}^{\text{R},k^n} - V_{h'+1}^{\text{R},K+1}| \lesssim \sum_{h'=1}^H \sum_{j=1}^K \mathbb{P}_{s_{h'}, a_{h'}, h'} |\hat{V}_{h'+1}^{\text{R},j} - V_{h'+1}^{\text{R},K+1}|$$

 2524
 2525

2526 by Equation (99). Combining these two inequalities, we have:
 2527

2528
$$\sum_{h'=1}^H \sum_{k=1}^K \left(\mathbb{P}_{h,k}^{\text{ref}} \hat{V}_{h'+1}^{\text{R},k^n} - \mathbb{P}_{h,k}^{\text{adv}} \hat{V}_{h'+1}^{\text{R},k^n} \right) \lesssim \log(T) \sum_{h'=1}^H \sum_{j=1}^K \mathbb{P}_{s_{h'}, a_{h'}, h'} |\hat{V}_{h'+1}^{\text{R},j} - V_{h'+1}^{\text{R},K+1}|. \quad (112)$$

 2529
 2530

2531 According to Lemma D.4, the following inequality holds:
 2532

2533
$$\begin{aligned} & \sum_{h'=1}^H \sum_{j=1}^K |\hat{V}_{h'+1}^{\text{R},j}(s) - V_{h'+1}^{\text{R},K+1}(s)| \\ &\leq \sum_{h'=1}^H \sum_{j=1}^K \left(V_{h'+1}^k(s_{h'+1}^j) - V_{h'+1}^{\text{LCB},k}(s_{h'+1}^j) \right) \mathbb{I} [V_{h'+1}^j(s_{h'+1}^j) - V_{h'+1}^{\text{LCB},j}(s_{h'+1}^j) > \beta] \lesssim \frac{H^6 S A \iota}{\beta}. \end{aligned}$$

 2534
 2535
 2536
 2537

Combining Equation (112) with the event \mathcal{E}_7 in Lemma D.2, we have:

$$\sum_{h'=1}^H \sum_{k=1}^K \left(\mathbb{P}_{h,k}^{\text{ref}} \hat{V}_{h'+1}^{\mathbf{R},k^n} - \mathbb{P}_{h,k}^{\text{adv}} \hat{V}_{h'+1}^{\mathbf{R},k^n} \right) \lesssim \frac{H^6 S A \log(T) \iota}{\beta}. \quad (113)$$

Now we move to the fourth term in $R_{\text{else}}^{h,k}$. By Lemma D.4 we have:

$$\begin{aligned} \Psi_{h'}^k(s_{h'}^k, a_{h'}^k) &= \sum_{n=1}^{N_{h'}^k} \left(V_{h'+1}^{\mathbf{R},k^n}(s_{h'+1}^{k^n}) - \hat{V}_{h'+1}^{\mathbf{R},k^n}(s_{h'+1}^{k^n}) \right) \\ &\leq \sum_{n=1}^{N_{h'}^k} \left(V_{h'+1}^{k^n}(s_{h'+1}^{k^n}) - V_{h'+1}^{\text{LCB},k^n}(s_{h'+1}^{k^n}) \right) \cdot \mathbb{I} \left[V_{h'+1}^{k^n}(s_{h'+1}^{k^n}) - V_{h'+1}^{\text{LCB},k^n}(s_{h'+1}^{k^n}) > \beta \right] \\ &\triangleq \Phi_{h'}^k(s_{h'}^k, a_{h'}^k) \end{aligned}$$

Then it holds that:

$$\begin{aligned} \sum_{k=1}^K \frac{\sqrt{\Psi_{h'}^k(s_{h'}^k, a_{h'}^k)}}{N_{h'}^k(s_{h'}^k, a_{h'}^k)} &\leq \sum_{k=1}^K \frac{\sqrt{\Phi_{h'}^k(s_{h'}^k, a_{h'}^k)}}{N_{h'}^k(s_{h'}^k, a_{h'}^k)} = \sum_{s,a} \sum_{n=1}^{N_{h'}^K(s,a)} \frac{\sqrt{\Phi_{h'}^{k^n}(s,a)}}{n} \\ &\leq \log T \sum_{s,a} \sqrt{\Phi_{h'}^K(s,a)} \leq \log T \sqrt{S A \sum_{s,a} \Phi_{h'}^K(s,a)} \end{aligned} \quad (114)$$

The first inequality is because of the monotonicity of $\Phi_{h'}^n(s,a)$. The second inequality is by Cauchy's inequality. To continue, note that:

$$\begin{aligned} &\sum_{h'=1}^H \sqrt{\sum_{s,a} \Phi_{h'}^{N_{h'}^K(s,a)}(s,a)} \\ &= \sum_{h'=1}^H \sqrt{\sum_{k=1}^K \left(V_{h'+1}^k(s_{h'+1}^k) - V_{h'+1}^{\text{LCB},k}(s_{h'+1}^k) \right) \cdot \mathbb{I} \left[V_{h'+1}^k(s_{h'+1}^k) - V_{h'+1}^{\text{LCB},k}(s_{h'+1}^k) > \beta \right]} \\ &\leq \sqrt{H \sum_{h'=1}^H \sum_{k=1}^K \left(V_{h'+1}^k(s_{h'+1}^k) - V_{h'+1}^{\text{LCB},k}(s_{h'+1}^k) \right) \cdot \mathbb{I} \left[V_{h'+1}^k(s_{h'+1}^k) - V_{h'+1}^{\text{LCB},k}(s_{h'+1}^k) > \beta \right]} \\ &\leq \sqrt{\frac{H^7 S A \iota}{\beta}} \end{aligned}$$

Together with Equation (114), it holds:

$$\sum_{h'=1}^H \sum_{k=1}^K \frac{\sqrt{H \Psi_{h'}^k(s_{h'}^k, a_{h'}^k) \iota}}{N_{h'}^k(s_{h'}^k, a_{h'}^k)} \leq \log T \sum_{h'=1}^H \sqrt{S A \sum_{s,a} \Phi_{h'}^{N_{h'}^K(s,a)}(s,a) \iota} \lesssim \frac{H^{\frac{7}{2}} S A \log(T) \iota}{\sqrt{\beta}}. \quad (115)$$

By Lemma D.3 with $\alpha = 1$, we have:

$$\sum_{h'=1}^H \sum_{k=1}^K \frac{H \iota}{N_{h'}^k(s_{h'}^k, a_{h'}^k)} \leq H^2 S A \log(T) \iota. \quad (116)$$

By summing Equation (108), Equation (111), Equation (113), Equation (115) and Equation (116), since $\beta \in (0, H]$, we can conclude that:

$$\sum_{h'=1}^H \sum_{k=1}^K R_{\text{else}}^{h',k} \lesssim \frac{H^6 S A \log(T) \iota}{\beta}.$$

Then we have

$$\sum_{k=1}^K \text{clip}[(Q_h^k - Q_h^*)(s_h^k, a_h^k) \mid \Delta_{\min}] = O \left(\frac{(\mathbb{Q}^* + \beta^2 H) S A H^2 \iota}{\Delta_{\min}} + \frac{H^4 S A \iota^{\frac{4}{3}}}{(\Delta_{\min})^{\frac{1}{3}}} + \frac{H^6 S A \log(T) \iota}{\beta} \right)$$

$$\leq O\left(\frac{(\mathbb{Q}^* + \beta^2 H) SAH^2 \iota}{\Delta_{\min}} + \frac{H^6 SA \iota^2}{\beta}\right). \quad (117)$$

The last inequality is because

$$\frac{H^4 SA \iota^{\frac{4}{3}}}{(\Delta_{\min})^{\frac{1}{3}}} \leq \frac{\beta^2 H^3 SA \iota}{\Delta_{\min}} + \frac{H^5 SA \iota}{\beta} + \frac{H^5 SA \iota^2}{\beta}$$

by AM-GM inequality.

D.6 STEP 4: BOUNDING THE EXPECTED GAP-DEPENDENT REGRET

By Equation (9), $Q_h^k(s_h^k, a_h^k) = V_h^k(s_h^k) \geq V_h^*(s_h^k)$. Thus, for any episode-step pair (k, h)

$$\Delta_h(x_h^k, a_h^k) = \text{clip}[V_h^*(x_h^k) - Q_h^*(x_h^k, a_h^k) | \Delta_{\min}] \leq \text{clip}[(Q_h^k - Q_h^*)(x_h^k, a_h^k) | \Delta_{\min}].$$

By Equation (4) in Yang et al. (2021), we have $\mathbb{E}(\text{Regret}(K)) = \mathbb{E}\left[\sum_{k=1}^K \sum_{h=1}^H \Delta_h(x_h^k, a_h^k)\right]$, which further implies

$$\mathbb{E}(\text{Regret}(K)) \leq \mathbb{E}\left[\sum_{k=1}^K \sum_{h=1}^H \text{clip}[(Q_h^k - Q_h^*)(x_h^k, a_h^k) | \Delta_{\min}]\right].$$

Finally, let $\delta = \frac{1}{6T}$ and $\mathcal{E} = \bigcap_{i=1}^6 \mathcal{E}_i$. Then the event \mathcal{E} holds with probability at least $1 - 6\delta = 1 - \frac{1}{T}$. Then we have:

$$\begin{aligned} \mathbb{E}(\text{Regret}(K)) &\leq \mathbb{E}\left[\sum_{k=1}^K \sum_{h=1}^H \text{clip}[(Q_h^k - Q_h^*)(x_h^k, a_h^k) | \Delta_{\min}] \middle| \mathcal{E}\right] \mathbb{P}(\mathcal{E}) \\ &\quad + \mathbb{E}\left[\sum_{k=1}^K \sum_{h=1}^H \text{clip}[(Q_h^k - Q_h^*)(x_h^k, a_h^k) | \Delta_{\min}] \middle| \mathcal{E}^c\right] \mathbb{P}(\mathcal{E}^c) \\ &\leq O\left(\frac{(\mathbb{Q}^* + \beta^2 H) H^3 SA \iota}{\Delta_{\min}} + \frac{H^6 SA \iota^2}{\beta}\right) + (1 - \frac{1}{T}) HT \\ &= O\left(\frac{(\mathbb{Q}^* + \beta^2 H) H^3 SA \iota}{\Delta_{\min}} + \frac{H^6 SA \iota^2}{\beta}\right). \end{aligned} \quad (118)$$

The last inequality is because under the event \mathcal{E} , we have proved that

$$\sum_{k=1}^K \sum_{h=1}^H \text{clip}[(Q_h^k - Q_h^*)(x_h^k, a_h^k) | \Delta_{\min}] \leq O\left(\frac{(\mathbb{Q}^* + \beta^2 H) H^3 SA \iota}{\Delta_{\min}} + \frac{H^6 SA \iota^2}{\beta}\right)$$

by Equation (117) and under the event \mathcal{E}^c ,

$$\sum_{k=1}^K \sum_{h=1}^H \text{clip}[(Q_h^k - Q_h^*)(x_h^k, a_h^k) | \Delta_{\min}] \leq HT.$$

E RELATED WORK

On-policy RL for finite-horizon tabular MDPs with worst-case regret. There are mainly two types of algorithms for reinforcement learning: model-based and model-free learning. Model-based algorithms learn a model from past experience and make decisions based on this model, while model-free algorithms only maintain a group of value functions and take the induced optimal actions. Due to these differences, model-free algorithms are usually more space-efficient and time-efficient compared to model-based algorithms. However, model-based algorithms may achieve better learning performance by leveraging the learned model.

Next, we discuss the literature on model-based and model-free algorithms for finite-horizon tabular MDPs with worst-case regret. Auer et al. (2008), Agrawal & Jia (2017), Azar et al. (2017), Kakade

et al. (2018), Agarwal et al. (2020), Dann et al. (2019), Zanette & Brunskill (2019), Zhang et al. (2021a), Zhou et al. (2023) and Zhang et al. (2023) worked on model-based algorithms. Notably, Zhang et al. (2023) provided an algorithm that achieves a regret of $\tilde{O}(\min\{\sqrt{SAH^2T}, T\})$, which matches the information lower bound. Jin et al. (2018), Yang et al. (2021), Zhang et al. (2020), Li et al. (2021) and Ménard et al. (2021) work on model-free algorithms. The latter three have introduced algorithms that achieve minimax regret of $\tilde{O}(\sqrt{SAH^2T})$.

Suboptimality Gap. When there is a strictly positive suboptimality gap, it is possible to achieve logarithmic regret bounds. In RL, earlier work obtained asymptotic logarithmic regret bounds Auer & Ortner (2007); Tewari & Bartlett (2008). Recently, non-asymptotic logarithmic regret bounds were obtained (Jaksch et al. (2010); Ok et al. (2018); Simchowitz & Jamieson (2019); He et al. (2021)). Specifically, Jaksch et al. (2010) developed a model-based algorithm, and their bound depends on the policy gap instead of the action gap studied in this paper. Ok et al. (2018) derived problem-specific logarithmic type lower bounds for both structured and unstructured MDPs. Simchowitz & Jamieson (2019) extended the model-based algorithm by Zanette & Brunskill (2019) and obtained logarithmic regret bounds. Logarithmic regret bounds are obtained in linear function approximation settings (He et al., 2021). Nguyen-Tang et al. (2023) also provides a gap-dependent regret bounds for offline RL with linear function approximation.

Specifically, for model free algorithm, Yang et al. (2021) showed that the optimistic Q-learning algorithm by Jin et al. (2018) enjoyed a logarithmic regret $O(\frac{H^6SAT}{\Delta_{\min}})$, which was subsequently refined by Xu et al. (2021). In their work, Xu et al. (2021) introduced the Adaptive Multi-step Bootstrap (AMB) algorithm.

There are also some other works focusing on gap-dependent sample complexity bounds (Jonsson et al., 2020; Marjani & Proutiere, 2020; Al Marjani et al., 2021; Tirinzoni et al., 2022; Wagenmaker et al., 2022b; Wagenmaker & Jamieson, 2022; Wang et al., 2022; Tirinzoni et al., 2023).

Variance reduction in RL. The reference-advantage decomposition used in Zhang et al. (2020) and Li et al. (2021) is a technique of variance reduction that was originally proposed for finite-sum stochastic optimization (see e.g. Gower et al. (2020); Johnson & Zhang (2013); Nguyen et al. (2017)). Later on, model-free RL algorithms also used variance reduction to improve the sample efficiency. For example, it was used in learning with generative models Sidford et al. (2018; 2023); Wainwright (2019), policy evaluation Du et al. (2017); Khamaru et al. (2021); Wai et al. (2019); Xu et al. (2020), offline RL Shi et al. (2022); Yin et al. (2021), and Q-learning Li et al. (2020); Zhang et al. (2020); Li et al. (2021); Yan et al. (2023).

RL with low switching cost. Research in RL with low switching costs aims to minimize the number of policy switches while maintaining comparable regret bounds to fully adaptive counterparts. Bai et al. (2019) first introduced the problem of RL with low-switching cost and proposed a Q-learning algorithm with lazy updates, achieving $\tilde{O}(SAH^3 \log T)$ switching costs. This work was advanced by Zhang et al. (2020), which improved the regret upper bound and the switching cost. Additionally, Wang et al. (2021) studied RL under the adaptivity constraint. Recently, Qiao et al. (2022) proposed a model-based algorithm with $\tilde{O}(\log \log T)$ switching costs.

Other problem-dependent performance. In practice, RL algorithms often perform far more appealingly than what their worst-case performance guarantees would suggest. This motivates a recent line of works that investigate optimal performance in various problem-dependent settings (Fruit et al., 2018; Jin et al., 2020; Talebi & Maillard, 2018; Wagenmaker et al., 2022a; Zhao et al., 2023; Zhou et al., 2023).

F NUMERICAL EXPERIMENTS

In this section, we conduct experiments¹. All the experiments are conducted in a synthetic environment to demonstrate the better gap-dependent regret of UCB-Advantage and Q-EarlySettled-Advantage compared to other two model-free algorithms: UCB-Hoeffding (Jin et al., 2018) and

¹All the experiments are run on a server with Intel Xeon E5-2650v4 (2.2GHz) and 100 cores. Each replication is limited to a single core and 4GB RAM. The total execution time is less than 2 hours. The code for the numerical experiments is included in the supplementary materials along with the submission.

AMB (Xu et al., 2021). We will consider two different scales of experiments across two cases: a general MDP and a deterministic MDP.

We first set $H = 5$, $S = 3$, and $A = 2$. The reward $r_h(s, a)$ for each (s, a, h) is generated independently and uniformly at random from $[0, 1]$. For general MDP, $\mathbb{P}_h(\cdot | s, a)$ is generated on the S -dimensional simplex independently and uniformly at random for (s, a, h) . For deterministic MDP, $\mathbb{P}_h(\cdot | s, a)$ is a randomly generated vector with only one element equal to 1, and all others equal to 0 for each (s, a, h) . Under the given MDP, we generate 3×10^5 episodes. For each episode, we randomly choose the initial state uniformly from the S states. For all four algorithms, we set $\iota = 1$ and the hyper-parameter c_1 in the Hoeffding-type bonus to $\sqrt{2}$. Here, c_1 represents the only undefined constant in the bonus terms of the UCB-Hoeffding and AMB algorithms, as well as the multipliers in the bonus expressions in line 10 of Algorithm 1 (UCB-Advantage) and lines 2 and 4 of Algorithm 2 (Q-EarlySettled-Advantage). In both the UCB-Advantage and Q-EarlySettled-Advantage algorithms, we set the hyper-parameters c_2 to 2. Here, c_2 denotes the constant in the variance estimators of the advantage-type bonus, which is the undefined constant in line 16 of Algorithm 2. In addition, we set c_3 to 1. Here, c_3 denotes the multiplier in the last term in line 9 of Algorithm 1 and the last term in line 8 of Algorithm 2. For UCB-Advantage, we set $N_0 = 200$, and for Q-EarlySettled-Advantage, we set $\beta = 0.05$.

To show error bars, we collect 10 sample paths for all algorithms under the same MDP environment and show the relationship between $\text{Regret}(T)/\log(K + 1)$ and the total number of episodes K in Figure 1. For both panels, the solid line represents the median of the 10 sample paths, while the shaded area shows the 10th and 90th percentiles.

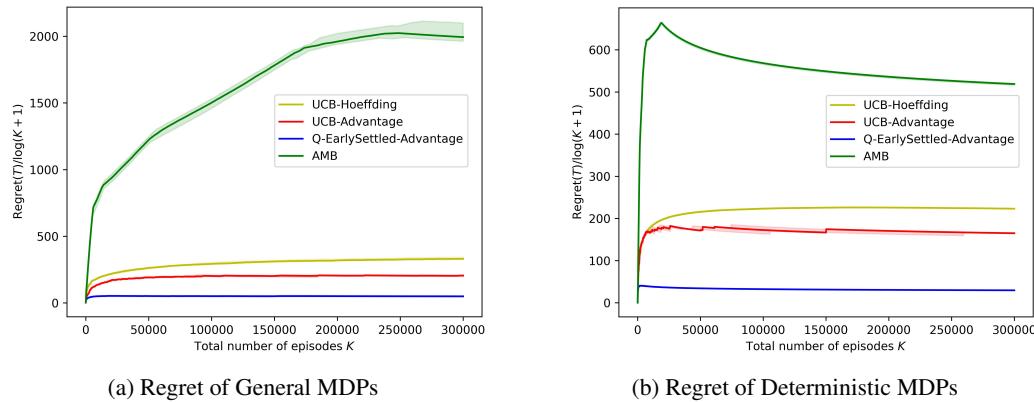


Figure 1: Numerical comparison of regrets with $H = 5$, $S = 3$, and $A = 2$

We also conduct a larger scale experiment with $H = 10$, $S = 5$, and $A = 5$ for 3×10^6 episodes in both types of MDPs. With all other settings unchanged, the result is shown in the following Figure 2:

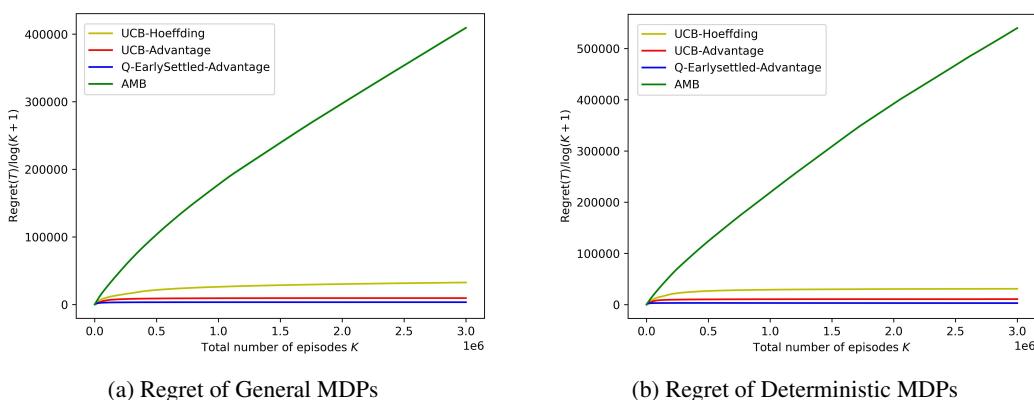


Figure 2: Numerical comparison of regrets with $H = 10$, $S = 5$, and $A = 5$

Next, we discuss the results. From the two figures, we observe that both UCB-Advantage and Q-EarlySettled-Advantage enjoy lower regret compared to UCB-Hoeffding and AMB. The y-axis represents $\text{Regret}(T)/\log(K+1)$, and we note that the curves for UCB-Advantage and Q-EarlySettled-Advantage approach horizontal lines as K becomes sufficiently large. This suggests that the regret for these two algorithms grows logarithmically with K . In particular, Q-EarlySettled-Advantage achieves even lower regret than UCB-Advantage when K is large. These features are consistent with our theoretical results.

We also conduct an experiment to evaluate the policy switching cost of the UCB-Advantage algorithm for these two different scales of (H, S, A) , under the same experimental setting. The results are presented in the following figures:

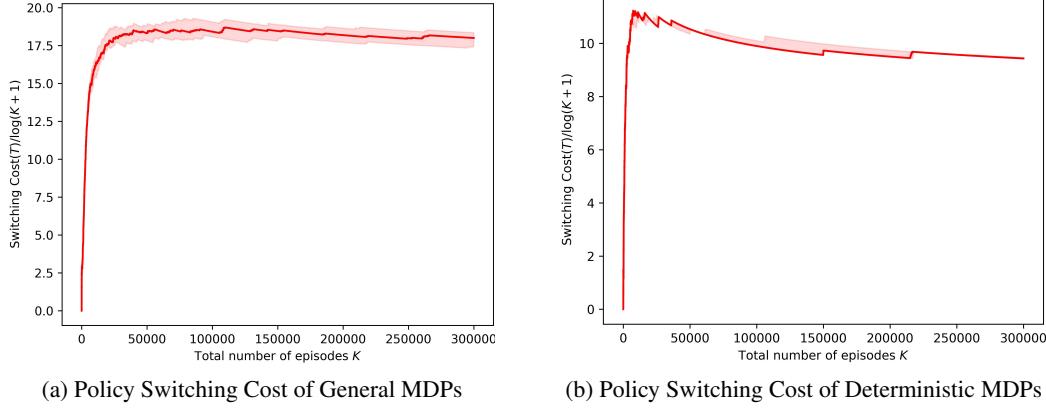


Figure 3: Policy switching cost of UCB-Advantage algorithm with $H = 5$, $S = 3$, and $A = 2$

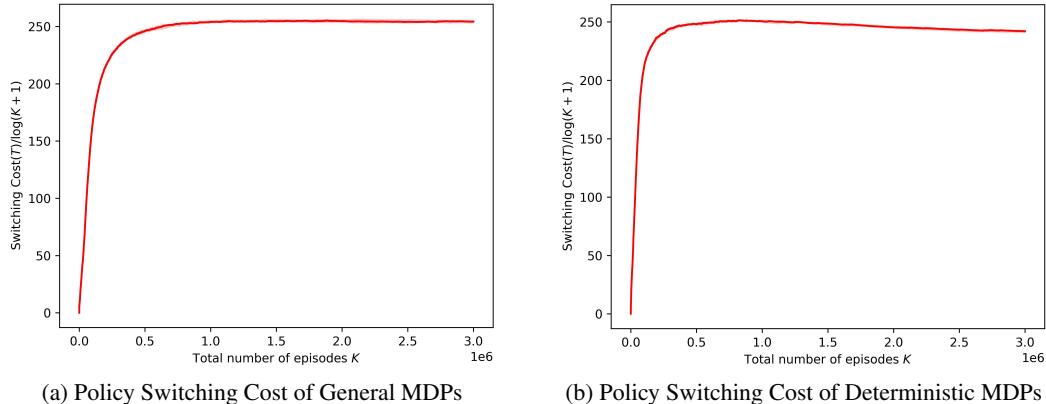


Figure 4: Policy switching cost of UCB-Advantage algorithm with $H = 10$, $S = 5$, and $A = 5$

In these two figures, the y-axis represents the ratio of policy switching cost to $\log(K + 1)$. We note that all these four curves approach horizontal lines as K becomes sufficiently large, which is consistent with our logarithmic policy switching cost shown in Equation (4).

G MATHEMATICAL EXPLANATION OF THE SURROGATE FUNCTION

Next, we explain the surrogate function in a more mathematical manner.

Our proof relies on relating the regret to multiple groups of estimation error sums that take the form $\sum_{k=1}^K \omega_{h,k}^{(i)} (Q_h^k - Q_h^*)(s_h^k, a_h^k)$. Here $\{\omega_{h,k}^{(i)}\}_k$ are nonnegative weights and i represents the group. Bounding the weighted sum via controlling each individual $Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)$ by

recursion on h is a common technique for model-free optimism-based algorithms, which was used by all of Zhang et al. (2020); Li et al. (2021); Yang et al. (2021). Yang et al. (2021) used it on gap-dependent regret analysis, and Zhang et al. (2020) and Li et al. (2021) used it to control the reference setting errors $\sum_{k=1}^K (V_h^{R,k+1}(s_h^k) - V_h^{R,K+1}(s_h^k))$. However, their techniques are only limited to the Hoeffding-type update. In detail, the Hoeffding-type update in Q -function is given by

$$Q_h^{k+1}(s_h^k, a_h^k) = r_h(s_h^k, a_h^k) + \sum_{n=1}^{N_h^{k+1}} \eta_n^{N_h^{k+1}} V_{h+1}^{k^n}(s_{h+1}^{k^n}) + \tilde{O}\left(\sqrt{H^3/N_h^{k+1}}\right),$$

which is the key update of Yang et al. (2021), and the update of $Q_h^{\text{UCB},k+1}$ for [2, 3]. Accordingly, we can find that

$$(Q_h^k - Q_h^*)(s_h^k, a_h^k) \leq H\eta_0^{N_h^k} + \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (V_{h+1}^{k^n} - V_{h+1}^*)(s_{h+1}^{k^n}) + \tilde{O}\left(\sqrt{H^3/N_h^k}\right),$$

which is the event in Definition 4.1 of Yang et al. (2021). Here, $\eta_0^{N_h^k} = 0$ when $N_h^k > 0$. After taking the weighted sum with regard to $k \in [K]$ on both sides, we can establish recursions on h where the main terms are $\sum_{k=1}^K \omega_{h,k}^{(i)} (Q_h^k - Q_h^*)(s_h^k, a_h^k)$ and $\sum_{k=1}^K \omega_{h,k}^{(i)} \sum_{n=1}^{N_h^{k+1}} \eta_n^{N_h^{k+1}} (V_{h+1}^{k^n} - V_{h+1}^*)(s_{h+1}^{k^n})$. With $\sum_{k=1}^K H\eta_0^{N_h^k}$ being easily controlled, the error generated by the recursion is mainly dominated by the weighted sum regarding the simple term $\tilde{O}\left(\sqrt{H^3/N_h^{k+1}}\right)$, which obviously vanishes when k is large so that N_h^k (the number of visit to (s_h^k, a_h^k, h)) is large.

Here, we explain why Zhang et al. (2020) and Li et al. (2021) only rely on the weighted sum $\sum_{k=1}^K \omega_{h,k}^{(i)} (Q_h^k - Q_h^*)(s_h^k, a_h^k)$ with simple Hoeffding-type errors though their algorithms involve reference-advantage decomposition. Both methods incorporate a Hoeffding-type update (see $Q_h^{\text{UCB},k+1}$ in Equation (7)), with which they bound the reference settling error by controlling the weighted sum. When analyzing the worst-case regret, they only need to relate the regret to $\sum_{k=1}^K (Q_h^k - Q_h^*)(s_h^k, a_h^k)$, i.e., the sum instead of the weighted sum. However, in our gap-dependent regret analysis, because the weights do not adapt to the learning process (see our proof sketch for more details), we have to analyze each item $(Q_h^k - Q_h^*)(s_h^k, a_h^k)$ individually in the weighted sum with complicated errors with new technical tools when we consider the reference-advantage update (Equation (8)).

The reference-advantage update is listed as follows

$$Q_h^{R,k+1}(s_h^k, a_h^k) = r_h^k(s_h^k, a_h^k) + \sum_{n=1}^{N_h^{k+1}} \left(\eta_n^{N_h^{k+1}} (V_{h+1}^{k^n} - V_{h+1}^{R,k^n}) + u_n^{N_h^{k+1}} V_{h+1}^{R,k^n} \right) (s_{h+1}^{k^n}) + \tilde{R}^{h,k+1}.$$

Here, $\{\eta_n^{N_h^{k+1}}\}_{n=1}^{N_h^{k+1}}$ are the corresponding nonnegative weights that sum to 1. $\{u_n^{N_h^{k+1}}\}_{n=1}^{N_h^{k+1}}$ that sum to 1 are nonnegative weights for the reference function. $\tilde{R}^{h,k+1}$ is the cumulative bonus that contains variance estimators and dominates the variances in reference estimations and advantage estimations. Accordingly, we can find that

$$\begin{aligned} (Q_h^k - Q_h^*)(s_h^k, a_h^k) &\leq H\eta_0^{N_h^k} + \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (V_{h+1}^{k^n} - V_{h+1}^*)(s_{h+1}^{k^n}) \\ &\quad + \sum_{n=1}^{N_h^k} \left(\eta_n^{N_h^k} (V_{h+1}^* - V_{h+1}^{R,k^n}) + u_n^{N_h^k} V_{h+1}^{R,k^n} \right) (s_{h+1}^{k^n}) - (1 - \eta_0^{N_h^k}) \mathbb{P}_{(s_h^k, a_h^k, h)} V_{h+1}^* + R^{h,k}. \end{aligned}$$

To establish the recursion on h in the same way, when keeping the main terms unchanged and neglecting the term $H\eta_0^{N_h^k}$, the error term in our iteration becomes the weighted summation for

$$\sum_{n=1}^{N_h^k} \left(\eta_n^{N_h^k} (V_{h+1}^* - V_{h+1}^{R,k^n}) + u_n^{N_h^k} V_{h+1}^{R,k^n} \right) (s_{h+1}^{k^n}) - (1 - \eta_0^{N_h^k}) \mathbb{P}_{(s_h^k, a_h^k, h)} V_{h+1}^* + R^{h,k}.$$

2862 It is much more complicated than $\tilde{O}(\sqrt{H^3/N_h^k})$ for the Hoeffding-type update.
 2863
 2864 To handle this error, we propose a decomposition method following the reference advantage struc-
 2865 ture. Naively, we can move towards advantage estimation errors (the first term), reference estimation
 2866 errors (the second term), reference settling errors (the third term), the cumulative bonus (the fourth
 2867 term), and a negative term (the last term), i.e.

$$\begin{aligned} & \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(\mathbb{P}_{s_h^k, a_h^k, h} - \mathbb{1}_{s_{h+1}^{k^n}} \right) (V_{h+1}^{R, K+1} - V_{h+1}^*) + \sum_{n=1}^{N_h^k} u_n^{N_h^k} \left(\mathbb{1}_{s_{h+1}^{k^n}} - \mathbb{P}_{s_h^k, a_h^k, h} \right) V_{h+1}^{R, K+1}(s_{h+1}^{k^n}) \\ & + \sum_{n=1}^{N_h^k} u_n^{N_h^k} (V_{h+1}^{R, k^n} - V_{h+1}^{R, K+1})(s_{h+1}^{k^n}) + R^{h, k} + \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (V_{h+1}^{R, K+1} - V_{h+1}^{R, k^n})(s_{h+1}^{k^n}) \end{aligned}$$

2875 because the properties of the settled reference function $V_{h+1}^{R, K+1}$ is well-studied in Zhang et al.
 2876 and Li et al. (2021). However, it will cause a non-martingale issue when we try to ap-
 2877 ply concentration inequalities as $V_{h+1}^{R, K+1}$ depends on the whole learning process. To solve this
 2878 issue, we propose our **surrogate reference function** $\hat{V}_h^{R, k}$ and decompose the error above as
 2879
 $\mathcal{G}_1 := \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (\mathbb{P}_{s_h^k, a_h^k, h} - \mathbb{1}_{s_{h+1}^{k^n}}) (\hat{V}_{h+1}^{R, k^n} - V_{h+1}^*), \mathcal{G}_2 := \sum_{n=1}^{N_h^k} u_n^{N_h^k} (\mathbb{1}_{s_{h+1}^{k^n}} - \mathbb{P}_{s_h^k, a_h^k, h}) \hat{V}_{h+1}^{R, k^n},$
 $\mathcal{G}_3 := \sum_{n=1}^{N_h^k} (u_n^{N_h^k} - \eta_n^{N_h^k}) \mathbb{P}_{s_h^k, a_h^k, h} \hat{V}_{h+1}^{R, k^n} + \sum_{n=1}^{N_h^k} u_n^{N_h^k} (V_{h+1}^{R, k^n} - \hat{V}_{h+1}^{R, k^n})(s_{h+1}^{k^n}),$ the bonus term
 $\mathcal{G}_4 = R^{h, k}$, and a negative negligible term $\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (\hat{V}_{h+1}^{R, k^n} - V_{h+1}^{R, k^n})(s_{h+1}^{k^n})$. The first three terms
 2883 correspond to advantage estimation error, reference estimation error, and reference settling error, re-
 2884 spectively. Here, we creatively use the surrogate $\hat{V}_{h+1}^{R, k}$ as it is determined before the start of episode
 2885 k . Thus, $\mathcal{G}_1, \mathcal{G}_2$ are martingale sums and can be controlled by concentration inequalities that are
 2886 given in Equation (16), so the non-martingale challenge can be addressed. \mathcal{G}_3 corresponds to the
 2887 reference settling error and can also be controlled given the settling conditions and properties of
 2888 $\hat{V}_h^{R, k}(s)$. The bonus \mathcal{G}_4 is controlled using the same idea of bounding $\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3$.

2890 Our decomposition above expands the technique of bounding the weighted sum of estimation errors
 2891 to reference-advantage type estimations. In addition, we are the first to use the novel construction
 2892 of the reference surrogates for reference-advantage decomposition in the literature, which makes a
 2893 separate contribution to future work on off-policy methods and offline methods.

2894
 2895
 2896
 2897
 2898
 2899
 2900
 2901
 2902
 2903
 2904
 2905
 2906
 2907
 2908
 2909
 2910
 2911
 2912
 2913
 2914
 2915