# Reason for Future, Act for Now: A Principled Framework for Autonomous LLM Agents with Provable Sample Efficiency

Zhihan Liu\*† Hao Hu\*‡ Shenao Zhang\*†
Hongyi Guo† Shuqi Ke § Boyi Liu† Zhaoran Wang†
July 2, 2024

#### Abstract

Large language models (LLMs) demonstrate impressive reasoning abilities, but translating reasoning into actions in the real world remains challenging. In particular, it is unclear how to complete a given task provably within a minimum number of interactions with the external environment, e.g., through an internal mechanism of reasoning. To this end, we propose the first framework with provable regret guarantees to orchestrate reasoning and acting, which we call "reason for future, act for now" (RAFA). Specifically, we design a prompt template for reasoning that learns from the memory buffer and plans a future trajectory over a long horizon ("reason for future"). At each step, the LLM agent takes the initial action of the planned trajectory ("act for now"), stores the collected feedback in the memory buffer, and reinvokes the reasoning routine to replan the future trajectory from the new state. The key idea is to cast reasoning in LLMs as learning and planning in Bayesian adaptive Markov decision processes (MDPs). Correspondingly, we prompt LLMs with the memory buffer to estimate the unknown environment (learning) and generate an optimal trajectory for multiple future steps that maximize a value function (planning). The learning and planning subroutines are performed in an "incontext" manner to emulate the actor-critic update for MDPs. Our theoretical analysis establishes a  $\sqrt{T}$  regret, while our experimental validation demonstrates superior empirical performance. Here, T denotes the number of online interactions. Project page: https://agentification.github.io/RAFA.

<sup>\*</sup>Equal contribution.

 $<sup>^\</sup>dagger Northwestern~University.~\{ \verb|zhihanliu| 2027, \verb|shenaoz| hang| 2028, \verb|hongyiguo| 2025, \verb|boyiliu| 2018 \} \\ @u.northwestern.edu, \verb|zhaoran| wang| @gmail.com|$ 

<sup>&</sup>lt;sup>‡</sup>Tsinghua University. huh22@mails.tsinghua.edu.cn

<sup>§</sup>The Chinese University of Hong Kong. shuqike@link.cuhk.edu.cn

# Contents

1	Introduction	4
	1.1 Literature	5 8
2	Bridging LLM and RL	8
3	Algorithm	11
4	Experiment         4.1 Game of 24          4.2 ALFWorld          4.3 BlocksWorld          4.4 Tic-Tac-Toe	15 15 16 17 18
5	Theoretical Analysis  5.1 Planning Optimality  5.2 LLMs with Posterior Alignments Perform Bayesian Model Averaging (BMA)  5.3 Regret Bound of RAFA  5.4 RAFA with Efficient Exploration Strategies  5.4.1 Optimistic Bonus  5.4.2 Posterior Sampling	19 19 20 21 23 24 25
6	Conclusions	27
A	Notations	35
В	More Algorithms	36
C	Main Proofs C.1 Proof of Proposition 5.2 C.2 LLMs with Posterior Alignments Perform BMA C.3 Contraction Property of the Posterior Variance C.4 Proof of Theorem 5.7 C.5 Proof of Theorem 5.8 C.6 Proof of Theorem 5.10 C.7 Relaxing Assumption 5.3 for Theorem 5.7	38 38 38 39 41 46 50 53
D	Missing Proofs in Appendix C  D.1 Proof of Lemma C.2	<b>57</b> 57
${f E}$	Linear Special Case	60

${f F}$	More Experiments	63
	F.1 Game of 24	64
	F.2 ALFWorld	65
	F.3 BlocksWorld	67
	F.4 Tic-Tac-Toe	68
$\mathbf{G}$	Prompts	69
	G.1 Game of 24	70
	G.2 ALFWorld	73
	G.3 Blocksworld	79
	G.4 Tic-Tac-Toe	86

## 1 Introduction

Large language models (LLMs) exhibit remarkable reasoning abilities, which open a new avenue for agents to interact with the real world autonomously. However, turning reasoning into actions remains challenging. Specifically, although LLMs are equipped with the prior knowledge obtained through pretraining, it is stateless in nature and ungrounded in the real world, which makes the resulting action suboptimal. To bridge the reasoning-acting gap, we aim to design an internal mechanism of reasoning on top of LLMs, which optimizes actions iteratively by incorporating feedback from the external environment. In particular, we focus on the sample efficiency of autonomous LLM agents in interactive decision-making tasks, which plays a key role in their practical adoption, especially when interactions are costly and risky. Our primary goal is to enable agents to complete a given task in a guaranteed manner through reasoning within a minimum number of interactions with the external environment.

Reinforcement learning (RL) is a well-studied paradigm for improving actions by collecting feedback. However, to tailor existing RL techniques for autonomous LLM agents, we lack a rigorous mapping between RL and LLMs, which leads to various conceptual discrepancies. For example, RL operates in a numerical system, where rewards and transitions are defined by scalars and probabilities. In comparison, the inputs and outputs of LLMs are described by tokens in a linguistic system. As another example, LLMs are trained on a general-purpose corpus and remain fixed throughout the interactive process. In contrast, RL trains actors and critics via parameter updates on the collected feedback iteratively. Thus, it appears inappropriate to treat LLMs as actors or critics under the RL framework, although all of them are parameterized by deep neural networks. Moreover, it remains unclear what reasoning with LLMs means under the RL framework, e.g., what are the inputs and outputs of a reasoning routine and how reasoning should be coordinated with acting. Such conceptual discrepancies prevent us from establishing a principled framework beyond borrowing the "trial and error" concept from RL straightforwardly and make it difficult to establish the theoretical guarantee.

To address such conceptual discrepancies, we formalize reasoning and acting with LLMs under a Bayesian adaptive Markov decision process (MDP) framework, where the latent variable of interest is the unknown environment. The starting point is to cast the full history of states (of the external environment), actions, rewards, and their linguistic summaries in the memory buffer as the information state of Bayesian adaptive MDPs. Throughout the interactive process, the information state accumulates a growing collection of feedback from the external environment, which is mapped to an optimized action at each step by an internal mechanism of reasoning. As detailed below, we construct the reasoning routine through two key subroutines, namely learning and planning, which are instantiated by LLMs with specially designed prompts. (a) The learning subroutine forms an estimate of the external environment given the memory buffer, where LLMs are prompted to infer the transition and reward models (model) or/and the value function (critic). (b) The planning subroutine generates an optimal policy (actor) or trajectory for multiple future steps, which maximizes the value function (up to a certain error). Depending on the specific configuration of the

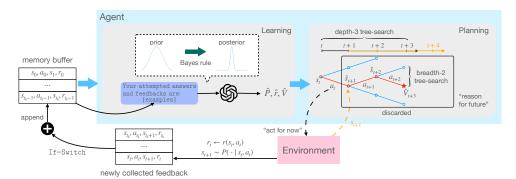


Figure 1: Illustration of the RAFA ("reason for future, act for now") framework.

state and action spaces (continuous versus discrete) and the transition and reward models (stochastic versus deterministic), the planning subroutine emulates the value iteration algorithm, the random shooting algorithm, or the Monte-Carlo tree-search algorithm.

Although LLMs remain fixed throughout the interactive process, we can reduce their estimation uncertainty by prompting the growing collection of feedback from the external environment as contexts, which is verified both theoretically and empirically in this paper. From the perspective of Bayesian adaptive MDPs, LLMs can be considered as some functional of the posterior of the environment (for example, Bayesian model averaging (Wasserman, 2000)), hence the estimation uncertainty is reduced with increasing information via interactions. For several tasks, we demonstrate that LLMs can make a more precise prediction when prompted with more data as contexts. Hence, LLMs can play a similar role of model estimators in the design of online RL algorithms for interactions. We improve the accuracy of LLMs by simply adding the new feedback to the memory buffer as contexts, instead of performing explicit parameter updates (such as gradient descent) on deep neural networks as in existing RL methods.

We conclude our contributions in this paper from three perspectives. (a) We establish the LLM-RL correspondence and design a principled framework RAFA for orchestrating the reasoning and acting of LLMs. (b) Our empirical validation shows that RAFA outperforms various existing frameworks in interactive decision-making tasks, including ALFWorld, BlocksWorld, Game of 24, and a new benchmark based on Tic-Tac-Toe. (c) Our theoretical analysis proves that RAFA achieves a  $\sqrt{T}$  regret, explaining why RAFA demonstrates strong empirical performance. Here, T denotes the number of online interactions. We also provide two provably efficient variants of RAFA to implement efficient exploration for more complex tasks.

### 1.1 Literature

Reasoning with LLM. We build on a recent line of work that develops various prompting schemes to improve the reasoning performance of LLMs. "Chain of thoughts" ("CoT") (Wei et al., 2022) decomposes a challenging problem into several reasoning stages and guides LLMs to solve them one by one. As generalizations, "tree of thoughts" (Yao et al., 2023a), "graph of thoughts" (Yao et al., 2023b), "algorithm of thoughts" (Sel et al., 2023), and "cumulative reasoning" (Zhang et al., 2023a) provide different graph-search schemes to guide LLMs. See

also Wang et al. (2022a); Creswell et al. (2022); Creswell and Shanahan (2022); Guo et al. (2024); Zhang et al. (2024). Also, "reasoning via planning" ("RAP") (Hao et al., 2023) emulates the Monte-Carlo tree-search (MCTS) algorithm to reduce the search complexity. Pouplin et al. (2024) improve LLM reasoning process with MCTS and formulate the reasoning process as an MDP. Sun et al. (2023a) use offline inverse RL to optimize the prompts for arithmetic problems. For embodied LLM agents, Huang et al. (2022a) propose to decompose a complex task into multiple executable steps. Most of them focus on general reasoning tasks, e.g., solving a mathematical or logic puzzle, where LLMs generate a detailed trace (trajectory) of arguments through an internal mechanism to reach a final answer. Here, LLMs play the same role as the planning subroutine in RAFA. In contrast, we focus on interactive decision-making tasks, where autonomous LLM agents collect feedback from the external environment to optimize actions iteratively. In particular, we aim to complete a given task within a minimum number of interactions with the external environment. To this end, it is essential to operate three interleaved modules, namely learning, planning, and acting, in a closed loop. While it is feasible to incorporate existing graph-search or MCTS schemes as the planning subroutine for generating trajectories, our core contribution is a principled framework that executes a selected subset of the planned trajectory to collect feedback ("act for now") and replans an improved trajectory from the new state by learning from feedback ("reason for future"). From an RL perspective, existing graph-search or MCTS schemes are analogous to an open-loop method, e.g., motion planning or trajectory optimization (Betts, 1998), which does not involve interactions with the external environment. To integrate them into a closed-loop approach, e.g., model predictive control (Rawlings, 2000), one has to specify how to act given the planned trajectory and when to reinvoke the reasoning (learning and planning) routine, which is the key technique of RAFA. Another recent line of work tackles more complex tasks by allowing LLMs to access various additional modules, e.g., tools, programs, and other learning algorithms (Ahn et al., 2022; Shen et al., 2023; Lu et al., 2023; Liu et al., 2023a; Cai et al., 2023), or by finetuning LLMs on the feedback (Zelikman et al., 2022; Li et al., 2022; Paul et al., 2023; Sun, 2023).

Acting (and Reasoning) with LLM. We build on a recent line of work that develops various closed-loop frameworks for interacting with the external environment. "Inner monologue" (Huang et al., 2022b) and "ReAct" (Yao et al., 2022) combine reasoning and acting to refine each other for the first time. In comparison, RAFA provides a specific schedule for orchestrating reasoning and acting (as discussed above). As generalizations, "Reflexion" (Shinn et al., 2023) enables autonomous LLM agents to revise the current action of a pregenerated trajectory by learning from feedback, especially when they make mistakes. See also Kim et al. (2023). However, making a local revision to the pre-generated trajectory is myopic because it fails to consider the long-term consequences of actions. Consequently, the obtained policy may get trapped by a local optimum. From an RL perspective, "Reflexion" (Shinn et al., 2023) is an oversimplified version of RAFA, where the planning subroutine revises the current action to maximize the reward function ("reason for now") instead of planning mul-

Closed-Loop Mechanisms	No Parameter Update	Theoretical Guarantee
RAFA	<b>✓</b>	✓
Model-Based Deep RL	X	✓
Model Predictive Control	X	✓
Thompson Sampling	×	✓
"React", "Reflexion", and "Adaplanner"	✓	Х

Table 1: Comparison between RAFA and other mechanisms.

tiple future steps to maximize the value function ("reason for future"), which measures the expected cumulative future reward. To remedy this issue, "AdaPlanner" (Sun et al., 2023b) regenerates the whole trajectory at each step, which yields a global improvement. See also Wang et al. (2023b). However, the reasoning routine of "AdaPlanner" requires a handcrafted set of programs to reject suboptimal candidate trajectories. Without the domain knowledge of the current task, the regenerated trajectory is not necessarily optimal, i.e., maximizing the value function (up to a certain error). In contrast, the reasoning routine of RAFA is designed following the principled approach in RL. In particular, the learning subroutine infers the transition and reward models (model) or/and the value function (critic), while the planning subroutine emulates the value iteration algorithm, the random shooting algorithm, or the MCTS algorithm, none of which use any domain knowledge. RAFA also achieves provable sample efficiency guarantees for the first time and outperforms those existing frameworks empirically.

Large Language Model (LLM) and In-Context Learning (ICL). LLMs (Radford et al., 2019; Brown et al., 2020; Hoffmann et al., 2022; Chowdhery et al., 2022; OpenAI, 2023; Touvron et al., 2023) display notable reasoning abilities. A pivotal aspect of reasoning is the ICL ability (Liang et al., 2022; Razeghi et al., 2022; Shin et al., 2022; Olsson et al., 2022; Akyürek et al., 2022; Kirsch et al., 2022; Garg et al., 2022; Von Oswald et al., 2023; Li et al., 2023; Abernethy et al., 2023), which allows LLMs to solve a broad range of tasks with only a few in-context examples instead of finetuning parameters on a specific dataset. We focus on harnessing the ICL ability of LLMs to optimize actions in the real world, which is crucial to autonomous LLM agents. In particular, we build on a recent line of work (Xie et al., 2021; Zhang et al., 2022, 2023b; Wang et al., 2023a; Wies et al., 2023; Jiang, 2023; Lee et al., 2023) that attributes the ICL ability to implicit Bayesian inference, i.e., an implicit mechanism that enables LLMs to infer a latent concept from those in-context examples, which is verified both theoretically and empirically. In RAFA, the latent concept is the transition and reward models (model) of the unknown environment or/and the value function (critic), which is inferred from the memory buffer in the learning subroutine. Claim 2.1 can also be considered as a result of ICL ability.

Reinforcement Learning (RL) under a Bayesian Framework. We build on a recent line of work on the infinite-horizon (Abbasi-Yadkori and Szepesvári, 2015; Dong et al., 2019;

Wei et al., 2020; Zhou et al., 2021a,b; Chen et al., 2022; Chua et al., 2018; Hafner et al., 2019; Sekar et al., 2020) and Bayesian (Strens, 2000; Osband et al., 2013; Russo and Van Roy, 2014a,b, 2016; Lu and Van Roy, 2019) settings of RL, which include model-based deep RL (Janner et al., 2019; Liu et al., 2023b; Wang et al., 2022b; Liu et al., 2024), model predictive control (Morari and Lee, 1999), and Thompson sampling (Russo and Van Roy, 2014b). The infinite-horizon setting allows RAFA to interact with the external environment continuously without resetting to an initial state, while the Bayesian setting allows us to connect RAFA with BMA and establish the theoretical guarantee. RL operates in a numerical system, where rewards and transitions are defined by scalars and probabilities, and trains actors and critics on the collected feedback iteratively. We focus on emulating the actor-model or actor-critic update in RL through an internal mechanism of reasoning on top of LLMs, which allows data and actions to be tokens in a linguistic system while bypassing the explicit update of parameters in model-based RL (Chua et al., 2018; Hafner et al., 2019; Sekar et al., 2020; Liu et al., 2022b; Zhong et al., 2022; Zheng et al., 2022; Liu et al., 2022a). In particular, the learning and planning subroutines of RAFA emulate the posterior update and various planning algorithms in RL. Moreover, RAFA orchestrates reasoning (learning and planning) and acting following the principled approach in RL, i.e., (re)planning a future trajectory over a long horizon ("reason for future") at the new state and taking the initial action of the planned trajectory ("act for now"). As a result, RAFA inherits provable sample efficiency guarantees from RL. We summarize the comparison between RAFA and other closed-loop mechanisms in Table 1.

#### 1.2 Notations

We provide a table of notations in Appendix A.

# 2 Bridging LLM and RL

Interaction Protocol. We use Bayesian adaptive Markov decision processes (MDPs) (Ghavamzadeh et al., 2015) to model how autonomous LLM agents interact with the external environment. We consider an infinite-horizon MDP  $M = (S, A, P, r, \rho, \gamma, \mathbb{P}_0)$ , where S is the state space, A is the action space,  $P: S \times A \mapsto \Delta(S)$  is the transition kernel,  $r: S \times A \mapsto \mathbb{R}$  is the reward function,  $\rho$  is the initial distribution of states,  $\gamma \in (0,1)$  is the discount factor, and  $\mathbb{P}_0$  is the prior distribution of the transition kernel and the reward function. Here, P gives the probability distribution of the next state given the current state and action, while r is assumed to be deterministic without loss of generality. For notational simplicity, we parameterize P and r by a shared parameter  $\theta \in \Theta$  and denote them as  $P_{\theta}$  and  $r_{\theta}$ . At the beginning of the interaction, the data-generating parameter  $\theta^*$  is sampled from the prior  $\mathbb{P}_0$ . At the t-th step during the interaction, the LLM agent receives a state  $s_t \in S$ , takes an action  $a_t \in A$  following the current policy  $\pi_t: S \mapsto A$ , and receives a reward  $r_t = r_{\theta^*}(s_t, a_t)$ . Subsequently, the external environment transits to the next state

 $s_{t+1} \sim P_{\theta^*}(\cdot \mid s_t, a_t)$ , while the LLM agent computes the updated policy  $\pi_{t+1}$  through an internal mechanism of reasoning (as discussed below). Note that  $\mathcal{S}$  and  $\mathcal{A}$  are represented by tokens in a linguistic system. Here,  $\pi \in \Pi$  is assumed to be deterministic without loss of generality, where  $\Pi$  is the feasible set of policies.

**Value Function.** For a policy  $\pi$  and a parameter  $\theta$  of the transition and reward models, we define the state-value and action-value functions as

$$V_{\theta}^{\pi}(s) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} r_{\theta}(s_{t}, a_{t}) \middle| s_{0} = s\right],$$

$$Q_{\theta}^{\pi}(s, a) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} r_{\theta}(s_{t}, a_{t}) \middle| s_{0} = s, a_{0} = a\right],$$

$$(2.1)$$

where  $\mathbb{E}$  is taken with respect to  $a_t = \pi(s_t)$  and  $s_{t+1} \sim P_{\theta}(\cdot \mid s_t, a_t)$  for all  $t \geq 0$ . In other words,  $V_{\theta}^{\pi}$  (and  $Q_{\theta}^{\pi}$ ) gives the expected cumulative future reward from the current state s (and action a).

We define the optimal policy  $\pi_{\theta}^{\star}$  with respect to a given parameter  $\theta$  as  $\pi_{\theta}^{\star} = \operatorname{argmax}_{a \in \mathcal{A}} Q_{\theta}^{\star}$ , where  $Q_{\theta}^{\star}$  is the fixed point of the following Bellman optimality equation,

$$Q_{\theta}^{\star}(s, a) = (B_{\theta}V_{\theta}^{\star})(s, a),$$
  

$$V_{\theta}^{\star}(s) = \max_{a \in \mathcal{A}} Q_{\theta}^{\star}(s, a),$$
(2.2)

where  $Q_{\theta}^{\star}$  and  $V_{\theta}^{\star}$  are the fixed-point solutions. Here, we define  $(B_{\theta}V)(s,a) = r_{\theta}(s,a) + \gamma \cdot (P_{\theta^{\star}}V)(s,a)$  and  $(P_{\theta}V)(s,a) = \mathbb{E}_{s' \sim P_{\theta}(\cdot \mid s,a)}[V(s')]$  for any value function V. See Sutton and Barto (2018) for the existence and uniqueness guarantees for  $Q_{\theta}^{\star}$ ,  $V_{\theta}^{\star}$ , and  $\pi_{\theta}^{\star}$ .

Posterior, Entropy, and Information Gain. By Bayes' rule, the posterior of  $\theta^*$  given any in-context dataset  $\mathcal{D}$  is

$$\mathbb{P}_{\text{post}}(\theta \mid \mathcal{D}) \propto \mathbb{P}_{0}(\theta) L(\mathcal{D} \mid \theta), \tag{2.3}$$

where we denote by  $L(\mathcal{D} | \theta)$  the likelihood of  $\mathcal{D}$  conditioned on  $\theta$ . We define the random variable  $\xi_{(s,a)}$  as the pair of the next state and the current reward (s',r) given the query state-action pair (s,a). Given any in-context dataset  $\mathcal{D}$  and query state-action pair (s,a), the posterior of  $\xi_{(s,a)}$  can be specified as

$$\mathbb{P}_{\text{post}}(\xi_{(s,a)} \mid \mathcal{D}, s, a) = \mathbb{E}_{\theta \sim \mathbb{P}_{\text{post}}(\cdot \mid \mathcal{D})} [P_{\theta}(s' \mid s, a) \cdot \mathbf{1}(r = r_{\theta}(s, a))], \tag{2.4}$$

where we use Bayes' rule and the fact that the query state-action pair (s, a) is conditionally independent of  $\theta^*$  given  $\mathcal{D}$ . To characterize the uncertainty of  $\theta^*$  conditioned on  $\mathcal{D}$ , we define the posterior entropy  $H(\theta \mid \mathcal{D})$  as

$$H(\theta \mid \mathcal{D}) = \mathbb{E}_{\theta \sim \mathbb{P}_{post}(\cdot \mid \mathcal{D})} \left[ -\log \left( p_{post}(\theta \mid \mathcal{D}) \right) \right], \tag{2.5}$$

where  $p_{\text{post}}$  is the probability mass (or density) function of  $\mathbb{P}_{\text{post}}$ . High posterior entropy  $H(\theta \mid \mathcal{D})$  means high uncertainty of  $\theta^*$ , which suggests that it is hard for the agent to make a precise prediction given  $\mathcal{D}$ . We also define the information gain  $I(\theta; \xi \mid \mathcal{D})$  as  $H(\theta \mid \mathcal{D}) - H(\theta \mid \mathcal{D}, \xi)$ , which characterizes how much information  $\xi_{(s,a)}$  carries to reduce the uncertainty of  $\theta^*$  conditioned on  $\mathcal{D}$ .

**Sample Efficiency.** As the performance metric, we define the Bayesian regret after T steps of interactions,

$$\mathfrak{R}(T) = \mathbb{E}\left[\sum_{t=0}^{T-1} V_{\theta^{\star}}^{\pi^{\star}}(s_t) - V_{\theta^{\star}}^{\pi_t}(s_t)\right],\tag{2.6}$$

where  $\pi^* = \pi_{\theta^*}^*$ ,  $\mathbb{E}$  is taken with respect to the prior distribution  $\mathbb{P}_0$  of  $\theta^*$ , the stochastic outcome of  $s_t$ , and the iterative update of  $\pi_t$ , which involves states, actions, and rewards until the t-th step, i.e., the full history  $\mathcal{D}_t = \{(s_i, a_i, s_{i+1}, r_i)\}_{i=0}^{t-1}$ . We aim to design a sample-efficient agent that satisfies  $\mathfrak{R}(T) = o(T)$ , i.e., the Bayesian regret is sublinear in the total number of interactions T.

What Reasoning Means and Role of LLM. To bridge LLM mechanisms with online RL algorithms, we claim that LLMs can play a similar role of model estimators in the design of online RL algorithms for interactions, which is one aspect of In-Context Learning (ICL) ability of LLMs.

Claim 2.1. LLMs provide a more accurate estimate for the environment with more feedback from online interactions.

In Proposition 5.4 in Section 5, we prove that LLMs with posterior alignment perform Bayesian model averaging (BMA). This theoretical result supports Claim 2.1, as the estimation uncertainty of BMA is reduced given more feedback from interactions with the environment (Wasserman, 2000). We also provide empirical evidence on three tasks for Claim 2.1 as follows.

- (a) Information Bandit. The goal of our 100-arm bandit experiment is to find the arm with the highest reward. There is an informative arm whose reward is the index of the best arm. We prompt the LLM (gpt-4) to pull the arm by providing it with the historical data of several bandit instances that share the same informative arm. It can be observed from the right figure that the LLM can learn the best arm with only 6 examples, and is thus an effective reward estimator.
- (b) Concept Learning. We evaluate LLMs (Llama 2-7B) in three tasks (Todd et al., 2023) with hidden concepts: (1) Antonym: Generate the word with the opposite meaning given an input word; (2) Country-Capital: Generate the capital city of a given country; and (3) Present-Past: Generate the verb's past inflection given a verb in the present tense. We observe that with more in-context examples provided to the LLM, the accuracy of the test instance monotonically increases, indicating that the hidden concepts of the tasks are learned.

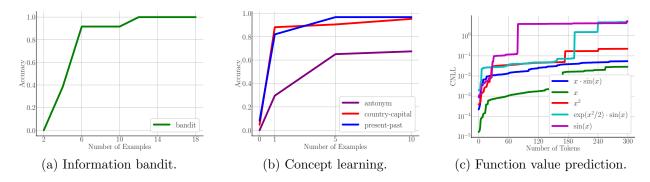


Figure 2: Empirical evidences for Claim 2.1 on different tasks. LLMs demonstrate improving prediction abilities as the number of in-context samples grows.

(c) Function Value Prediction. The goal of this experiment is to let the LLM (gpt-3) predict the values of a function on unseen data points given the values on the points with fixed intervals. Following Gruver et al. (2023), we report the t-interval cumulative negative log-likelihood CNLL =  $-\sum_{i}^{t} \log P(v_i|\text{prompt}_{i-1})$ , where  $v_i$  is the value of the function at data point i. It can be observed that the LLMs are good time series forecasters.

Under Claim 2.1, we establish the correspondence between LLMs and RL by using LLMs as model estimators in RL algorithms, which opens the door to creating a practical algorithm that combines the strengths of both LLMs and RL. LLMs excel in accuracy with minimal feedback, which improves the sample efficiency. LLMs can also refine estimates using new feedback as prompts, which avoids explicit parameter updates. RL algorithms benefit from online interaction to improve estimates and policies and have theoretical guarantees with optimal planning algorithms like value iteration. This LLM-RL correspondence inspires us to introduce a new framework in the next section, aiming to orchestrate the reasoning (learning and planning) and acting of LLMs.

# 3 Algorithm

Architecture of RAFA. By leveraging the LLM-RL correspondence in Section 2, we provide a principled framework for orchestrating reasoning and acting, namely "reason for future, act for now" (RAFA), in Algorithms 1 and 2. At the t-th step of Algorithm 1, the LLM agent invokes the reasoning routine, which learns from the memory buffer and plans a future trajectory over a long horizon ("reason for future" in Line 6), takes the initial action of the planned trajectory ("act for now" in Line 7), and stores the collected feedback (state, action, and reward) in the memory buffer (Line 8). Upon the state transition of the external environment, the LLM agent reinvokes the reasoning routine to replan another future trajectory from the new state (Line 6 following Line 9). To ensure the learning and planning stability, we impose the switching condition (Line 10) to decide whether to incorporate the newest chunk of history, i.e., the set difference  $\mathcal{D}_t - \mathcal{D}_{t_k}$ , into the information state, which is used in the reasoning routine as contexts. In other words, the reasoning routine uses the same

## Algorithm 1 Reason for future, act for now (RAFA): The LLM version.

- 1: **input**: An LLM learner-planner LLM-LR-PL, which aims at generating an optimal trajectory given an initial state and returns the initial action (e.g., Algorithm 2), and a switching condition If-Switch.
- 2: **initialization**: Sample the initial state  $s_0 \sim \rho$ , set t = 0, and initialize the memory buffer  $\mathcal{D}_0 = \emptyset$ .
- 3: **for** k = 0, 1, ...,**do**
- 4: Set  $t_k \leftarrow t$ .
- 5: repeat
- 6: Learn and plan given memory  $\mathcal{D}_{t_k}$  to get action  $a_t \leftarrow \text{LLM-LR-PL}(\mathcal{D}_{t_k}, s_t)$ . ("reason for future")
- 7: Execute action  $a_t$  to receive reward  $r_t$  and state  $s_{t+1}$  from environment. ("act for now")
- 8: Update memory  $\mathcal{D}_{t+1} \leftarrow \mathcal{D}_t \cup \{(s_t, a_t, s_{t+1}, r_t)\}.$
- 9: Set  $t \leftarrow t + 1$ .
- 10: until If-Switch( $\mathcal{D}_t$ ) is True.

(the switching condition is satisfied)

11: end for

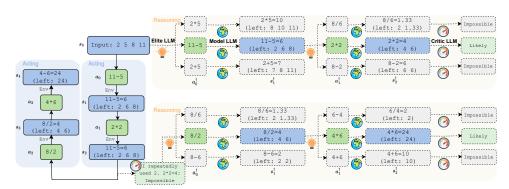


Figure 3: RAFA  $\operatorname{Game}$ of 24. Actions are proposed (dotted) and selected (green). Hallucinations that same number can be reused are mitigated through interactions.

history  $\mathcal{D}_{t_k}$  for all  $t_k \leq t < t_{k+1}$  until the (k+1)-th switch at the  $(t_{k+1}-1)$ -th step, which guarantees that the posterior distribution and the optimized action or the corresponding policy are updated in a conservative manner. We specify the switching condition in Sections 4 and 5.

"Reason for Future" (Line 6 in Algorithm 1 and Lines 3-11 in Algorithm 2). As detailed below, the reasoning routine composes the learning and planning subroutines to map the full history  $\mathcal{D}_{t_k}$  (until the  $t_k$ -th step) to an optimized action  $a_t$ . Note that the reasoning routine does not interact with the external environment throughout the learning and planning subroutines.

• The learning subroutine (Lines 3-4 in Algorithm 2) maps the memory buffer  $\mathcal{D}_{t_k}$  to a transition kernel (Model) and a value function (Critic), which are used in the planning subroutine. In practice, we prompt LLMs to form an estimate of the external environment based on the memory buffer. Here, the estimate is instantiated by two LLMs: Model and Critic, which estimate their ground-truth counterparts in association with the data-generating pa-

rameter. Under Claim 2.1, the learning subroutine yields more accurate versions of Model and Critic when we prompt them with a growing collection of feedback from the external environment. Consequently, the planning subroutine can use them to assess the long-term outcome of actions with a higher accuracy. Depending on whether we emulate the model-based or model-free approach of RL, we may choose to emulate Model or Critic individually. Compared with the learning subroutine in RL, we replace the parameterized function approximation (usually deep neural networks) with LLMs and use an "in-context" manner to update the LLMs, which eliminates the need for explicit parameter updates. Because LLMs are pretrained and undergo supervised fine-tuning, they provide much better estimates compared to learning from scratch, leading to an improvement in sample efficiency for online interactions.

**Algorithm 2** The LLM learner-planner (LLM-LR-PL): A tree-search example. (the deterministic case)

- 1: **input**: The memory buffer  $\mathcal{D}$ , the initial state s, the search breadth B, and the search depth U.
- 2: **initialization**: Initialize the state array  $S_0 \leftarrow \{s\}$  and the action array  $A_0 \leftarrow \emptyset$ .

  (the learning subroutine)
- 3: Set Model as an LLM instance prompted to use  $\mathcal{D}$  as contexts to generate the next state.
- 5: Set Elite as an LLM instance prompted to use  $\mathcal{D}$  as contexts to generate multiple candidate actions.
- 6: **for** u = 0, ..., U **do**
- 7: For each current state in  $S_u$ , invoke Elite to generate B candidate actions and store them in  $A_u$ .
- 8: For each candidate action in  $A_u$ , invoke Model to generate the next state and store it in  $S_{u+1}$ .
- 9: end for
- 10: For all resulting rollouts in  $S_0 \times A_0 \times \cdots \times S_U \times A_U$ , invoke Critic to evaluate the expected cumulative future reward and select the best one  $(s_0^{\dagger}, a_0^{\dagger}, \dots, s_U^{\dagger}, a_U^{\dagger})$ , where  $s_0^{\dagger} = s$ .
- 11: **output**: The initial action  $a_0^{\dagger}$  of the selected rollout.
- The planning subroutine (Lines 5-11 in Algorithm 2) maps Model and Critic to a future trajectory  $(s_0^{\dagger}, a_0^{\dagger}, \dots, s_U^{\dagger}, a_U^{\dagger})$ , where  $s_0^{\dagger}$  is the current state  $s_t$  and  $a_0^{\dagger}$  is executed in the external environment as the current action  $a_t$  during the acting phase. Intuitively, we prompt LLMs to generate an optimal policy (actor) for multiple future steps, which maximizes the value function (Critic). From an RL perspective (Sections 2 and 5), the planning subroutine approximately solves the Bellman equation (Sutton and Barto, 2018), where we solve the optimal policy (or the corresponding action) given the estimated transition kernel and reward function (or critic) by LLMs. As two LLM instances from the learning subroutine, Model and

Critic instantiate the estimated transition kernel and the estimated value function. Hence, we can simulate a given number of trajectories with Model, evaluate them with Critic, and obtain an improved policy, which is achieved by specially designed prompts instead of a numerical algorithm. By maximizing the expected cumulative future reward (instead of the immediate reward), the planning subroutine returns an optimized action that improves the long-term outcome. There are two error sources that affect the planning subroutine, namely the posterior uncertainty, which is inherited from Model and Critic due to the finite size of  $\mathcal{D}_{t_k}$ , and the planning suboptimality, which is induced by the limited capacity for computation, e.g., the bounded width and depth of tree-search (Lines 6-9 in Algorithm 2). Depending on the specific configuration of the state and action spaces (continuous versus discrete) and the transition and reward models (stochastic versus deterministic), we may choose to emulate the value iteration algorithm, the random shooting algorithm, or the Monte-Carlo tree-search algorithm. All of them allow RAFA to achieve provable sample efficiency guarantees as long as they satisfy a specific requirement of optimality (Definition 5.1). For illustration, we emulate the tree-search algorithm and defer its stochastic variant to Appendix B.

"Act for Now" (Lines 7-10 in Algorithm 1). At the current state  $s_t$ , the LLM agent executes the optimized action  $a_t$  in the external environment, which is obtained from the reasoning routine. Specifically, we take the initial action  $a_0^{\dagger}$  of the planned trajectory  $(s_0^{\dagger}, a_0^{\dagger}, \dots, s_U^{\dagger}, a_U^{\dagger})$ , where  $s_0^{\dagger} = s_t$  and  $a_0^{\dagger} = a_t$ , and discard the remaining subset. At the next state  $s_{t+1}$ , the LLM agent replans another future trajectory  $(s_0^{\dagger}, a_0^{\dagger}, \dots, s_U^{\dagger}, a_U^{\dagger})$  with  $s_0^{\dagger} = s_{t+1}$  and  $a_0^{\dagger} = a_{t+1}$ . In other words, the acting phase follows a short-term subset of the long-term plan, which is regenerated at every new state. The LLM agent stores the collected feedback  $(s_t, a_t, r_t, s_{t+1})$  in the memory buffer  $\mathcal{D}_t$  and queries a switching condition If-Switch to decide when to update the memory buffer  $\mathcal{D}_{t_k}$ , which is used in the reasoning routine as contexts for learning and planning. Intuitively, we incorporate the newest chunk of history  $\mathcal{D}_t - \mathcal{D}_{t_k}$  to improve the current policy only when it carries significant novel information, e.g., when the LLM agent loses for the first time following a winning streak.

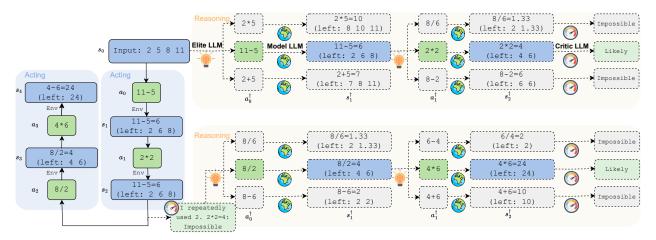


Figure 4: RAFA for Game of 24. Actions are proposed (dotted) and selected (green). Hallucinations that the same number can be reused are mitigated through interactions.

# 4 Experiment

We evaluate RAFA in several text-based benchmarks, e.g., Game of 24, ALFWorld, BlocksWorld, and Tic-Tac-Toe. The detailed setups, results, and ablations are provided in Appendix F, while the detailed prompts are found in Appendix G. We release all the codes on the page: https://agentification.github.io/RAFA.

#### 4.1 Game of 24

Game of 24 (Yao et al., 2023a) is a mathematical puzzle to obtain 24 from four natural numbers through basic arithmetic operations. The state is the (possibly unfinished) current formula and the action is the next formula (or the modified part).

**Setup.** We emulate the tree-search algorithm to plan  $(B \in \{1, 2\})$ . At the t-th step, RAFA learns from the memory buffer and switches to a new policy upon receiving an unexpected reward, which is the switching condition. After the t-th step, RAFA digests the collected feedback and generates a linguistic summary, which is saved into the memory buffer to avoid similar previous mistakes.

	$   \ {\tt RAFA} \ (B=1)$	${\tt RAFA}\ (B=2)$	$\operatorname{ToT} (B=1)$	$\operatorname{ToT}\ (B=2)$	Reflexion
gpt-4	89%	93%	73%	81%	21%
gpt-3.5	29%	46%	10%	17%	16%

Table 2: Game of 24 results.

**Result.** RAFA attains SOTA performances as shown in Table 2. RAFA achieves superior sample efficiency by mitigating hallucinations and avoid careless trials (Figures 4 and 5).

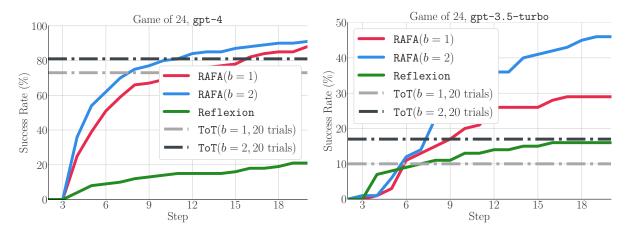


Figure 5: Sample efficiency on Game of 24.

	Pick	Clean	Heat	Cool	Examine	PickTwo	Total
BUTLER	46.00	39.00	74.00	100.00	22.00	24.00	37.00
ReAct	66.67	41.94	91.03	80.95	55.56	35.29	61.94
AdaPlanner	100.00	96.77	95.65	100.00	100.00	47.06	91.79
Reflexion	100.00	90.32	82.61	90.48	100.00	94.12	92.54
RAFA	100.00	96.77	100.00	100.00	100.00	100.00	99.25

Table 3: ALFWorld results (success rates %).

### 4.2 ALFWorld

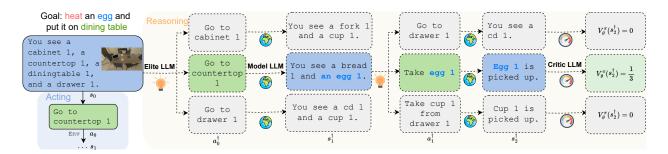


Figure 6: An illustration of RAFA in the ALFWorld environment.

ALFWorld (Shridhar et al., 2020) is an interactive environment for embodied agent simulations, which encompasses 134 household tasks in six overall categories (Table 3). We use gpt-3 (text-davinci-003).

**Setup.** We emulate the tree-search algorithm to plan (B=2). RAFA invokes Critic to evaluate the completed portion of the desired goal and switches to a new policy after 20 consecutive failures.

Result. RAFA outperforms various existing frameworks (Figure 7). The better performance of AdaPlanner at the initial episode is attributed to a handcrafted set of programs for rejecting suboptimal candidate trajectories, which is challenging to construct without the domain knowledge of a specific task. One such example is the Pick-Two category.

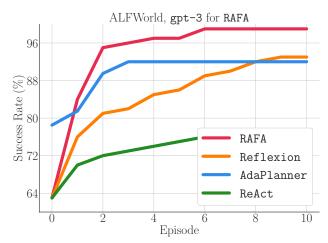


Figure 7: Sample efficiency on ALFWorld.

## 4.3 BlocksWorld

BlocksWorld (Hao et al., 2023) contains tasks to arrange blocks in specific configurations.

**Setup.** We use the Vicuna (Zheng et al., 2023) model and emulate the MCTS algorithm to plan.

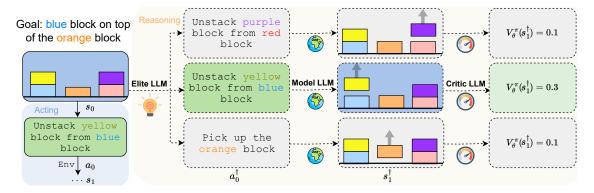


Figure 8: RAFA for BlocksWorld.

**Result.** RAFA achieves superior success rates across multiple Vicuna versions (Figure 9). Comparisons with CoT and RAP demonstrate how the learning subroutine improves the planning optimality.

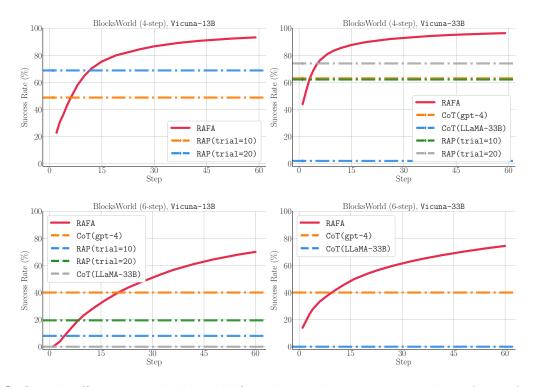


Figure 9: Sample efficiency on BlocksWorld (4 and 6 are the minimum numbers of steps for solving a specific task). CoT is prompted by four in-context examples.

## 4.4 Tic-Tac-Toe

Tic-Tac-Toe (Beck, 2008) is a competitive game where the X and O sides take turns to place marks. RAFA invokes Model to simulate the transition and opponent dynamics (Figure 18).

**Setup.** We use gpt-4 and emulate the tree-search algorithm to plan  $(B \in \{3,4\})$ . RAFA switches to a new policy when (a) the predicted state differs from the observed one, (2) the predicted action of opponents differs from the observed one, or (3) Critic gives the wrong prediction of the game status. Here, X has an asymmetric advantage (winning surely if played properly).

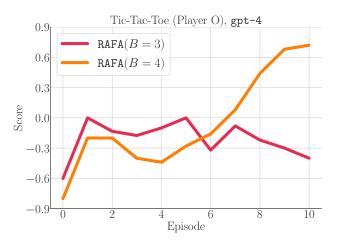


Figure 10: Sample efficiency on Tic-Tac-Toe (0 means tie).

**Result.** RAFA (playing O) matches and beats gpt-4 for T = 5 and T = 7 (Table 4), although O is destined to lose. The ablation study (B = 3 versus B = 4) illustrates how the planning suboptimality affects the sample efficiency (Figure 10).

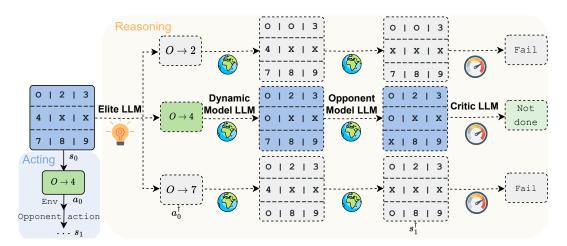


Figure 11: RAFA (playing O) for Tic-Tac-Toe.

XO	gpt-4	$\mathtt{RAFA}(T\!=\!1)$	$\mathtt{RAFA}(T\!=\!5)$	$\mathtt{RAFA}(T\!=\!7)$
gpt-4	90%, 0%, <b>10</b> %	90%, 0%, <b>10</b> %	50%, 0%, <b>50</b> %	0%, 0%, <b>100</b> %

Table 4: Tic-Tac-Toe Results. We set B=4 and report the winning rate of X, the tie rate, and the winning rate of O.

# 5 Theoretical Analysis

In this section, we provide the theoretical results in this paper. In Section 5.1, we characterize the requirement for the planning subroutine in RAFA and show the value iteration algorithm with a truncated horizon can be an example of the desired planner. In Section 5.2, we show that the LLM with a posterior alignment performs BMA, which supports Claim 2.1 in theory. We present the regret analysis for RAFA in Section 5.3 to explain its superior empirical performance, where we provide necessary assumptions and the regret bound of RAFA. In Section 5.4, we show that RAFA can be modified to encourage efficient exploration for more complex tasks such that RAFA is still sample-efficient without the concentrability assumption in Section 5.3.

## 5.1 Planning Optimality

To characterize the requirement for the planning subroutine in RAFA (Algorithm 1), we define the  $\epsilon$ -optimal planner as follows.

**Definition 5.1** ( $\epsilon$ -Optimal Planner). Denote  $\{V \mid V \text{ is a value function}\}$  as  $\mathcal{V}$ . A planning algorithm  $\mathsf{PL}^{\epsilon}: \mathcal{P} \times \mathcal{R} \mapsto \Pi \times \mathcal{V} \text{ is an } \epsilon\text{-optimal planner if } \mathsf{PL}^{\epsilon}(P,r) = (\pi,V), \text{ where } |Q(s,a) - r(s,a) - \gamma \cdot (PV)(s,a)| \leq \epsilon \text{ and } V(s) = \max_{a} Q(s,a) = Q(s,\pi(s)) \text{ for all } (s,a) \in \mathcal{S} \times \mathcal{A}.$ 

In other words, an  $\epsilon$ -optimal planner with a model (transition kernel and reward function) can generate a policy to approximately maximize the corresponding long-term value function instead of the myopic reward with an approximate error limit  $\epsilon$ . As an instance of the planner satisfying Definition 5.1, we present the value iteration algorithm (Algorithm 3) with a truncated horizon U, i.e., a finite length of the lookahead window as the  $\epsilon$ -optimal planner in Algorithm 4. The following proposition ensures that Algorithm 3 satisfies Definition 5.1.

**Proposition 5.2.** Algorithm 3 is an  $\epsilon$ -optimal planner as long as we set  $U \ge 1 + \lceil \log_{\gamma}(\epsilon/L) \rceil$  and any value function is bounded by  $L \ge 0$ .

*Proof.* See Appendix C.1 for a detailed proof.

#### **Algorithm 3** $\epsilon$ -Optimal planner: The value iteration algorithm with a truncated horizon.

- 1: **input**: The model (P, r) and the truncated horizon U.
- 2: **initialization**: Set the value function  $V_{\theta}^{(U)}(\cdot) \leftarrow 0$ .
- 3: **for** u = U 1, ..., 1 **do**
- 4: Set the value function  $V^{(u)}(\cdot) \leftarrow \max_{a \in \mathcal{A}} Q^{(u)}(\cdot, a)$ , where  $Q^{(u)}(\cdot, \cdot) \leftarrow r(\cdot, \cdot) + \gamma(PV^{(u+1)})(\cdot, \cdot)$ .
- 5: end for
- 6: **output**: The greedy policy  $\pi(\cdot) = \arg \max_{a \in \mathcal{A}} Q^{(1)}(\cdot, a)$  and the value function  $V^{(1)}$ .

Alternatively, we may choose to emulate the tree-search algorithm, the random shooting algorithm, or the Monte-Carlo tree-search algorithm. In the tree-search example (Lines 5-11 in Algorithm 2),  $\epsilon$  decreases as the search breadth B and depth U increase. Note that, as long as we emulate an  $\epsilon$ -optimal planner, we are able to establish provable sample efficiency guarantees.

# 5.2 LLMs with Posterior Alignments Perform Bayesian Model Averaging (BMA)

In the following, we analyze Claim 2.1 from the theoretical perspective. For LLMs used in RAFA, we denote  $P^{\text{LLM}}(\xi_{(s,a)} | \mathcal{D}, s, a)$  as the probability measure of the predicted state-reward pair given the query state-action pair and the memory buffer  $\mathcal{D}$  as the in-context dataset. Induced by  $P^{\text{LLM}}$ , we also denote  $P_{\text{LLM}(\mathcal{D})}$  and  $r_{\text{LLM}(\mathcal{D})}$  as the estimated transition kernel and reward function by LLMs, respectively.

For the simplicity of analysis, we assume that all LLMs have posterior alignments in the tasks that we study, that is, their posterior distributions of the reward and the next state given the current state-action pair and any in-context dataset match the posteriors in these tasks. We formulate this assumption as follows.

**Assumption 5.3** (Posterior Alignment). We assume that LLMs are aligned with the posterior of the state and reward in the underlying MDP, which is formulated as

$$P^{\text{LLM}}(\xi_{(s,a)} | \mathcal{D}, s, a) = \mathbb{P}_{post}(\xi_{(s,a)} | \mathcal{D}, s, a),$$

for any in-context dataset  $\mathcal{D} = \{(s_i, a_i, r_i, s_i')\}_{i=0}^{I}$  with size I, query state-action pair (s, a), reward r, and state s'. Here, the posterior  $\mathbb{P}_{post}$  is defined in (2.4).

We remark that the posterior alignment in Assumption 5.3 comes from the in-context ability of LLMs, which is widely studied in Lee et al. (2023); Wies et al. (2024); Xie et al. (2021). We also remark that Assumption 5.3 does not mean that LLMs can trivially make the optimal decision at each step in the underlying MDP: (1) Though the posterior distributions of state and reward are aligned, LLMs still need to be instructed to maximize the long-term value (via explicit planning) instead of the myopic reward. (2) LLMs still require online interactions to enlarge the in-context dataset  $\mathcal{D}$  such that their prediction uncertainty can be reduced from the prior uncertainty. In Appendix C.7, we also discuss how to relax Assumption 5.3 to accommodate an additional error term in the regret bound derived by our analysis, where we assume that LLMs are maximum likelihood estimatiors (MLEs) on the pretraining dataset with uniform coverage. Based on Assumption 5.3, we prove that LLMs with posterior alignments perform BMA in the model estimation in the following proposition.

**Proposition 5.4** (LLMs with Posterior Alignments Perform BMA). *Under Assumption 5.3*, the *LLM predictions satisfy* 

$$r_{\mathit{LLM}(\mathcal{D})}(s,a) + \gamma \cdot (P_{\mathit{LLM}(\mathcal{D})}V)(s,a) = \mathbb{E}_{\theta \sim \mathbb{P}_{post}(\cdot \mid \mathcal{D})}[(B_{\theta}V)(s,a)]$$

for any dataset  $\mathcal{D} = \{(s_i, a_i, r_i, s_i')\}_{i=0}^{I}$  with size I, value function V, and query state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . Here,  $\mathbb{P}_{post}(\theta \mid \mathcal{D})$  is the posterior of  $\theta^*$  given  $\mathcal{D}$  in the underlying MDP.

Proof of Proposition 5.4. See the detailed proof in Appendix C.2.

The proof of Proposition 5.4 can be found in Appendix C.2. Some variants of Proposition 5.4 can be found in various literature (Lee et al., 2023; Zhang et al., 2022, 2023b). In particular, Zhang et al. (2022) establish the theoretical equivalence between BMA and the ideal attention architecture and analyze the generalization error rate of LLMs. By Proposition 5.4, LLMs can provide a more certain and accurate estimate for the data-generating model with more collected feedback, as the uncertainty in the posterior is reduced with more data. Thus, Proposition 5.4 supports Claim 2.1 in theory.

## 5.3 Regret Bound of RAFA

Algorithm 4 Reason for future, act for now (RAFA): The theoretical version.

- 1: **input**: An  $\epsilon$ -optimal planner PL $^{\epsilon}$ , which returns an  $\epsilon$ -optimal policy that maximizes the value function up to an  $\epsilon$  accuracy (Definition 5.1), and LLMs with posterior alignments.
- 2: **initialization**: Sample the initial state  $s_0 \sim \rho$ , set t = 0, and initialize the memory buffer  $\mathcal{D}_0 = \emptyset$ .
- 3: **for** k = 0, 1, ...,**do**
- 4: Set  $t_k \leftarrow t$ .
- 5: repeat
- 6: Plan ahead with the ε-optimal planner and LLMs  $(\pi_t, V_t) \leftarrow \text{PL}^{\epsilon}(P_{\text{LLM}(\mathcal{D}_{t_k})}, r_{\text{LLM}(\mathcal{D}_{t_k})})$ .

  ("reason for future")
- 7: Execute action  $a_t = \pi_t(s_t)$  to receive reward  $r_t$  and state  $s_{t+1}$  from environment. ("act for now")
- 8: Update the memory buffer  $\mathcal{D}_{t+1} \leftarrow \mathcal{D}_t \cup \{(s_t, a_t, s_{t+1}, r_t)\}.$
- 9: Set  $t \leftarrow t + 1$ .
- 10: **until**  $H_{t_k} H_t > \log 2$ , where  $H_t$  denotes posterior entropy of  $\theta^*$  conditioned on  $\mathcal{D}_t$ . (the switching condition is satisfied)

## 11: end for

To analyze RAFA in theory, we propose the theoretical version of RAFA in Algorithm 4, where we instantiate the switching condition of RAFA in Line 10 by measuring the reduction of the posterior entropy. At the t-th step and the k-th switching times, Algorithm 6 only makes the (k+1)-th switch when the reduction of posterior entropy  $H_{t_k} - H_t$  is greater than log 2. In Line 6 of Algorithm 4, we describe the planning subroutine in RAFA (Algorithm 1) by an  $\epsilon$ -planner PL $^{\epsilon}$  (defined in Definition 5.1). We specify the terminating condition for Algorithm 4. Let (K-1) be the total number of switches until t reaches (T-1). Let  $t_K = T$ . At the (T-1)-th step, Algorithm 4 executes  $a_{T-1} = \pi_{T-1}(s_{T-1})$ , where

we have  $\pi_{T-1} = \operatorname{PL}^{\epsilon}(P_{\operatorname{LLM}(\mathcal{D}_{t_{K-1}})}, r_{\operatorname{LLM}(\mathcal{D}_{t_{K-1}})})$ . Upon receiving  $r_{T-1}$  and  $s_T$  from the external environment, Algorithm 4 updates  $\mathcal{D}_T = \{(s_t, a_t, s_{t+1}, r_t)\}_{t=0}^{T-1}$  and terminates. Since the agent in Algorithm 4 executes the same policy until making a switch, we have  $\pi_t = \pi_{t_k}$  for any  $t_k \leq k < t_{k+1}$ . We denote by  $\pi^k = \pi_{t_k}$  for the notational simplicity. Next, we impose a regularity assumption on the structure of MDPs to measure the learning difficulty. Recall that we define the posterior entropy  $H_t$  in (2.5), the information gain  $I(\theta; \xi \mid \mathcal{D})$ , and  $\xi_{(s,a)}$  as the pair of the next state and the current reward (s', r) given the query state-action pair (s, a) in Section 2. Define  $H_t$  as the posterior entropy  $H(\theta \mid \mathcal{D}_t)$  given the dataset  $\mathcal{D}_t = \{(s_i, a_i, r_i, s_{i+1})\}_{i=0}^{t-1}$ .

**Assumption 5.5** (MDPs Regularity). We assume that there exists a coefficient  $\eta > 0$  such that, if  $H_{t_1} - H_{t_2} \leq \log 2$ , then it holds that

$$I(\theta; \xi_{(s,a)} \mid \mathcal{D}_{t_1}) \leq 4\eta \cdot I(\theta; \xi_{(s,a)} \mid \mathcal{D}_{t_2})$$

for any given value function V,  $t_1 < t_2$  and  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .

Assumption 5.5 is a regularity assumption on MDPs and is intrinsic to the agent design. In Appendix E, we prove that d-dimensional Bayesian linear kernel MDPs (defined in Definition E.1), satisfy Assumption 5.5 with the coefficient  $\eta = d/\log(1+d)$ . Intuitively, Assumption 5.5 restricts the increase of the information gain given one bit (log 2) reduction of the posterior entropy.

Similar to other theoretical work on deep RL (Lazaric et al., 2010; Fan et al., 2020; Zhang et al., 2020), we introduce the concentrability coefficient  $\kappa$  to bound the distribution shift between the current policy and the optimal policy. For the simplicity of discussions, we define the optimal  $\gamma$ -discounted visitation measure  $\nu^*$  starting from state s as

$$\nu^{\star}(s'|s) = \frac{1}{1-\gamma} \cdot \sum_{\tau=0}^{\infty} \gamma^{\tau} \cdot \mathbb{P}(s_{\tau} = s'|s_{0} = s, s_{i+1} \sim P_{\theta^{\star}}(\cdot|s_{i}, \pi^{\star}(s_{i})) \text{ for any } 0 \leq i < \tau),$$

$$(5.1)$$

for any state  $s, s' \in \mathcal{S}$ . Here,  $\nu^*(\cdot | s)$  describes the discounted average probability measure of the state that the optimal policy  $\pi^*$  visits starting from state s in the underlying MDP. Now, we are ready to provide the full statement of the concentrability coefficient as follows.

**Assumption 5.6** (Concentrability). For RAFA (Algorithm 4), we assume that there exists a constant  $\kappa < \infty$  such that

$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \mathbb{E}_{\theta^{\star} \sim \mathbb{P}_{t_k}} \left[\frac{\mathbb{E}_{s \sim \nu^{\star}(\cdot \mid s_t)} \left[\left((B_k - B_{\theta^{\star}})V_t\right)^2 (s, \pi^{\star}(s))\right]}{\left((B_k - B_{\theta^{\star}})V_t\right)^2 (s_t, \pi^k(s_t))} \middle| \mathcal{D}_{t_k}\right]\right]\right]$$

is bounded by  $\kappa^2 \cdot T$ , where we define  $(B_k V)(s,a) = r_{\text{LLM}(\mathcal{D}_{t_k})}(s,a) + \gamma \cdot (P_{\text{LLM}(\mathcal{D}_{t_k})}V)(s,a)$  and denote by  $\mathbb{P}_{t_k}$  the posterior of  $\theta^*$  given  $\mathcal{D}_{t_k}$ .

Intuitively,  $\kappa$  measures the hardness to generalize the low prediction error  $(B_k - B_{\theta^*})V_t$  on the current trajectory induced by  $\pi^*$  in the underlying MDP. We remark that we can drop the dependency of the concentrability coefficient  $\kappa$  (Assumption 5.6) if we modify RAFA to encourage efficient exploration in MDPs. We will discuss the variants of RAFA with efficient exploration strategies in Section 5.4.

In the following theorem, we give the bound of the Bayesian regret of RAFA (Algorithm 4) as follows.

**Theorem 5.7.** Under Assumptions 5.3, 5.5, and 5.6, the Bayesian regret of RAFA (Algorithm 4) satisfies

$$\mathfrak{R}(T) = \mathcal{O}\left(\frac{(\kappa+1)L \cdot \sqrt{\mathbb{E}[H_0 - H_T]}}{1 - \gamma} \cdot \sqrt{T} + \frac{\epsilon}{1 - \gamma} \cdot T + \frac{L \cdot \mathbb{E}[H_0 - H_T]}{1 - \gamma}\right),$$

where  $\kappa$  is the concentrability coefficient defined in Assumption 5.6,  $H_t$  is the posterior entropy of  $\theta^*$  given the history  $\mathcal{D}_t = \{(s_i, a_i, r_i, s_{i+1})\}_{i=0}^{t-1}$ , and L is the bound of |r + V(s)| for any reward r, state s, and value function V.

Proof of Theorem 5.7. See the detailed proof in Appendix C.4.

Theorem 5.10 establishes the  $\sqrt{T}$  regret of RAFA (Algorithm 4) for a proper choice of the planning suboptimality  $\epsilon$ , e.g.,  $\epsilon = \mathcal{O}(1/\sqrt{T})$ , which shows that RAFA is sample efficient and explains its strong empirical performance in Section 4. Here, the first term in the upper bound in Theorem 5.10 is the leading term and involves several multiplicative factors, namely the effective horizon  $1/(1-\gamma)$ , the value bound L, and the cumulative posterior entropy reduction  $H_0 - H_T$  throughout the T steps, which are common in the RL literature (Abbasi-Yadkori and Szepesvári, 2015; Osband et al., 2013; Russo and Van Roy, 2014a,b, 2016; Lu and Van Roy, 2019). In particular,  $H_0$  highlights the prior knowledge obtained through pretraining, as  $H_0$  quantifies the prior uncertainty of LLMs before incorporating any collected feedback. Hence,  $H_0 - H_T$  highlights the uncertainty reduction achieved by reasoning and acting, as  $H_T$  quantifies the posterior uncertainty of LLMs after incorporating the collected feedback. In Appendix  $\mathbf{E}$ , we prove that  $H_0 - H_T = \mathcal{O}(d \cdot \log T)$  and the  $1 - \delta$  probability bound on value functions  $L = \mathcal{O}(\sqrt{d} \cdot \log(dT/\delta))$  for the d-dimensional Bayesian linear kernel MDPs, which implies  $\Re(T) = \widetilde{\mathcal{O}}((1-\gamma)^{-1}(\kappa+1) \cdot \sqrt{d^3T})$  with probability at least  $1 - \delta$ . Here  $\widetilde{\mathcal{O}}$  hides the logarithmic factor.

# 5.4 RAFA with Efficient Exploration Strategies

To drop the dependency of Assumption 5.6 (Concentrability) and solve more complex tasks, we provide two variants of RAFA (Algorithm 4): (1) RAFA with an optimistic bonus (Algorithm 5) and (2) RAFA with posterior sampling (Algorithm 6). We also prove the bound of the Bayesian regret of each variant, which demonstrates the effectiveness of these efficient exploration strategies without Assumption 5.6 (Concentrability).

Algorithm 5 Reason for future, act for now (RAFA): The theoretical version with an optimistic bonus

- 1: **input**: An  $\epsilon$ -optimal planner PL $^{\epsilon}$ , which returns an  $\epsilon$ -optimal policy that maximizes the value function up to an  $\epsilon$  accuracy (Definition 5.1), and LLMs with posterior alignments.
- 2: **initialization**: Sample the initial state  $s_0 \sim \rho$ , set t = 0, and initialize the memory buffer  $\mathcal{D}_0 = \emptyset$ .
- 3: **for** k = 0, 1, ...,**do**
- 4: Set  $t_k \leftarrow t$ .
- 5: repeat
- 6: Design optimistic bonus  $\Gamma_k(s, a) = \sqrt{2L} \cdot \sqrt{I(\theta; \xi_{(s,a)} \mid \mathcal{D}_{t_k})}$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .
- 7: Plan ahead with the  $\epsilon$ -optimal planner and LLMs  $(\pi_t, V_t) \leftarrow \text{PL}^{\epsilon}(P_{\text{LLM}(\mathcal{D}_{t_k})}, r_{\text{LLM}(\mathcal{D}_{t_k})} + \Gamma_k).$

("reason for future")

- 8: Execute action  $a_t = \pi_t(s_t)$  to receive reward  $r_t$  and state  $s_{t+1}$  from environment. ("act for now")
- 9: Update memory  $\mathcal{D}_{t+1} \leftarrow \mathcal{D}_t \cup \{(s_t, a_t, s_{t+1}, r_t)\}.$
- 10: Set  $t \leftarrow t + 1$ .
- 11: **until**  $H_{t_k} H_t > \log 2$ , where  $H_t$  denotes posterior entropy of  $\theta^*$  conditioned on  $\mathcal{D}_t$ . (the switching condition is satisfied)

12: end for

## 5.4.1 Optimistic Bonus

We incorporate the *Optimism in Face of Uncertainty* (OFU) principle (Cai et al., 2020; Zhou et al., 2021b; Jin et al., 2020; Liu et al., 2022b; Wang et al., 2023c) to encourage efficient exploration by adding an optimistic bonus on the reward function in the planning subroutine of RAFA. We design the optimistic bonus by the information gain and implement a variant of RAFA in Algorithm 5. In particular, the bonus  $\Gamma_k(s, a)$  takes the following form

$$\Gamma_k(s, a) = \sqrt{2L} \cdot \sqrt{I(\theta; \xi_{(s,a)} \mid \mathcal{D}_{t_k})}$$
(5.2)

for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$  and k < K. In Line 7 of Algorithm 5, we generate the policy  $\pi^t$  by  $PL^{\epsilon}(P_{\text{LLM}(\mathcal{D}_{t_k})}, r_{\text{LLM}(\mathcal{D}_{t_k})} + \Gamma_k)$  for any  $t_k \leq t < t_{k+1}$ . Intuitively, the bonus is higher at the state-action pair with higher information gain, which incentivizes the agent to explore those less visited states (with higher information gain). In the following theorem, we prove the regret bound of RAFA with an optimistic bonus (Algorithm 5).

**Theorem 5.8.** Under Assumptions 5.3 and 5.5, the Bayesian regret of RAFA with an optimistic bonus (Algorithm 5) satisfies

$$\mathfrak{R}(T) = \mathcal{O}\left(\frac{L \cdot \sqrt{\mathbb{E}[H_0 - H_T]}}{1 - \gamma} \cdot \sqrt{T} + \frac{\epsilon}{1 - \gamma} \cdot T + \frac{L \cdot \mathbb{E}[H_0 - H_T]}{1 - \gamma}\right),$$

where all the variables have the same definitions in Theorem 5.7.

*Proof of Theorem 5.8.* See the detailed proof in Appendix C.5.

Compared with Theorem 5.7, the regret bound in Theorem 5.8 is not dependent on the concentrability coefficient  $\kappa$ , which demonstrates the effectiveness of the optimistic bonus in Algorithm 5. In Appendix E, we prove that  $H_0 - H_T = \mathcal{O}(d \cdot \log T)$  and the  $1 - \delta$  probability bound on value functions  $L = \mathcal{O}(\sqrt{d} \cdot \log(dT/\delta))$  for the d-dimensional Bayesian linear kernel MDPs, which implies  $\mathfrak{R}(T) = \widetilde{\mathcal{O}}((1-\gamma)^{-1} \cdot \sqrt{d^3T})$  with probability at least  $1 - \delta$ . Here  $\widetilde{\mathcal{O}}$  hides the logarithmic factor.

#### 5.4.2 Posterior Sampling

**Algorithm 6** Reason for future, act for now (RAFA): The theoretical version with posterior sampling.

- 1: **input**: An  $\epsilon$ -optimal planner PL $^{\epsilon}$ , which returns an  $\epsilon$ -optimal policy that maximizes the value function up to an  $\epsilon$  accuracy (Definition 5.1), and LLMs satisfying Assumption 5.9.
- 2: **initialization**: Sample the initial state  $s_0 \sim \rho$ , set t = 0, and initialize the memory buffer  $\mathcal{D}_0 = \emptyset$ .
- 3: **for** k = 0, 1, ..., do
- 4: Set  $t_k \leftarrow t$ .
- 5: repeat
- 6: Plan ahead with the  $\epsilon$ -optimal planner and the posterior sampling mechanism of LLMs (defined in Assumption 5.9)  $(\pi_t, V_t) \leftarrow \text{PL}^{\epsilon}(P_{\text{LLM+PS}(\mathcal{D}_{t_k})}, r_{\text{LLM+PS}(\mathcal{D}_{t_k})}).$

("reason for future")

- 7: Execute action  $a_t = \pi_t(s_t)$  to receive reward  $r_t$  and state  $s_{t+1}$  from environment. ("act for now")
- 8: Update memory  $\mathcal{D}_{t+1} \leftarrow \mathcal{D}_t \cup \{(s_t, a_t, s_{t+1}, r_t)\}.$
- 9: Set  $t \leftarrow t + 1$
- 10: **until**  $H_{t_k} H_t > \log 2$ , where  $H_t$  denotes posterior entropy of  $\theta^*$  conditioned on  $\mathcal{D}_t$ . (the switching condition is satisfied)

#### 11: end for

As another method for efficient exploration, we assume that there exists a mechanism that deploys posterior sampling and we use this mechanism to encourage exploration for RAFA.

**Assumption 5.9** (LLMs with Posterior Sampling Mechanism). We assume that there exists a mechanism LLM+PS that maps the memory buffer  $\mathcal{D}$  to the transition kernel and the reward function, such that  $(r_{LLM+PS(\mathcal{D})}(s,a) + \gamma \cdot (P_{LLM+PS(\mathcal{D})}V)(s,a)) \mid \mathcal{D}$  and  $(B_{\theta^*}V(s,a)) \mid \mathcal{D}$  are identically independent distributed for any  $(s,a) \in \mathcal{S} \times \mathcal{A}$ , in-context dataset  $\mathcal{D}$ , and value function V. Here,  $\theta^*$  is the data-generating parameter.

We remark that the bootstrap method (Efron, 1982) can approximate the posterior sampling mechanism satisfying Assumption 5.9. Widely used in applied statistics (Davison and Hinkley, 1997) and the design of RL algorithms (Osband et al., 2016; Hao et al., 2019), the bootstrap method takes a dataset D and a functional estimator g as inputs. Depending on the configuration of bootstrap, we generate the bootstrapped dataset  $\widetilde{D}$  from D by uniform sampling with replacement (Efron, 1982) or weighted sampling with replacement (Newton and Raftery, 1994). Viewing the LLM as the functional estimator g and the memory buffer  $\mathcal{D}$  as the dataset D, we can use this bootstrap method to approximate the mechanism LLM+PS that is introduced in Assumption 5.9. From the statistics literature (Bickel and Freedman, 1981; Singh, 1981; Newton and Raftery, 1994), we also know that bootstrap distribution recovers the posterior distribution asymptotically.

Based on the mechanism satisfying Assumption 5.9, we propose a variant of RAFA in Algorithm 6, where we use the mechanism LLM+PS as the model estimator in the learning subroutine of RAFA. In Line 7 of Algorithm 5, we generate the policy  $\pi^t$  by  $\text{PL}^{\epsilon}(P_{\text{LLM+PS}(\mathcal{D}_{t_k})}, r_{\text{LLM+PS}(\mathcal{D}_{t_k})})$ . In the following, we give a simple explanation of how this mechanism helps the agent to explore efficiently. By the Bayes' rule, we have  $p(\theta \mid \mathcal{D}) \propto L(\mathcal{D} \mid \theta) \mathbb{P}_0(\theta)$ , where  $L(\mathcal{D} \mid \theta)$  is the likelihood of  $\mathcal{D}$  given  $\theta$  and  $\mathbb{P}_0$  is the prior of  $\theta^*$ . Taking the logarithm, we have  $\log(p(\theta \mid \mathcal{D})) = c + \log(\mathbb{P}_0(\theta)) + \log(L(\mathcal{D} \mid \theta))$  for some constant c. Hence, the uncertainty of the posterior is higher  $(p(\theta \mid \mathcal{D}))$  is closer to 0) at the less visited states (the likelihood of these states is closer to 0). Suppose we sample the model estimator from the posterior. In that case, the agent has more incentives to explore the less visited states, which explains why the mechanism LLM+PS encourages the efficient exploration.

In the following theorem, we prove the regret bound of RAFA with posterior sampling (Algorithm 6).

**Theorem 5.10** (Bayesian Regret). Under Assumptions 5.5 and 5.9, the Bayesian regret of RAFA with posterior sampling (Algorithm 6) satisfies

$$\mathfrak{R}(T) = \mathcal{O}\left(\frac{L \cdot \sqrt{\mathbb{E}[H_0 - H_T]}}{1 - \gamma} \cdot \sqrt{T} + \frac{\epsilon}{1 - \gamma} \cdot T + \frac{L \cdot \mathbb{E}[H_0 - H_T]}{1 - \gamma}\right),$$

where all the variables have the same definitions in Theorem 5.7.

*Proof of Theorem 5.10.* See the detailed proof in Appendix C.6.

Compared with Theorem 5.7, the regret bound in Theorem 5.10 is not dependent on the concentrability coefficient  $\kappa$ , which demonstrates the effectiveness of the posterior sampling mechanism in Algorithm 6. In Appendix E, we prove that  $H_0 - H_T = \mathcal{O}(d \cdot \log T)$  and the  $1 - \delta$  probability bound on value functions  $L = \mathcal{O}(\sqrt{d} \cdot \log(dT/\delta))$  for the d-dimensional Bayesian linear kernel MDPs, which implies  $\mathfrak{R}(T) = \widetilde{\mathcal{O}}((1 - \gamma)^{-1} \cdot \sqrt{d^3T})$  with probability at least  $1 - \delta$ . Here  $\widetilde{\mathcal{O}}$  hides the logarithmic factor.

# 6 Conclusions

In this paper, we establish the LLM-RL correspondence and propose a principled framework RAFA for orchestrating reasoning and acting, which achieves provable sample efficiency guarantees in autonomous LLM agents for the first time. We prove the  $\sqrt{T}$  regret bound of RAFA to highlight the synergy between prior knowledge from pretraining and the iterative process of reasoning and acting. RAFA's outstanding empirical performance underscores its potential for autonomous and adaptive decision-making in various complex tasks, which we remain for future work.

# Acknowledgement

Zhaoran Wang acknowledges National Science Foundation (Awards 2048075, 2008827, 2015568, 1934931), Simons Institute (Theory of Reinforcement Learning), Amazon, J.P. Morgan, and Two Sigma for their supports.

## References

- Abbasi-Yadkori, Y., Pál, D. and Szepesvári, C. (2011). Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*.
- Abbasi-Yadkori, Y. and Szepesvári, C. (2015). Bayesian optimal control of smoothly parameterized systems. In *Uncertainty in Artificial Intelligence*.
- Abernethy, J., Agarwal, A., Marinov, T. V. and Warmuth, M. K. (2023). A mechanism for sample-efficient in-context learning for sparse retrieval tasks. arXiv preprint arXiv:2305.17040.
- Agarwal, A., Kakade, S., Krishnamurthy, A. and Sun, W. (2020). Flambe: Structural complexity and representation learning of low rank mdps. *Advances in neural information processing systems*, **33** 20095–20107.
- Ahn, M., Brohan, A., Brown, N., Chebotar, Y., Cortes, O., David, B., Finn, C., Fu, C., Gopalakrishnan, K., Hausman, K. et al. (2022). Do as I can, not as I say: Grounding language in robotic affordances. arXiv preprint arXiv:2204.01691.
- Akyürek, E., Schuurmans, D., Andreas, J., Ma, T. and Zhou, D. (2022). What learning algorithm is in-context learning? Investigations with linear models. arXiv preprint arXiv:2211.15661.
- Beck, J. (2008). Combinatorial games: Tic-Tac-Toe theory.
- Betts, J. T. (1998). Survey of numerical methods for trajectory optimization. *Journal of Guidance, Control, and Dynamics*.

- Bickel, P. J. and Freedman, D. A. (1981). Some asymptotic theory for the bootstrap. *The annals of statistics*, **9** 1196–1217.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A. et al. (2020). Language models are few-shot learners. *Advances in neural information processing systems*, **33** 1877–1901.
- Cai, Q., Yang, Z., Jin, C. and Wang, Z. (2020). Provably efficient exploration in policy optimization. In *International Conference on Machine Learning*.
- Cai, T., Wang, X., Ma, T., Chen, X. and Zhou, D. (2023). Large language models as tool makers. arXiv preprint arXiv:2305.17126.
- Chen, Y., He, J. and Gu, Q. (2022). On the sample complexity of learning infinite-horizon discounted linear kernel MDPs. In *International Conference on Machine Learning*.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S. et al. (2022). PaLM: Scaling language modeling with pathways. arXiv preprint arXiv:2204.02311.
- Chua, K., Calandra, R., McAllister, R. and Levine, S. (2018). Deep reinforcement learning in a handful of trials using probabilistic dynamics models. *Advances in neural information processing systems*, **31**.
- Creswell, A. and Shanahan, M. (2022). Faithful reasoning using large language models. arXiv preprint arXiv:2208.14271.
- Creswell, A., Shanahan, M. and Higgins, I. (2022). Selection-inference: Exploiting large language models for interpretable logical reasoning. arXiv preprint arXiv:2205.09712.
- Davison, A. C. and Hinkley, D. V. (1997). *Bootstrap methods and their application*. 1, Cambridge university press.
- Dong, K., Wang, Y., Chen, X. and Wang, L. (2019). Q-learning with UCB exploration is sample efficient for infinite-horizon MDP. arXiv preprint arXiv:1901.09311.
- Efron, B. (1982). The jackknife, the bootstrap and other resampling plans. SIAM.
- Fan, J., Wang, Z., Xie, Y. and Yang, Z. (2020). A theoretical analysis of deep q-learning. In *Learning for dynamics and control*. PMLR.
- Garg, S., Tsipras, D., Liang, P. S. and Valiant, G. (2022). What can transformers learn incontext? A case study of simple function classes. In *Advances in Neural Information Processing Systems*.
- Ghavamzadeh, M., Mannor, S., Pineau, J., Tamar, A. et al. (2015). Bayesian reinforcement learning: A survey. Foundations and Trends® in Machine Learning, 8 359–483.

- Ghosh, M. (2021). Exponential tail bounds for chisquared random variables. *Journal of Statistical Theory and Practice*, **15** 35.
- Gruver, N., Finzi, M., Qiu, S. and Wilson, A. G. (2023). Large language models are zero-shot time series forecasters. arXiv preprint arXiv:2310.07820.
- Guo, H., Liu, Z., Zhang, Y. and Wang, Z. (2024). Can large language models play games? a case study of a self-play approach. arXiv preprint arXiv:2403.05632.
- Hafner, D., Lillicrap, T., Fischer, I., Villegas, R., Ha, D., Lee, H. and Davidson, J. (2019). Learning latent dynamics for planning from pixels. In *International conference on machine learning*. PMLR.
- Hao, B., Abbasi Yadkori, Y., Wen, Z. and Cheng, G. (2019). Bootstrapping upper confidence bound. Advances in neural information processing systems, 32.
- Hao, S., Gu, Y., Ma, H., Hong, J. J., Wang, Z., Wang, D. Z. and Hu, Z. (2023). Reasoning with language model is planning with world model. arXiv preprint arXiv:2305.14992.
- Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., Casas, D. d. L., Hendricks, L. A., Welbl, J., Clark, A. et al. (2022). Training compute-optimal large language models. arXiv preprint arXiv:2203.15556.
- Huang, W., Abbeel, P., Pathak, D. and Mordatch, I. (2022a). Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*.
- Huang, W., Xia, F., Xiao, T., Chan, H., Liang, J., Florence, P., Zeng, A., Tompson, J., Mordatch, I., Chebotar, Y. et al. (2022b). Inner monologue: Embodied reasoning through planning with language models. arXiv preprint arXiv:2207.05608.
- Janner, M., Fu, J., Zhang, M. and Levine, S. (2019). When to trust your model: Model-based policy optimization. Advances in neural information processing systems, **32**.
- Jiang, H. (2023). A latent space theory for emergent abilities in large language models. arXiv preprint arXiv:2304.09960.
- Jin, C., Yang, Z., Wang, Z. and Jordan, M. I. (2020). Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*. PMLR.
- Kim, G., Baldi, P. and McAleer, S. (2023). Language models can solve computer tasks. arXiv preprint arXiv:2303.17491.
- Kirsch, L., Harrison, J., Sohl-Dickstein, J. and Metz, L. (2022). General-purpose in-context learning by meta-learning transformers. arXiv preprint arXiv:2212.04458.

- Lazaric, A., Ghavamzadeh, M. and Munos, R. (2010). Analysis of a classification-based policy iteration algorithm. In *ICML-27th International Conference on Machine Learning*. Omnipress.
- Lee, J. N., Xie, A., Pacchiano, A., Chandak, Y., Finn, C., Nachum, O. and Brunskill, E. (2023). Supervised pretraining can learn in-context reinforcement learning. arXiv preprint arXiv:2306.14892.
- Li, B. Z., Nye, M. and Andreas, J. (2022). Language modeling with latent situations. arXiv preprint arXiv:2212.10012.
- Li, Y., Ildiz, M. E., Papailiopoulos, D. and Oymak, S. (2023). Transformers as algorithms: Generalization and implicit model selection in in-context learning. arXiv preprint arXiv:2301.07067.
- Liang, P., Bommasani, R., Lee, T., Tsipras, D., Soylu, D., Yasunaga, M., Zhang, Y., Narayanan, D., Wu, Y., Kumar, A. et al. (2022). Holistic evaluation of language models. arXiv preprint arXiv:2211.09110.
- Liu, B., Jiang, Y., Zhang, X., Liu, Q., Zhang, S., Biswas, J. and Stone, P. (2023a). LLM+P: Empowering large language models with optimal planning proficiency. arXiv preprint arXiv:2304.11477.
- Liu, Z., Lu, M., Wang, Z., Jordan, M. and Yang, Z. (2022a). Welfare maximization in competitive equilibrium: Reinforcement learning for markov exchange economy. In *International Conference on Machine Learning*. PMLR.
- Liu, Z., Lu, M., Xiong, W., Zhong, H., Hu, H., Zhang, S., Zheng, S., Yang, Z. and Wang, Z. (2023b). Maximize to explore: One objective function fusing estimation, planning, and exploration. In *Advances in Neural Information Processing Systems*, vol. 36.
- Liu, Z., Lu, M., Xiong, W., Zhong, H., Hu, H., Zhang, S., Zheng, S., Yang, Z. and Wang, Z. (2024). Maximize to explore: One objective function fusing estimation, planning, and exploration. *Advances in Neural Information Processing Systems*, **36**.
- Liu, Z., Zhang, Y., Fu, Z., Yang, Z. and Wang, Z. (2022b). Learning from demonstration: Provably efficient adversarial policy imitation with linear function approximation. In *International conference on machine learning*. PMLR.
- Lu, P., Peng, B., Cheng, H., Galley, M., Chang, K.-W., Wu, Y. N., Zhu, S.-C. and Gao, J. (2023). Chameleon: Plug-and-play compositional reasoning with large language models. arXiv preprint arXiv:2304.09842.
- Lu, X. and Van Roy, B. (2019). Information-theoretic confidence bounds for reinforcement learning. Advances in Neural Information Processing Systems.

- Morari, M. and Lee, J. H. (1999). Model predictive control: past, present and future. Computers & chemical engineering, 23 667–682.
- Newton, M. A. and Raftery, A. E. (1994). Approximate bayesian inference with the weighted likelihood bootstrap. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, **56** 3–26.
- Olsson, C., Elhage, N., Nanda, N., Joseph, N., DasSarma, N., Henighan, T., Mann, B., Askell, A., Bai, Y., Chen, A. et al. (2022). In-context learning and induction heads. arXiv preprint arXiv:2209.11895.
- OpenAI (2023). GPT-4 technical report.
- Osband, I., Blundell, C., Pritzel, A. and Van Roy, B. (2016). Deep exploration via bootstrapped dqn. Advances in neural information processing systems, 29.
- Osband, I., Russo, D. and Van Roy, B. (2013). (More) efficient reinforcement learning via posterior sampling. In *Advances in Neural Information Processing Systems*.
- Paul, D., Ismayilzada, M., Peyrard, M., Borges, B., Bosselut, A., West, R. and Faltings, B. (2023). REFINER: Reasoning feedback on intermediate representations. arXiv preprint arXiv:2304.01904.
- Pouplin, T., Sun, H., Holt, S. and Van der Schaar, M. (2024). Retrieval-augmented thought process as sequential decision making. arXiv preprint arXiv:2402.07812.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D. and Sutskever, I. (2019). Language models are unsupervised multitask learners.
- Rawlings, J. B. (2000). Tutorial overview of model predictive control. *IEEE Control Systems Magazine*.
- Razeghi, Y., Logan IV, R. L., Gardner, M. and Singh, S. (2022). Impact of pretraining term frequencies on few-shot reasoning. arXiv preprint arXiv:2202.07206.
- Russo, D. and Van Roy, B. (2014a). Learning to optimize via information-directed sampling. In Advances in Neural Information Processing Systems.
- Russo, D. and Van Roy, B. (2014b). Learning to optimize via posterior sampling. *Mathematics of Operations Research*.
- Russo, D. and Van Roy, B. (2016). An information-theoretic analysis of Thompson sampling. Journal of Machine Learning Research.
- Sekar, R., Rybkin, O., Daniilidis, K., Abbeel, P., Hafner, D. and Pathak, D. (2020). Planning to explore via self-supervised world models. In *International Conference on Machine Learning*. PMLR.

- Sel, B., Al-Tawaha, A., Khattar, V., Wang, L., Jia, R. and Jin, M. (2023). Algorithm of thoughts: Enhancing exploration of ideas in large language models. arXiv preprint arXiv:2308.10379.
- Shen, Y., Song, K., Tan, X., Li, D., Lu, W. and Zhuang, Y. (2023). HuggingGPT: Solving AI tasks with ChatGPT and its friends in HuggingFace. arXiv preprint arXiv:2303.17580.
- Shin, S., Lee, S.-W., Ahn, H., Kim, S., Kim, H., Kim, B., Cho, K., Lee, G., Park, W., Ha, J.-W. et al. (2022). On the effect of pretraining corpora on in-context learning by a large-scale language model. *arXiv* preprint arXiv:2204.13509.
- Shinn, N., Cassano, F., Labash, B., Gopinath, A., Narasimhan, K. and Yao, S. (2023). Reflexion: Language agents with verbal reinforcement learning. arXiv preprint arXiv:2303.11366.
- Shridhar, M., Yuan, X., Côté, M.-A., Bisk, Y., Trischler, A. and Hausknecht, M. (2020). Alfworld: Aligning text and embodied environments for interactive learning. arXiv preprint arXiv:2010.03768.
- Singh, K. (1981). On the asymptotic accuracy of efron's bootstrap. *The Annals of Statistics* 1187–1195.
- Strens, M. (2000). A Bayesian framework for reinforcement learning. In *International Conference on Machine Learning*.
- Sun, H. (2023). Reinforcement learning in the era of llms: What is essential? what is needed? an rl perspective on rlhf, prompting, and beyond.
- Sun, H., Hüyük, A. and van der Schaar, M. (2023a). Query-dependent prompt evaluation and optimization with offline inverse rl. In *The Twelfth International Conference on Learning Representations*.
- Sun, H., Zhuang, Y., Kong, L., Dai, B. and Zhang, C. (2023b). AdaPlanner: Adaptive planning from feedback with language models. arXiv preprint arXiv:2305.16653.
- Sutton, R. S. and Barto, A. G. (2018). Reinforcement learning: An introduction.
- Todd, E., Li, M. L., Sharma, A. S., Mueller, A., Wallace, B. C. and Bau, D. (2023). Function vectors in large language models. *arXiv preprint arXiv:2310.15213*.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F. et al. (2023). LLaMa: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
- Von Oswald, J., Niklasson, E., Randazzo, E., Sacramento, J., Mordvintsev, A., Zhmoginov, A. and Vladymyrov, M. (2023). Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*.

- Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A. and Zhou, D. (2022a). Self-consistency improves chain of thought reasoning in language models. arXiv preprint arXiv:2203.11171.
- Wang, X., Zhu, W. and Wang, W. Y. (2023a). Large language models are implicitly topic models: Explaining and finding good demonstrations for in-context learning. arXiv preprint arXiv:2301.11916.
- Wang, Z., Cai, S., Liu, A., Ma, X. and Liang, Y. (2023b). Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. arXiv preprint arXiv:2302.01560.
- Wang, Z., Pan, T., Zhou, Q. and Wang, J. (2023c). Efficient exploration in resource-restricted reinforcement learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, **37** 10279–10287.
- Wang, Z., Wang, J., Zhou, Q., Li, B. and Li, H. (2022b). Sample-efficient reinforcement learning via conservative model-based actor-critic. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36.
- Wasserman, L. (2000). Bayesian model selection and model averaging. *Journal of mathematical psychology*, **44** 92–107.
- Wei, C.-Y., Jahromi, M. J., Luo, H., Sharma, H. and Jain, R. (2020). Model-free reinforcement learning in infinite-horizon average-reward Markov decision processes. In *International Conference on Machine Learning*.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V. and Zhou, D. (2022). Chain-of-thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems.
- Wies, N., Levine, Y. and Shashua, A. (2023). The learnability of in-context learning. arXiv preprint arXiv:2303.07895.
- Wies, N., Levine, Y. and Shashua, A. (2024). The learnability of in-context learning. Advances in Neural Information Processing Systems, 36.
- Xie, S. M., Raghunathan, A., Liang, P. and Ma, T. (2021). An explanation of in-context learning as implicit Bayesian inference. arXiv preprint arXiv:2111.02080.
- Yang, L. and Wang, M. (2019). Sample-optimal parametric q-learning using linearly additive features. In *International Conference on Machine Learning*.
- Yang, L. and Wang, M. (2020). Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine Learning*.

- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T. L., Cao, Y. and Narasimhan, K. (2023a). Tree of thoughts: Deliberate problem solving with large language models. arXiv preprint arXiv:2305.10601.
- Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K. and Cao, Y. (2022). ReAct: Synergizing reasoning and acting in language models. arXiv preprint arXiv:2210.03629.
- Yao, Y., Li, Z. and Zhao, H. (2023b). Beyond chain-of-thought, effective graph-of-thought reasoning in large language models. arXiv preprint arXiv:2305.16582.
- Zelikman, E., Wu, Y., Mu, J. and Goodman, N. (2022). STaR: Bootstrapping reasoning with reasoning. In *Advances in Neural Information Processing Systems*.
- Zhang, S., Zheng, S., Ke, S., Liu, Z., Jin, W., Yuan, J., Yang, Y., Yang, H. and Wang, Z. (2024). How can llm guide rl? a value-based approach. arXiv preprint arXiv:2402.16181.
- Zhang, Y., Cai, Q., Yang, Z. and Wang, Z. (2020). Generative adversarial imitation learning with neural network parameterization: Global optimality and convergence rate. In *International Conference on Machine Learning*. PMLR.
- Zhang, Y., Liu, B., Cai, Q., Wang, L. and Wang, Z. (2022). An analysis of attention via the lens of exchangeability and latent variable models. arXiv preprint arXiv:2212.14852.
- Zhang, Y., Yang, J., Yuan, Y. and Yao, A. C.-C. (2023a). Cumulative reasoning with large language models. arXiv preprint arXiv:2308.04371.
- Zhang, Y., Zhang, F., Yang, Z. and Wang, Z. (2023b). What and how does in-context learning learn? Bayesian model averaging, parameterization, and generalization. arXiv preprint arXiv:2305.19420.
- Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E. et al. (2023). Judging LLM-as-a-judge with MT-bench and chatbot arena. arXiv preprint arXiv:2306.05685.
- Zheng, S., Wang, L., Qiu, S., Fu, Z., Yang, Z., Szepesvari, C. and Wang, Z. (2022). Optimistic exploration with learned features provably solves markov decision processes with neural dynamics. In *The Eleventh International Conference on Learning Representations*.
- Zhong, H., Xiong, W., Zheng, S., Wang, L., Wang, Z., Yang, Z. and Zhang, T. (2022). Gec: A unified framework for interactive decision making in mdp, pomdp, and beyond. arXiv preprint arXiv:2211.01962.
- Zhou, D., Gu, Q. and Szepesvári, C. (2021a). Nearly minimax optimal reinforcement learning for linear mixture Markov decision processes. In *Conference on Learning Theory*.
- Zhou, D., He, J. and Gu, Q. (2021b). Provably efficient reinforcement learning for discounted MDPs with feature mapping. In *International Conference on Machine Learning*.

# A Notations

We provide a table of notations here.

Notations	Explanation
Ċ	the pair of the next state and the current reward $(s', r)$ ,
$\xi_{(s,a)}$	given the query state-action pair $(s, a)$
$P^{\mathtt{LLM}}(\xi_{(s,a)}   \mathcal{D}, s, a)$	the probability measure of the LLM predicted $\xi_{(s,a)}$
$\Gamma = (\zeta(s,a) \mid \mathcal{D}, S, a)$	given the memory buffer $\mathcal{D}$ as contexts
	the LLM reward estimator
$r_{ ext{LLM}(\mathcal{D})}$	with the memory buffer $\mathcal{D}$ prompted as contexts
D	the LLM transition kernel estimator
$P_{\mathrm{LLM}(\mathcal{D})}$	with the memory buffer $\mathcal{D}$ prompted as contexts
$\mathcal{D}_t$	the history at the t-th step, which includes $\{(s_i, a_i, r_i, s_{i+1})\}_{i=0}^{t-1}$
$\mathbb{P}_0(\theta)$	the prior of $\theta^{\star}$
$\mathbb{P}_{\mathrm{post}}( heta   \mathcal{D})$	the posterior of $\theta^{\star}$ conditioned on $\mathcal{D}$
$\mathbb{P}_t( heta)$	the abbreviation of $\mathbb{P}_{\text{post}}(\theta \mid \mathcal{D}_t)$
$H( heta   \mathcal{D})$	the posterior entropy of the posterior of $\theta$ conditioned on $\mathcal{D}$
$I(\theta; \xi \mid \mathcal{D})$	the information gain of $\xi$ , defined by $H(\theta \mid \mathcal{D}) - H(\theta, \xi \mid \mathcal{D})$
$H_t$	the abbreviation of $H(\theta \mid \mathcal{D}_t)$
$t_k$	the timestep when RAFA switches the policy for the $k$ -th time
$\pi^k$	the abbreviation of $\pi_{t_k}$
$B_{ heta}$	the Bellman operator induced by $\theta$
$B_k$	the Bellman operator induced by $\mathtt{LLM}(\mathcal{D}_{t_k})$
$d_{\mathrm{TV}}(p\ q)$	total variation (TV) between two probability measures $p$ and $q$
$d_{\mathrm{KL}}(p\ q)$	Kullback–Leibler divergence between two probability measures $p$ and $q$
$\mathbb{V}_{x \sim p}[f(x)]$	the variance of $f(X)$ , where X follows the distribution p
(PV)(s,a)	$\mathbb{E}_{s' \sim P(\cdot \mid s, a)}[V(s')]$
L	the bound of $ r + V(s) $ for any $r \in \mathcal{R}$ , $s \in \mathcal{S}$ , and value V
$ u^{\star}(\cdot   s)$	the optimal $\gamma$ -discounted visitation measure starting from state $s$
$\mathbb{N}$	the set of natural numbers
<b>1</b> (x=y)	the indicator with value 1 if $x$ equals $y$ and value 0 otherwise
$\mathbb E$	the expectation
$\mathbb{V}$	the variance

# B More Algorithms

Depending on the specific configuration of the state and action spaces (continuous versus discrete) and the transition and reward models (stochastic versus deterministic), we may choose to emulate the tree-search algorithm, the value iteration algorithm, the random shooting algorithm, or the MCTS algorithm. All of them allow RAFA to achieve provable sample efficiency guarantees as long as they satisfy a specific requirement of optimality (Definition 5.1). For illustration, we emulate the beam-search algorithm (an advanced version of the tree-search algorithm) in Algorithm 7 and the MCTS algorithm in Algorithm 8. For the theoretical discussion, we present the value iteration algorithm in Algorithm 3.

**Algorithm 7** The LLM learner-planner (LLM-LR-PL): A beam-search example (for the deterministic case).

- 1: **input**: The memory buffer  $\mathcal{D}$ , the initial state s, the proposal width L, the search breadth B, and the search depth U.
- 2: **initialization**: Initialize the state array  $S_0 \leftarrow \{s\}$  and the action array  $A_0 \leftarrow \emptyset$ .

  (the learning subroutine)
- 3: Set Model as an LLM instance prompted to use  $\mathcal{D}$  as contexts to generate the next state.
- 4: Set Critic as an LLM instance prompted to use  $\mathcal{D}$  as contexts to estimate the value function.

- 5: Set Elite as an LLM instance prompted to use  $\mathcal{D}$  as contexts to generate multiple candidate actions.
- 6: **for** u = 0, ..., U **do**
- 7: For each current state  $s_u$  in  $S_u$ , invoke Elite to generate L candidate actions.
- 8: For each candidate action  $a_u^{(\ell)}$ , invoke Model to generate the next state  $s_{u+1}^{(\ell)}$  and the received reward  $r_u^{(\ell)}$ .
- 9: For each resulting tuple  $(s_u, a_u^{(\ell)}, s_{u+1}^{(\ell)}, r_u^{(\ell)})$ , invoke Critic to evaluate the expected cumulative future reward  $\widehat{Q}(s_u, a_u^{(\ell)}) \leftarrow r_u^{(\ell)} + \gamma \widehat{V}(s_{u+1}^{(\ell)})$ , where  $\widehat{V}$  is given by Critic.
- 10: Select B best tuples  $(s_u, a_u^{(\ell)}, s_{u+1}^{(\ell)})$  with the highest value  $\widehat{Q}(s_u, a_u^{(\ell)})$  and write them to  $\mathcal{S}_u \times \mathcal{A}_u \times \mathcal{S}_{u+1}$ .
- 11: end for
- 12: For B preserved rollouts in  $S_0 \times A_0 \times \cdots \times S_U \times A_U \times S_{U+1}$ , invoke Critic to evaluate the expected cumulative future reward  $\sum_{u=0}^{U} \gamma^u r_u^{(b)} + \gamma^{U+1} \widehat{V}(s_{U+1}^{(b)})$  and select the best one  $(s_0^{\dagger}, a_0^{\dagger}, \dots, s_U^{\dagger}, a_U^{\dagger}, s_{U+1}^{\dagger})$ , where  $\widehat{V}$  is given by Critic and  $s_0^{\dagger} = s$ .
- 13: **output**: The initial action  $a_0^{\dagger}$  of the selected rollout.

Algorithm 8 LLM learner-planner (LLM-PL) for RAFA: A Monte-Carlo tree-search example (for the stochastic case).

- 1: **input**: The memory buffer  $\mathcal{D}$ , the initial state s, the proposal width L, L', and the expansion budget E.
- 2: **initialization**: Initialize the root node  $n \leftarrow s$  and the child function  $c(\cdot) \leftarrow \emptyset$ .

(the learning subroutine)

- 3: Set Model as an LLM instance prompted to use  $\mathcal{D}$  as contexts to generate the next state.
- 4: Set Critic as an LLM instance prompted to use  $\mathcal{D}$  as contexts to estimate the value function.

(the planning subroutine)

- 5: Set Elite as an LLM instance prompted to use  $\mathcal{D}$  as contexts to generate multiple candidate actions.
- 6: **for** e = 0, ..., E **do**
- 7: Set  $s_e \leftarrow n$ .
- 8: **while**  $s_e$  is not a leaf node, i.e.,  $c(s_e) \neq \emptyset$ , **do**
- 9: Invoke Critic to evaluate the expected cumulative future reward and select the child node  $a_e$  in  $c(s_e)$  with the highest value  $\widehat{Q}(s_e, a_e)$ .
- 10: Set  $s_e$  as a child node in  $c(a_e)$ .
- 11: end while
- 12: For the current state  $s_e$ , invoke Elite to generate L candidate actions.
- 13: Write each candidate action  $a_e^{(\ell)}$  to  $c(s_e)$ , i.e.,  $c(s_e) \leftarrow \{a_e^{(\ell)}\}_{\ell=1}^L$ .
- 14: For each candidate action  $a_e^{(\ell)}$ , invoke Model to sample L' next states.
- 15: Write each next state  $s_e^{(\ell,\ell')}$  to  $c(a_e^{(\ell)})$ , i.e.,  $c(a_e^{(\ell)}) \leftarrow \{s_e^{(\ell,\ell')}\}_{\ell'=1}^{L'}$ .
- 16: For each generated state  $s_e^{(\ell,\ell')}$ , invoke Critic to evaluate the expected cumulative future reward and update the estimated value  $\hat{V}$  for all ancestor nodes. (Optional)
- 17: end for
- 18: Set  $s_0^{\dagger} \leftarrow n$  and  $i \leftarrow 0$ .
- 19: while  $s_i^{\dagger}$  is not a leaf node, i.e.,  $c(s_i^{\dagger}) \neq \emptyset$ , do
- 20: Invoke Critic to evaluate the expected cumulative future reward and select the child node  $a_{i+1}^{\dagger}$  in  $c(s_i^{\dagger})$  with the highest value  $\widehat{Q}(s_i^{\dagger}, a_i^{\dagger})$ .
- 21: Set  $s_{i+1}^{\dagger}$  as a child node in  $c(a_i^{\dagger})$  and  $i \leftarrow i+1$ .
- 22: end while
- 23: **output**: The initial action  $a_0^{\dagger}$  of the selected rollout  $(s_0^{\dagger}, a_0^{\dagger}, \dots, s_i^{\dagger}, a_i^{\dagger})$ .

### C Main Proofs

### C.1 Proof of Proposition 5.2

Proof of Proposition 5.2. We now prove that the value iteration algorithm with a truncated horizon U (Algorithm 3) satisfies the definition of the  $\epsilon$ -optimal planner (Definition 5.1), where U is dependent on  $\epsilon$ . For notational simplicity, we denote  $\max_{s \in \mathcal{S}}$  and  $\max_{a \in \mathcal{A}}$  as  $\max_s$  and  $\max_a$ .

Let

$$\epsilon^{\dagger} = \max_{s,a} |Q^{(1)}(s,a) - r(s,a) - \gamma(PV^{(1)})(s,a)|.$$
 (C.1)

Note that the convergence analysis of the value iteration algorithm in Sutton and Barto (2018) gives

$$\max_{s,a} \left| Q^{(1)}(s,a) - Q^{(2)}(s,a) \right| \le \gamma^{U-2} \max_{s,a} \left| Q^{(U-1)}(s,a) - Q^{(U)}(s,a) \right|,$$

which implies

$$\max_{s,a} |Q^{(1)}(s,a) - Q^{(2)}(s,a)| \le \gamma^{U-2} L. \tag{C.2}$$

We have

$$\epsilon^{\dagger} = \max_{s,a} \left| Q_{\theta}^{(1)}(s,a) - r(s,a) - \gamma(PV^{(2)})(s,a) \right. \\
+ \gamma \mathbb{E}_{s' \sim P(\cdot \mid s,a)} \left[ V^{(1)}(s') - V^{(2)}(s') \right] \left| \right. \\
= \gamma \cdot \max_{s,a} \left| \mathbb{E}_{s' \sim P(\cdot \mid s,a)} \left[ V^{(1)}(s') - V^{(2)}(s') \right] \right| \\
= \gamma \cdot \max_{s,a} \left| \mathbb{E}_{s' \sim P(\cdot \mid s,a)} \left[ \max_{a} Q^{(1)}(s',a) - \max_{a} Q^{(2)}(s',a) \right] \right| \\
\leq \gamma \cdot \max_{s,a} \left| \mathbb{E}_{s' \sim P(\cdot \mid s,a)} \left[ \max_{a} \left| Q^{(1)}(s',a) - Q^{(2)}(s',a) \right| \right] \right| \\
\leq \gamma^{U-1} L, \tag{C.3}$$

where the first and third equalities are based on Algorithm 3, the second last inequality uses the contraction property of the maximum operator, and the last inequality uses (C.2). To let  $\epsilon^{\dagger} < \epsilon$ , it suffices to set  $U \ge 1 + \lceil \log_{\gamma}(\epsilon/L) \rceil$ . Note that the policy  $\pi$  returned by Algorithm 3 satisfies  $\pi(s) = \operatorname{argmax}_{a} Q^{(1)}(s, a)$ . Thus, we prove Proposition 5.2.

# C.2 LLMs with Posterior Alignments Perform BMA

Proof of Proposition 5.4. Recall that  $P_{\text{LLM}(\mathcal{D})}$  and  $r_{\text{LLM}(\mathcal{D})}$  are the estimated transition kernel and reward function induced by  $P^{\text{LLM}}$  that satisfies Assumption 5.3. For any query state-

action pair (s, a) and in-context dataset  $\mathcal{D}$ , it holds that

$$(P_{\text{LLM}(\mathcal{D})}V)(s,a) = \int_{\mathcal{S}} V(s')P_{\text{LLM}(\mathcal{D})}(ds' \mid s, a)$$

$$= \int_{\mathcal{S}} V(s') \left( \int_{\Theta} P_{\theta}(ds' \mid s, a)P_{\text{post}}(d\theta \mid \mathcal{D}) \right)$$

$$= \int_{\Theta} P_{\text{post}}(d\theta \mid \mathcal{D}) \left( \int_{\mathcal{S}} P_{\theta}(ds' \mid s, a)V(s') \right)$$

$$= \mathbb{E}_{\theta \sim \mathbb{P}_{\text{post}}(\cdot \mid \mathcal{D})}[(P_{\theta}V)(s, a)], \tag{C.4}$$

where the second equality uses Assumption 5.3 (Posterior Alignment), the third equality uses Fubin's theorem, and the last equality uses (2.4). For any query state-action pair (s, a) and in-context dataset  $\mathcal{D}$ , it holds that

$$r_{\text{LLM}(\mathcal{D})}(s, a) = \mathbb{E}_{P^{\text{LLM}}}[r \mid \mathcal{D}, s, a]$$

$$= \mathbb{E}_{P_{\text{post}}}[r \mid \mathcal{D}, s, a]$$

$$= \mathbb{E}_{\theta \sim \mathbb{P}_{\text{post}}(\cdot \mid \mathcal{D})}[r_{\theta}(s, a)], \qquad (C.5)$$

where the second equality uses Assumption 5.3 (Posterior Alignment) and the last equality uses (2.4). By the linearity of expectation, we combine (C.4) and (C.5) to obtain

$$\begin{split} r_{\text{LLM}(\mathcal{D})}(s, a) + \gamma \cdot (P_{\text{LLM}(\mathcal{D})}V)(s, a) &= \mathbb{E}_{\theta \sim \mathbb{P}_{\text{post}}(\cdot \mid \mathcal{D})}[r_{\theta}(s, a) + \gamma \cdot (P_{\theta}V)(s, a)] \\ &= \mathbb{E}_{\theta \sim \mathbb{P}_{\text{post}}(\cdot \mid \mathcal{D})}[(B_{\theta}V)(s, a)], \end{split}$$

where the last equality uses the definition of  $B_{\theta}$ . Thus, we finish the proof of Proposition 5.4.

# C.3 Contraction Property of the Posterior Variance

**Proposition C.1** (Contraction Property of the Posterior Variance). Under Assumptions 5.5, the posterior variance in Algorithms 4, 5, and 6 satisfies the following two properties:

(i) 
$$\mathbb{V}_{\theta \sim \mathbb{P}_{t_k}} [(B_{\theta} V_t)(s, a) | \mathcal{D}_{t_k}] \leq 2L^2 \cdot I(\theta; \xi_{(s, a)} | \mathcal{D}_{t_k})$$

(ii) 
$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \mathbb{V}_{\theta \sim \mathbb{P}_{t_k}} \left[ (B_{\theta} V_t)(s_t, a_t) \middle| \mathcal{D}_{t_k} \right] \right] \right] \leq 8\eta L^2 \cdot \mathbb{E}[H_0 - H_T],$$

where we denote the upper bound of the sum of any value function and the reward by a positive constant L, that is,  $|r + V(s)| \le L$  for any reward r, state s, and estimated value function V.

Proof of Proposition C.1. We begin with the proof of the first property in Proposition C.1. Recall the definition that  $\xi_{(s,a)}$  denotes random variables (s',r) in the underlying MDP given the current state s and action a. Define that  $g_t(\xi_{(s,a)}) = (r + V_t(s'))/(2L)$ . Since the sum of any reward and value function is bounded by L, we know that  $|g_t| \leq 1/2$ . for any  $t_k \leq t < t_{k+1}$ , we have

$$2L \cdot \mathbb{E}[g_t(\xi_{(s,a)}) \mid \theta, \mathcal{D}_{t_k}] = (B_{\theta}V_t)(s, a)$$
  

$$2L \cdot \mathbb{E}[g_t(\xi_{(s,a)}) \mid \mathcal{D}_{t_k}] = \mathbb{E}_{\theta \sim \mathbb{P}_t}[(B_{\theta}V_t)(s, a)],$$
(C.6)

for any query state s and action a. By the variational form of total variation (TV) distance  $d_{\text{TV}}$ , we have

$$d_{\text{TV}}^{2}(\mathbb{P}(\xi_{(s,a)} \mid \theta, \mathcal{D}_{t}) \| \mathbb{P}(\xi_{(s,a)} \mid \mathcal{D}_{t})) = \left( \sup_{g:|g| \le 1/2} \mathbb{E}[g(\xi_{(s,a)}) \mid \theta, \mathcal{D}_{t_{k}}] - \mathbb{E}[g(\xi_{(s,a)}) \mid \mathcal{D}_{t_{k}}] \right)^{2}$$

$$\geq \left( \mathbb{E}[g_{t}(\xi_{(s,a)}) \mid \theta, \mathcal{D}_{t_{k}}] - \mathbb{E}[g_{t}(\xi_{(s,a)}) \mid \mathcal{D}_{t_{k}}] \right)^{2}$$

$$= \frac{1}{4L^{2}} \cdot \left( (B_{\theta}V_{t})(s,a) - \mathbb{E}_{\theta \sim \mathbb{P}_{t}} \left[ (B_{\theta}V_{t})(s,a) \right] \right)^{2}, \quad (C.7)$$

where the last equality is the result of (C.6). By taking the expectation with respect to  $\theta \sim \mathbb{P}_t$  on (C.7), we have

$$\mathbb{V}_{\theta \sim \mathbb{P}_{t_k}} \left[ (B_{\theta} V_t)(s, a) \middle| \mathcal{D}_{t_k} \right] = \mathbb{E}_{\theta \sim \mathbb{P}_t} \left[ \left( (B_{\theta} V_t)(s, a) - \mathbb{E}_{\theta \sim \mathbb{P}_t} \left[ (B_{\theta} V_t)(s, a) \right] \right)^2 \right] \\
\leq 4L^2 \cdot \mathbb{E}_{\theta \sim \mathbb{P}_t} \left[ d_{\text{TV}}^2 \left( \mathbb{P}(\xi_{(s, a)} \middle| \theta, \mathcal{D}_{t_k}) \middle| \mathbb{P}(\xi_{(s, a)} \middle| \mathcal{D}_{t_k}) \right) \right] \\
\leq 2L^2 \cdot \mathbb{E}_{\theta \sim \mathbb{P}_t} \left[ d_{\text{KL}} \left( \mathbb{P}(\xi_{(s, a)} \middle| \theta, \mathcal{D}_{t_k}) \middle| \mathbb{P}(\xi_{(s, a)} \middle| \mathcal{D}_{t_k}) \right) \right] \\
= 2L^2 \cdot \left( H(\xi_{(s, a)} \middle| \mathcal{D}_{t_k}) - H(\xi_{(s, a)} \middle| \theta, \mathcal{D}_{t_k}) \right) \\
= 2L^2 \cdot I(\xi_{(s, a)}; \theta \middle| \mathcal{D}_{t_k}) \\
= 2L^2 \cdot I(\theta; \xi_{(s, a)} \middle| \mathcal{D}_{t_k}), \tag{C.8}$$

where the first equality uses the definition of variance, the first inequality uses (C.7) by taking the expectation with respect to  $\theta \sim \mathbb{P}_t$ , and the second inequality uses Pinsker's inequality. Here, the second equality uses the definition of entropy and the second last equality uses the definition of the information gain. Here, the last equality uses the fact that I(X;Y) = I(Y;X) for any two random variables X and Y. Thus, we finish the proof of the first property in Proposition C.1.

Next, we prove the second property in Proposition C.1. By the fact that  $a_t = \pi^k(s_t)$  for any  $t_k \le t < t_{k+1}$ , we have

$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{V}_{\theta \sim \mathbb{P}_{t_{k}}} \left[ (B_{\theta}V_{t})(s_{t}, a_{t}) \middle| \mathcal{D}_{t_{k}} \right] \right] \right] \\
= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{V}_{\theta \sim \mathbb{P}_{t_{k}}} \left[ (B_{\theta}V_{t})(s_{t}, \pi^{k}(s_{t})) \middle| \mathcal{D}_{t_{k}} \right] \right] \right] \\
\leq 2L^{2} \cdot \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} I(\theta; \xi_{(s_{t}, \pi^{k}(s_{t}))} \middle| \mathcal{D}_{t_{k}}) \right] \right],$$

where the inequality invokes (C.8). Under Assumption 5.5 and the same switching condition in Algorithms 4, 5, and 6, we have

$$\mathbb{E}\Big[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \Big[ \sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{V}_{\theta \sim \mathbb{P}_{t_{k}}} \Big[ (B_{\theta}V_{t})(s_{t}, a_{t}) \big| \mathcal{D}_{t_{k}} \Big] \Big] \Big]$$

$$\leq 8\eta L^{2} \cdot \mathbb{E}\Big[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \Big[ \sum_{t=t_{k}}^{t_{k+1}-1} I(\theta; \xi_{(s_{t}, \pi^{k}(s_{t}))} | \mathcal{D}_{t}) \Big] \Big]$$

$$\leq 8\eta L^{2} (H_{0} - H_{T}),$$

where the last inequality uses the chain rule of information gain. Thus, we finish the proof of the second property in Proposition C.1.

### C.4 Proof of Theorem 5.7

Proof of Theorem 5.7. Recall that we denote by  $\pi^k = \pi_{t_k}$  and  $V_t$  is the estimated value function returned by the  $\epsilon$ -optimal planner  $PL^{\epsilon}$  in Algorithm 4. By the definition of the Bayesian regret  $\Re(T)$  and the tower property of the conditional expectation, we have

$$\mathfrak{R}(T) = \mathbb{E}\left[\sum_{k=0}^{K-1} \sum_{t=t_{k}}^{t_{k+1}-1} V_{\theta^{\star}}^{\pi^{\star}}(s_{t}) - V_{\theta^{\star}}^{\pi_{t}}(s_{t})\right] \\
= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{\theta^{\star}}^{\pi^{\star}}(s_{t}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t})\right]\right] \\
= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{t}(s_{t}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t})\right]\right] + \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{\theta^{\star}}^{\pi^{\star}}(s_{t}) - V_{t}(s_{t})\right]\right], \quad (C.9)$$
term (a)

where the first equality uses the fact that  $\pi_t = \pi_{t_k}$  for any  $t_k \leq t < t_{k+1}$  in Algorithm 4. By Definition 5.1, we know that  $V_t(s) = Q_t(s, \pi^k(s))$  for any  $t_k \leq t < t_{k+1}$ . Then, we introduce the following performance difference lemma to bound terms (a) and (b) in (C.9), respectively.

**Lemma C.2** (Performance Difference). For an algorithm ALG with switching times K, estimated value functions  $\{(Q_t, V_t)\}_{t=0}^{T-1}$ , and the corresponding output policy  $\{\pi^k\}_{k=0}^{K-1}$  for T-steps interaction. We assume that ALG switches to the policy  $\pi^k$  at the  $t_k$ -th timestep for the k-th switch and  $V_t(s) = Q_t(s, \pi^k(s))$  for any  $s \in \mathcal{S}$  and k < K. Then, we have two parts of performance difference results for ALG as follows,

• (Part I)

$$(1 - \gamma) \cdot \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{t}(s_{t}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t})\right]\right]$$

$$= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} Q_{t}(s_{t}, a_{t}) - \left(B_{\theta^{\star}}V_{t}\right)(s_{t}, a_{t})\right]\right]$$

$$+ \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\left(V_{t}(s_{t_{k+1}}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t_{k+1}})\right) - \left(V_{t}(s_{t_{k}}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t_{k}})\right)\right]\right], \quad (C.10)$$

$$\text{term (B): value inconsistency}$$

• (Part II)

$$(1 - \gamma) \cdot \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{\theta^{\star}}^{\pi^{\star}}(s_{t}) - V_{t}(s_{t})\right]\right]$$

$$= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{\star}(\cdot \mid s_{t})} \left[ (B_{\theta^{\star}} V_{t})(s, \pi^{\star}(s)) - Q_{t}(s, \pi^{\star}(s))\right] \right]\right]$$
term (A): model prediction error
$$+ \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s' \sim \nu^{\star}(\cdot \mid s_{t})} \left[ (Q_{t}(s, \pi^{\star}(s)) - Q_{t}(s, \pi^{k}(s))\right] \right]\right]. \tag{C.11}$$

where  $\mathbb{E}_{\pi^k}$  is taken with respect to the state-action sequence following  $s_{t+1} \sim P_{\theta^*}(\cdot \mid s_t, a_t)$  and  $a_t = \pi^k(s_t)$  for any  $t_k \leq t < t_{k+1}$ , while  $\mathbb{E}$  is taken with respect to the prior distribution  $\mathbb{P}_0$  of  $\theta^*$ , the iterative update of  $\pi^k$ , and the randomness of  $\{(Q_t, V_t)\}_{t=0}^{T-1}$ . Here, the optimal  $\gamma$ -discounted visitation measure  $\nu^*$  is defined in (5.1).

By the first part of Lemma C.2, we analyze term (a) as follows,

$$(1 - \gamma) \cdot \text{term (a)} = \mathbb{E} \Big[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \Big[ \sum_{t=t_k}^{t_{k+1}-1} Q_t(s_t, a_t) - (B_{\theta^*} V_t)(s_t, a_t) \Big] \Big]$$

$$+ \mathbb{E} \Big[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \Big[ \left( V_t(s_{t_{k+1}}) - V_{\theta^*}^{\pi^k}(s_{t_{k+1}}) \right) - \left( V_t(s_{t_k}) - V_{\theta^*}^{\pi^k}(s_{t_k}) \right) \Big] \Big].$$
(C.12)

Recall that we define  $(B_k V)(s, a) = r_{\text{LLM}(\mathcal{D}_{t_k})}(s, a) + (P_{\text{LLM}(\mathcal{D}_{t_k})} V)(s, a)$  for any (s, a) and value V. By the definition of  $\epsilon$ -optimal planner (Definition 5.1) and the planning procedure

 $(\pi_t, V_t) \leftarrow \mathtt{PL}^{\epsilon}(P_{\mathtt{LLM}(\mathcal{D}_{t_k})}, r_{\mathtt{LLM}(\mathcal{D}_{t_k})})$  in Line 6 of Algorithm 4, we have

$$|Q_t(s_t, a_t) - (B_{\theta^*} V_t)(s_t, a_t)| \le \epsilon + ((B_k - B_{\theta^*}) V_t)(s_t, a_t)$$

$$= \epsilon + |((B_k - B_{\theta^*}) V_t)(s_t, a_t)|$$
(C.13)

for any  $t_k \leq t < t_{k+1}$ . Then, we plug (C.13) into (C.12) to obtain

$$(1 - \gamma) \cdot \text{term (a)} \leq \underbrace{\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \left| \left( (B_{k} - B_{\theta^{\star}}) V_{t} \right) (s_{t}, a_{t}) \right| \right] \right] + \epsilon \cdot T}_{\text{term (a1)}} + \underbrace{\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \left( V_{t}(s_{t_{k+1}}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t_{k+1}}) \right) - \left( V_{t}(s_{t_{k}}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t_{k}}) \right) \right] \right]}_{\text{term (a2)}}.$$

$$(C.14)$$

Under Assumption 5.3 and 5.5, we have

term (a1) = 
$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \mathbb{E}_{\theta \sim \mathbb{P}_{t_k}} \left[ \left| \left( (B_k - B_\theta) V_t \right) (s_t, a_t) \right| \right| \mathcal{D}_{t_k} \right] \right] \right]$$
  
 $\leq \sqrt{T} \cdot \left( \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \mathbb{E}_{\theta \sim \mathbb{P}_{t_k}} \left[ \left| \left( (B_k - B_\theta) V_t \right) (s_t, a_t) \right|^2 \right| \mathcal{D}_{t_k} \right] \right] \right] \right)^{1/2}$   
 $= \sqrt{T} \cdot \left( \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \mathbb{V}_{\theta \sim \mathbb{P}_{t_k}} \left[ (B_k V_t) (s_t, a_t) \right| \mathcal{D}_{t_k} \right] \right] \right)^{1/2},$  (C.15)

where the first equality uses the tower property of the conditional expectation and the definition that the posterior distribution of  $\theta^*$  given  $\mathcal{D}_{t_k}$  is  $\mathbb{P}_{t_k}$ . Here, the first inequality uses Cauchy-Schwarz inequality and the last equality invokes Proposition 5.4 and the definition of variance. Under Assumption 5.5, we apply the second property in Proposition C.1 on the right-hand side of (C.15) to have

term (a1) 
$$\leq 2\sqrt{2\eta}L \cdot \sqrt{\mathbb{E}[H_0 - H_T]} \cdot \sqrt{T}$$
. (C.16)

Since any value function is bounded by L, we bound term (a2) in (C.14) as follows,

term (a2) 
$$\leq 4L \cdot \mathbb{E}[K]$$
. (C.17)

To characterize the upper bound of the switching times K, we introduce the following lemma.

**Lemma C.3** (Upper bound of Switching Times). If  $H_{t_k} - H_{t_{k+1}} \ge \log 2$  for any k < K, then it holds that

$$K - 1 \le (H_0 - H_{t_{K-1}})/\log 2 \le (H_0 - H_T)/\log 2.$$

*Proof of Lemma C.3.* Since  $H_{t_k} - H_{t_{k+1}} \ge \log 2$ , we have

$$H_0 - H_{t_{K-1}} = \sum_{k=0}^{K-2} H_{t_k} - H_{t_{k+1}} \ge (K-1) \cdot \log 2,$$

which implies

$$K - 1 \le (H_0 - H_{t_{K-1}})/\log 2 \le (H_0 - H_T)/\log 2$$
.

Thus, we finish the proof of Lemma C.3.

As the switching condition in Algorithm 6 implies  $H_{t_k} - H_{t_{k+1}} \ge \log 2$ , we apply Lemma C.3 to have

$$\mathbb{E}[K] \le 1 + \mathbb{E}[H_0 - H_{t_{K-1}}]/\log 2 \le (H_0 - H_T)/\log 2. \tag{C.18}$$

Combining (C.17) and (C.18), we upper bound term (a2) in (C.14) as

term (a2) 
$$\leq 4L \cdot \mathbb{E}[K] \leq 4L + \frac{4L \cdot \mathbb{E}[H_0 - H_T]}{\log 2}$$
. (C.19)

Plugging (C.16) and (C.19) into (C.14), we have

$$(1 - \gamma) \cdot \text{term (a)} \le 2\sqrt{2\eta}L \cdot \sqrt{\mathbb{E}[H_0 - H_T]} \cdot \sqrt{T} + \epsilon \cdot T + 4L + \frac{4L \cdot \mathbb{E}[H_0 - H_T]}{\log 2}.$$
 (C.20)

Then, we invoke the second part of Lemma C.2 to decompose term (b) in (C.9) as follows,

$$(1 - \gamma) \cdot \text{term (b)} = \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ \left( B_{\theta^{*}} V_{t} \right) (s, \pi^{*}(s)) - Q_{t}(s, \pi^{*}(s)) \right] \right] \right]$$

$$+ \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ Q_{t}(s, \pi^{*}(s)) - Q_{t}(s, \pi^{k}(s)) \right] \right] \right]$$

$$\leq \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ \left| \left( B_{\theta^{*}} - B_{k} \right) V_{t} \right) (s, \pi^{*}(s)) \right| \right] \right] \right] + \epsilon \cdot T$$

$$+ \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ Q_{t}(s, \pi^{*}(s)) - Q_{t}(s, \pi^{k}(s)) \right] \right] \right], \quad (C.21)$$

$$+ \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ Q_{t}(s, \pi^{*}(s)) - Q_{t}(s, \pi^{k}(s)) \right] \right] \right], \quad (C.21)$$

where the inequality uses (C.13). For term (b1) in (C.21), we use the tower property of the

conditional expectation to obtain

term (b1) = 
$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \mathbb{E}_{\theta^* \sim \mathbb{P}_{t_k}} \left[\frac{\mathbb{E}_{s \sim \nu^*(\cdot \mid s_t)} \left[\left|\left((B_k - B_{\theta^*})V_t\right)(s, \pi^*(s))\right|\right]}{\left((B_k - B_{\theta^*})V_t\right)(s_t, \pi^k(s_t))}\right]\right]$$
 (C.22)  

$$\cdot \left(\left(B_k - B_{\theta^*}\right)V_t\right)(s_t, \pi^k(s_t)) \left|\mathcal{D}_{t_k}\right]\right]\right]$$

$$= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \mathbb{E}_{\theta^* \sim \mathbb{P}_{t_k}} \left[G_{t,k}(\theta^*) \cdot \left((B_k - B_{\theta^*})V_t\right)(s_t, \pi^k(s_t))\right| \mathcal{D}_{t_k}\right]\right]\right], \quad (C.23)$$

where we define

$$G_{t,k}(\theta^*) = \frac{\mathbb{E}_{s \sim \nu^*(\cdot \mid s_t)} \left[ \left| \left( (B_k - B_{\theta^*}) V_t \right) (s, \pi^*(s)) \right| \right]}{\left( (B_k - B_{\theta^*}) V_t \right) (s_t, \pi^k(s_t))}$$

and  $\mathbb{E}_{\theta^* \mid \mathcal{D}_{t_k}}[\cdot] = \mathbb{E}_{\theta^* \sim \mathbb{P}_{t_k}}[\cdot \mid \mathcal{D}_{t_k}]$  for notational simplicity.

Applying Cauchy Schwartz inequality on the left-hand side of (C.23) several times, we have

term (b1) 
$$\leq \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} \left( \mathbb{E}_{\theta^{*} \mid \mathcal{D}_{t_{k}}} \left[ G_{t,k}^{2}(\theta^{*}) \right] \right)^{1/2} \right] \right]$$

$$\cdot \left( \mathbb{E}_{\theta^{*} \mid \mathcal{D}_{t_{k}}} \left[ \left| \left( (B_{k} - B_{\theta^{*}}) V_{t} \right) (s_{t}, \pi^{k}(s_{t})) \right|^{2} \right] \right]^{1/2} \right]$$

$$\leq \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{\theta^{*} \mid \mathcal{D}_{t_{k}}} \left[ G_{t,k}^{2}(\theta^{*}) \right] \right]^{1/2} \right]$$

$$\cdot \left( \sum_{t=t_{k}}^{K-1} \mathbb{E}_{\theta^{*} \mid \mathcal{D}_{t_{k}}} \left[ \left| \left( (B_{k} - B_{\theta^{*}}) V_{t} \right) (s_{t}, \pi^{k}(s_{t})) \right|^{2} \right] \right]^{1/2} \right]$$

$$\leq \left( \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{\theta^{*} \mid \mathcal{D}_{t_{k}}} \left[ G_{t,k}^{2}(\theta^{*}) \right] \right] \right]^{1/2}$$

$$\cdot \left( \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{\theta^{*} \mid \mathcal{D}_{t_{k}}} \left[ \left| \left( (B_{k} - B_{\theta^{*}}) V_{t} \right) (s_{t}, \pi^{k}(s_{t})) \right|^{2} \right] \right] \right]^{1/2}$$

$$\leq \sqrt{\kappa^{2} \cdot T} \cdot \left( \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{\theta^{*} \mid \mathcal{D}_{t_{k}}} \left[ \left| \left( (B_{k} - B_{\theta^{*}}) V_{t} \right) (s_{t}, \pi^{k}(s_{t})) \right|^{2} \right] \right] \right]^{1/2} ,$$

where the first three inequalities are all based on Cauchy Schwartz inequality and the last inequality uses the definition of  $\kappa$  in Assumption 5.6. Under Assumptions 5.3 and 5.5, we have

term (b1) 
$$\leq \sqrt{\kappa^2 \cdot T} \cdot \left( \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[ \sum_{t=t_k}^{t_{k+1}-1} \mathbb{V}_{\theta \sim \mathbb{P}_{t_k}} \left[ (B_{\theta} V_t)(s_t, a_t) \middle| \mathcal{D}_{t_k} \right] \right] \right)^{1/2}$$

$$\leq 2\sqrt{2\eta} L \kappa \cdot \sqrt{\mathbb{E}[H_0 - H_T]} \cdot \sqrt{T}, \tag{C.26}$$

where the first inequality invokes Proposition 5.4 and the definition of variance. Here, the second inequality invokes Proposition C.1. By the definition of  $\epsilon$ -optimal planner (Definition 5.1), we know term (b2) in (C.21) is non-positive. Then, plugging (C.26) into (C.21), we have

$$(1 - \gamma) \cdot \text{term (b)} \le 2\sqrt{2\eta}L\kappa \cdot \sqrt{\mathbb{E}[H_0 - H_T]} \cdot \sqrt{T} + \epsilon \cdot T.$$
 (C.27)

Combining (C.9), (C.20), and (C.27), we have

$$\Re(T) = \frac{1}{1 - \gamma} \cdot \left( \text{term (a) + term (b)} \right)$$

$$\leq \frac{2\sqrt{2}(\kappa + 1)L \cdot \sqrt{\mathbb{E}[H_0 - H_T]}}{1 - \gamma} \cdot \sqrt{T} + \frac{2\epsilon}{1 - \gamma} \cdot T + \frac{4L}{1 - \gamma} + \frac{4L \cdot \mathbb{E}[H_0 - H_T]}{(1 - \gamma)\log 2}$$

$$= \mathcal{O}\left( \frac{(\kappa + 1)L \cdot \sqrt{\mathbb{E}[H_0 - H_T]}}{1 - \gamma} \cdot \sqrt{T} + \frac{\epsilon}{1 - \gamma} \cdot T + \frac{L \cdot \mathbb{E}[H_0 - H_T]}{1 - \gamma} \right).$$

Thus, we finish the proof of Theorem 5.7.

### C.5 Proof of Theorem 5.8

Proof of Theorem 5.8. Recall that we denote by  $\pi^k = \pi_{t_k}$ , and  $V_t$  is the estimated value function returned by the  $\epsilon$ -optimal planner  $PL^{\epsilon}$  in Algorithm 4. By the definition of the Bayesian regret  $\Re(T)$  and the tower property of the conditional expectation, we have

$$\mathfrak{R}(T) = \mathbb{E}\left[\sum_{k=0}^{K-1} \sum_{t=t_{k}}^{t_{k+1}-1} V_{\theta^{\star}}^{\pi^{\star}}(s_{t}) - V_{\theta^{\star}}^{\pi_{t}}(s_{t})\right] \\
= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{\theta^{\star}}^{\pi^{\star}}(s_{t}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t})\right]\right] \\
= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{t}(s_{t}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t})\right]\right] + \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{\theta^{\star}}^{\pi^{\star}}(s_{t}) - V_{t}(s_{t})\right]\right], \quad (C.28)$$
term (a)

where the first equality uses the fact that  $\pi_t = \pi_{t_k}$  for any  $t_k \leq t < t_{k+1}$  in Algorithm 4. By Definition 5.1, we know that  $V_t(s) = Q_t(s, \pi^k(s))$  for any  $t_k \leq t < t_{k+1}$ . Then, we apply the first part of Lemma C.2 to analyze term (a) as follows,

$$(1 - \gamma) \cdot \text{term (a)} = \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} Q_t(s_t, a_t) - \left(B_{\theta^*} V_t\right)(s_t, a_t)\right]\right] + \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\left(V_t(s_{t_{k+1}}) - V_{\theta^*}^{\pi^k}(s_{t_{k+1}})\right) - \left(V_t(s_{t_k}) - V_{\theta^*}^{\pi^k}(s_{t_k})\right)\right]\right].$$

Recall that we define  $(B_k V)(s, a) = r_{\text{LLM}(\mathcal{D}_{t_k})}(s, a) + (P_{\text{LLM}(\mathcal{D}_{t_k})} V)(s, a)$  for any (s, a) and value function V. By the definition of  $\epsilon$ -optimal planner (Definition 5.1) and the planning procedure  $(\pi_t, V_t) \leftarrow \text{PL}^{\epsilon}(P_{\text{LLM}(\mathcal{D}_{t_k})}, r_{\text{LLM}(\mathcal{D}_{t_k})} + \Gamma_k)$  in Algorithm 5, we have

$$|Q_t(s,a) - (B_{\theta^*}V_t)(s_t,a_t)| \le \epsilon + ((B_k - B_{\theta^*})V_t)(s,a) + \Gamma_k(s,a)$$
(C.29)

for any  $t_k \leq t < t_{k+1}$  and any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . Then, we plug (C.29) into (C.14) to obtain

$$(1 - \gamma) \cdot \text{term (a)} \leq \underbrace{\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \left( (B_{k} - B_{\theta^{\star}}) V_{t} \right) (s_{t}, a_{t}) \right] \right]}_{\text{term (a1)}} + \epsilon \cdot T + \underbrace{\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \Gamma_{k}(s_{t}, a_{t}) \right] \right]}_{\text{term (a2)}} + \underbrace{\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \left( V_{t}(s_{t_{k+1}}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t_{k+1}}) \right) - \left( V_{t}(s_{t_{k}}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t_{k}}) \right) \right] \right]}_{\text{term (a3)}}$$

$$(C.31)$$

Under Assumption 5.3, we have

term (a1) = 
$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \mathbb{E}_{\theta \sim \mathbb{P}_{t_k}} \left[ \left( (B_k - B_\theta) V_t \right) (s_t, a_t) \middle| \mathcal{D}_{t_k} \right] \right] \right]$$
  
 $\leq \sqrt{T} \cdot \left( \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \mathbb{E}_{\theta \sim \mathbb{P}_{t_k}} \left[ \left| \left( (B_k - B_\theta) V_t \right) (s_t, a_t) \middle|^2 \middle| \mathcal{D}_{t_k} \right] \right] \right] \right)^{1/2}$   
 $= \sqrt{T} \cdot \left( \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \mathbb{V}_{\theta \sim \mathbb{P}_{t_k}} \left[ (B_\theta V_t) (s_t, a_t) \middle| \mathcal{D}_{t_k} \right] \right] \right)^{1/2},$  (C.32)

where the first equality uses the tower property of the conditional expectation and the definition that the posterior distribution of  $\theta^*$  given  $\mathcal{D}_{t_k}$  is  $\mathbb{P}_{t_k}$ . Here, the first inequality uses Cauchy-Schwarz inequality and the last equality invokes Proposition 5.4 and the definition of variance. Under Assumption 5.5, we apply the second property in Proposition C.1 on the right-hand side of (C.32) to have

term (a1) 
$$\leq 2\sqrt{2\eta}L \cdot \sqrt{\mathbb{E}[H_0 - H_T]} \cdot \sqrt{T}$$
. (C.33)

Recall that the bonus  $\Gamma_k$  used in Algorithm 5 is defined by  $\Gamma_k(s,a) = \sqrt{2}L \cdot \sqrt{I(\theta; \xi_{(s,a)} \mid \mathcal{D}_{t_k})}$ 

. For term (a2) in (C.31), we have

term (a2) 
$$\leq \sqrt{T} \cdot \left( \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} \Gamma_{k}^{2}(s_{t}, a_{t}) \right] \right] \right)^{1/2}$$

$$= \sqrt{2}L \cdot \sqrt{T} \cdot \left( \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} I(\theta; \xi_{(s_{t}, a_{t})} \mid \mathcal{D}_{t_{k}}) \right] \right] \right)^{1/2}$$

$$\leq 2\sqrt{2\eta}L \cdot \sqrt{T} \cdot \left( \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \sum_{t=t_{k}}^{t_{k+1}-1} I(\theta; \xi_{(s_{t}, a_{t})} \mid \mathcal{D}_{t}) \right] \right] \right)^{1/2}, \quad (C.34)$$

where the equality uses the definition of  $\Gamma_k$  in Algorithm 5. Here, the last inequality invokes Assumption 5.5 and the switching condition in Algorithm 5.

As  $a_t = \pi^k(s_t)$ , we further bound the right-hand side of (C.34) as follows,

term (a2) 
$$\leq 2\sqrt{2\eta}L \cdot \sqrt{T} \cdot \left(\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k}\left[\sum_{t=t_k}^{t_{k+1}-1} I(\theta; \xi_{(s_t, \pi^k(s_t))} \mid \mathcal{D}_t)\right]\right]\right)^{1/2}$$
  
=  $2\sqrt{2\eta}L \cdot \sqrt{\mathbb{E}[H_0 - H_T]} \cdot \sqrt{T}$ , (C.35)

where the last inequality uses the chain rule of the information gain. Using the fact that any value function is bounded by L, we bound term (a3) in (C.14) as

term (a3) 
$$\leq 4L \cdot \mathbb{E}[K]$$
.

As the switching condition in Algorithm 5 implies  $H_{t_k} - H_{t_{k+1}} \ge \log 2$ , we apply Lemma C.3 to have

term (a3) 
$$\leq 4L + \frac{4L \cdot \mathbb{E}[H_0 - H_T]}{\log 2}$$
. (C.36)

Plugging (C.33), (C.35) and (C.36) into (C.14), we have

$$(1 - \gamma) \cdot \text{term (a)} \le 4\sqrt{2}L \cdot \sqrt{\mathbb{E}[H_0 - H_T]} \cdot \sqrt{T} + \epsilon \cdot T + 4L + \frac{4L \cdot \mathbb{E}[H_0 - H_T]}{\log 2}.$$
 (C.37)

Then, we invoke the second part of Lemma C.2 to decompose term (b) in (C.28) as follows,

$$(1 - \gamma) \cdot \text{term (b)} = \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ \left(B_{\theta^{*}} V_{t}\right)(s, \pi^{*}(s)) - Q_{t}(s, \pi^{*}(s)) \right] \right] \right]$$

$$+ \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ Q_{t}(s, \pi^{*}(s)) - Q_{t}(s, \pi^{k}(s)) \right] \right] \right]$$

$$\leq \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ \left( \left(B_{\theta^{*}} - B_{k}\right) V_{t}\right)(s, \pi^{*}(s)) - \Gamma_{k}(s, \pi^{*}(s)) \right] \right] \right]$$

$$+ \epsilon \cdot T + \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ Q_{t}(s, \pi^{*}(s)) - Q_{t}(s, \pi^{k}(s)) \right] \right] \right]$$

$$+ \epsilon \cdot T + \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ Q_{t}(s, \pi^{*}(s)) - Q_{t}(s, \pi^{k}(s)) \right] \right] \right]$$

$$+ \epsilon \cdot T + \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ Q_{t}(s, \pi^{*}(s)) - Q_{t}(s, \pi^{k}(s)) \right] \right] \right]$$

$$+ \epsilon \cdot T + \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ Q_{t}(s, \pi^{*}(s)) - Q_{t}(s, \pi^{k}(s)) \right] \right] \right]$$

$$+ \epsilon \cdot T + \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ Q_{t}(s, \pi^{*}(s)) - Q_{t}(s, \pi^{*}(s)) \right] \right] \right]$$

$$+ \epsilon \cdot T + \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s_{t})} \left[ Q_{t}(s, \pi^{*}(s)) - Q_{t}(s, \pi^{*}(s)) \right] \right] \right]$$

where the inequality uses (C.29). Under Assumption 5.3, we bound term (b1) in (C.38) as follows,

term (b1) = 
$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{\star}(\cdot \mid s_{t})} \left[ \mathbb{E}_{\theta \sim \mathbb{P}_{t_{k}}} \left[ \left( (B_{k} - B_{\theta}) V_{t} \right) (s, \pi^{\star}(s)) \middle| \mathcal{D}_{t_{k}} \right] - \Gamma_{k}(s, \pi^{\star}(s)) \right] \right]$$

$$\leq \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{\star}(\cdot \mid s_{t})} \left[ \sqrt{\mathbb{E}_{\theta \sim \mathbb{P}_{t_{k}}} \left[ \left| \left( (B_{k} - B_{\theta}) V_{t} \right) (s, \pi^{\star}(s)) \middle|^{2} \middle| \mathcal{D}_{t_{k}} \right] \right] \right]$$

$$- \Gamma_{k}(s, \pi^{\star}(s)) \right] \right]$$

$$= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{\star}(\cdot \mid s_{t})} \left[ \sqrt{\mathbb{V}_{\theta \sim \mathbb{P}_{t_{k}}} \left[ (B_{\theta} V_{t}) (s, \pi^{\star}(s)) \middle| \mathcal{D}_{t_{k}} \right] - \Gamma_{k}(s, \pi^{\star}(s)) \right] \right] \right],$$
(C.39)

where the first equality uses the tower property of the conditional expectation and the definition that the posterior distribution of  $\theta^*$  given  $\mathcal{D}_{t_k}$  is  $\mathbb{P}_{t_k}$ . Here, the first inequality uses Cauchy-Schwarz inequality and the last equality invokes Proposition 5.4 and the definition of variance. Under Assumption 5.5, we apply the first property in Proposition C.1 to have

$$\sqrt{\mathbb{V}_{\theta \sim \mathbb{P}_{t_k}} \left[ (B_{\theta} V_t)(s, \pi^{\star}(s)) \middle| \mathcal{D}_{t_k} \right]} - \Gamma_k(s, \pi^{\star}(s)) \leq \sqrt{2} L \cdot \sqrt{I(\theta; \xi_{(s, \pi^{\star}(s))})} - \Gamma_k(s, \pi^{\star}(s))$$

$$= \sqrt{2} L \cdot \sqrt{I(\theta; \xi_{(s, \pi^{\star}(s))})} - \sqrt{2} L \cdot \sqrt{I(\theta; \xi_{(s, \pi^{\star}(s))})}$$

$$\leq 0, \qquad (C.40)$$

for any  $t_k \leq t < t_{k+1}$ , k < K, and state  $s \in \mathcal{S}$ . Here, the equality uses the definition of  $\Gamma_k$  in Algorithm 5. Plugging (C.40) into (C.39), we have

$$term (b1) \le 0. \tag{C.41}$$

By the definition of  $\epsilon$ -optimal planner (Definition 5.1), we know term (b2) in (C.21) is non-positive. Then, plugging (C.38) into (C.38), we have

$$(1 - \gamma) \cdot \text{term (b)} \le \epsilon \cdot T.$$
 (C.42)

Combining (C.28), (C.37), and (C.42), we have

$$\mathfrak{R}(T) = \frac{1}{1 - \gamma} \cdot \left( \text{term (a)} + \text{term (b)} \right)$$

$$\leq \frac{4\sqrt{2}L \cdot \sqrt{\mathbb{E}[H_0 - H_T]}}{1 - \gamma} \cdot \sqrt{T} + \frac{2\epsilon}{1 - \gamma} \cdot T + \frac{4L}{1 - \gamma} + \frac{4L \cdot \mathbb{E}[H_0 - H_T]}{(1 - \gamma)\log 2}$$

$$= \mathcal{O}\left(\frac{L \cdot \sqrt{\mathbb{E}[H_0 - H_T]}}{1 - \gamma} \cdot \sqrt{T} + \frac{\epsilon}{1 - \gamma} \cdot T + \frac{L \cdot \mathbb{E}[H_0 - H_T]}{1 - \gamma}\right).$$

Thus, we finish the proof of Theorem 5.8.

### C.6 Proof of Theorem 5.10

Proof of Theorem 5.10. For notational simplicity, we denote by  $\theta^k$  the corresponding parameter for the mechanism LLM+PS given  $\mathcal{D}_{t_k}$  in Algorithm 6, which satisfies

$$(B_{\theta^k}V)(s,a) = r_{\text{LLM+PS}(\mathcal{D}_{t_k})}(s,a) + \gamma \cdot (P_{\text{LLM+PS}(\mathcal{D}_{t_k})}V)(s,a), \tag{C.43}$$

for any k < K,  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , and value function V. Recall the definition of optimal value  $V_{\theta}^{\star}$  given the parameter  $\theta$  in (2.2). Under Assumption 5.9, we know that  $(B_{\theta^k}V)(s, a) \mid \mathcal{D}_{t_k}$  and  $(B_{\theta^k}V)(s, a) \mid \mathcal{D}_{t_k}$  follows the same distribution for any k < K,  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , and value function V. By Bellman optimality equation in (2.2), we have that  $V_{\theta^k}^{\star}(s, a) \mid \mathcal{D}_{t_k}$  and  $V_{\theta^k}^{\star}(s, a) \mid \mathcal{D}_{t_k}$  follows the same distribution for any k < K,  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , and value function V. Recall that we denote by  $\pi^k = \pi_{t_k}$  and  $V_t$  is the estimated value function returned by the  $\epsilon$ -optimal planner  $PL^{\epsilon}$  in Algorithm 6. By the definition of the Bayesian regret  $\mathfrak{R}(T)$  and  $\pi_t = \pi^k$  for any  $t_k \leq t < t_{k+1}$  and k < K, we have

$$\mathfrak{R}(T) = \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{\theta^{\star}}^{\pi^{\star}}(s_{t}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t})\right]\right]$$

$$= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{\theta^{\star}}^{\star}(s_{t}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t})\right]\right]$$

$$= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{\theta^{\star} \sim \mathbb{P}_{t_{k}}} \left[V_{\theta^{\star}}^{\star}(s_{t}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t}) \mid \mathcal{D}_{t_{k}}\right]\right]\right], \qquad (C.44)$$

where the second equality uses the definition of optimal policy and (2.2). Here, the last equality uses the tower property of the conditional expectation. Recall that  $\theta^* \mid \mathcal{D}_{t_k}$  and

 $\theta^k \mid \mathcal{D}_{t_k}$  follows the same distribution  $\mathbb{P}_{t_k}$  for any  $t_k \leq t < t_{k+1}$ , which implies

$$\mathbb{E}_{\theta^{\star} \sim \mathbb{P}_{t_{k}}} \left[ V_{\theta^{\star}}^{\star}(s_{t}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t}) \,|\, \mathcal{D}_{t_{k}} \right] = \mathbb{E}_{\theta^{\star} \sim \mathbb{P}_{t_{k}}} \left[ V_{\theta^{\star}}^{\star}(s_{t}) \,|\, \mathcal{D}_{t_{k}} \right] - \mathbb{E}_{\theta^{\star} \sim \mathbb{P}_{t_{k}}} \left[ V_{\theta^{\star}}^{\pi^{k}}(s_{t}) \,|\, \mathcal{D}_{t_{k}} \right] \\
= \mathbb{E}_{\theta^{k} \sim \mathbb{P}_{t_{k}}} \left[ V_{\theta^{k}}^{\star}(s_{t}) \,|\, \mathcal{D}_{t_{k}} \right] - \mathbb{E}_{\theta^{\star} \sim \mathbb{P}_{t_{k}}} \left[ V_{\theta^{\star}}^{\pi^{k}}(s_{t}) \,|\, \mathcal{D}_{t_{k}} \right] \\
= \mathbb{E}_{\theta^{\star}, \theta^{k} \sim \mathbb{P}_{t_{k}}} \left[ V_{\theta^{k}}^{\star}(s_{t}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t}) \,|\, \mathcal{D}_{t_{k}} \right], \tag{C.45}$$

where the first and the second inequalities use the linear property of the conditional expectation. Plugging (C.45) into (C.44), we have

$$\mathfrak{R}(T) = \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{\theta^{\star},\theta^{k} \sim \mathbb{P}_{t_{k}}} \left[V_{\theta^{k}}^{\star}(s_{t}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t}) \mid \mathcal{D}_{t_{k}}\right]\right]\right]$$

$$= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{\theta^{k}}^{\star}(s_{t}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t})\right]\right], \tag{C.46}$$

where the last equality uses the tower property of the conditional expectation.

Meanwhile, by Definition 5.1, we have

$$\max_{s \in \mathcal{S}} |V_{\theta^{k}}^{\star}(s_{t}) - V_{t}(s_{t})| = \max_{s \in \mathcal{S}} |\max_{a \in \mathcal{A}} Q_{\theta^{k}}^{\star}(s, a) - \max_{a} Q_{t}(s, a)| 
\leq \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} |Q_{\theta^{k}}^{\star}(s, a) - Q_{t}(s, a)| 
= \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} |((B_{\theta^{k}} V_{\theta^{k}}^{\star})(s, a) - (B_{\theta^{k}} V_{t})(s, a)) + ((B_{\theta^{k}} V_{t})(s, a) - Q_{t}(s, a))| 
\leq \gamma \cdot \max_{s \in \mathcal{S}} |V_{\theta^{k}}^{\star}(s_{t}) - V_{t}(s_{t})| + \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} |(B_{\theta^{k}} V_{t})(s, a) - Q_{t}(s, a)|,$$
(C.47)

where the equality and the second equality uses the definitions of  $(Q_{\theta}^{\star}, V_{\theta}^{\star})$  in (2.2). Here, the first inequality uses the fact that the maximum operator is a contraction map, and the last inequality uses triangle inequality and (C.43). Rearranging (C.47), we have

$$\max_{s \in \mathcal{S}} |V_{\theta^k}^{\star}(s_t) - V_t(s_t)| \leq \frac{1}{1 - \gamma} \cdot \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} |(B_{\theta^k} V_t)(s, a) - Q_t(s, a)| \\
\leq \frac{\epsilon}{1 - \gamma}, \tag{C.48}$$

where the last inequality uses the definition of  $\epsilon$ -optimal planner (Definition 5.1), the planning procedure  $(\pi_t, V_t) \leftarrow \text{PL}^{\epsilon}(P_{\text{LLM+PS}(\mathcal{D}_{t_k})}, r_{\text{LLM+PS}(\mathcal{D}_{t_k})})$  in Algorithm 6 and the definition of  $\theta^k$  in (C.43). Then, we upper bound the right-hand side of (C.46) as

$$\Re(T) \le \frac{\epsilon}{1 - \gamma} \cdot T + \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} V_t(s_t) - V_{\theta^*}^{\pi^k}(s_t)\right]\right]. \tag{C.49}$$

By the first part of Lemma C.2, we analyze term (a) in (C.49) as follows,

$$(1 - \gamma) \cdot \text{term (a)} = \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} Q_t(s_t, a_t) - \left(B_{\theta^*} V_t\right)(s_t, a_t)\right]\right] + \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\left(V_t(s_{t_{k+1}}) - V_{\theta^*}^{\pi^k}(s_{t_{k+1}})\right) - \left(V_t(s_{t_k}) - V_{\theta^*}^{\pi^k}(s_{t_k})\right)\right]\right].$$
(C.50)

By the last inequality in (C.48), we have

$$|Q_t(s_t, a_t) - (B_{\theta^*} V_t)(s_t, a_t)| \le \epsilon + ((B_k - B_{\theta^*}) V_t)(s_t, a_t)$$

$$= \epsilon + |((B_k - B_{\theta^*}) V_t)(s_t, a_t)|$$
(C.51)

for any  $t_k \leq t < t_{k+1}$ . Then, we plug (C.51) into (C.50) to obtain

$$(1 - \gamma) \cdot \text{term (a)} \leq \underbrace{\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \left| \left( (B_{k} - B_{\theta^{\star}}) V_{t} \right) (s_{t}, a_{t}) \right| \right] \right] + \epsilon \cdot T}_{\text{term (a1)}} + \underbrace{\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[ \left( V_{t}(s_{t_{k+1}}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t_{k+1}}) \right) - \left( V_{t}(s_{t_{k}}) - V_{\theta^{\star}}^{\pi^{k}}(s_{t_{k}}) \right) \right] \right]}_{\text{term (a2)}}.$$

$$(C.52)$$

Recall that  $\mathbb{P}_{t_k}$  denotes the posterior distribution of  $\theta^*$  given  $\mathcal{D}_{t_k}$ . Under Assumption 5.9, we know that the distribution of  $\theta^k \mid \mathcal{D}_{t_k}$  is also  $\mathbb{P}_{t_k}$ , where  $\theta^k$  is defined in (C.43). Then, we have

term (a1) = 
$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \mathbb{E}_{\theta^*,\theta^k \sim \mathbb{P}_{t_k}} \left[ \left| \left( (B_{\theta^k} - B_{\theta^*}) V_t \right) (s_t, a_t) \right| \left| \mathcal{D}_{t_k} \right] \right] \right]$$
  

$$\leq \sqrt{T} \cdot \left( \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \mathbb{E}_{\theta^*,\theta^k \sim \mathbb{P}_{t_k}} \left[ \left| \left( (B_{\theta^k} - B_{\theta^*}) V_t \right) (s_t, a_t) \right|^2 \right| \mathcal{D}_{t_k} \right] \right] \right] \right)^{1/2},$$
(C.53)

where the first equality uses the tower property of the conditional expectation and the first inequality uses Cauchy-Schwarz inequality. Note that  $\mathbb{E}[|X-X'|^2] = 2\mathbb{V}[X]$ , if X and X' are two identically independently distributed variables. Recall that Assumption 5.9 tells that  $\theta^k$  (defined in (C.43)) and the data-generating parameter  $\theta^*$  are identically independently distributed given  $\mathcal{D}_{t_k}$ , which implies

term (a1) 
$$\leq \sqrt{2T} \cdot \left( \mathbb{E} \left[ \sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[ \sum_{t=t_k}^{t_{k+1}-1} \mathbb{V}_{\theta \sim \mathbb{P}_{t_k}} \left[ (B_{\theta^k} V_t)(s_t, a_t) \middle| \mathcal{D}_{t_k} \right] \right] \right] \right)^{1/2}$$
. (C.54)

Under Assumption 5.5, we apply the second property in Proposition C.1 on the right-hand side of (C.54) to have

term (a1) 
$$\leq 4L \cdot \sqrt{\mathbb{E}[H_0 - H_T]} \cdot \sqrt{T}$$
. (C.55)

Using the fact that any value function is bounded by L, we bound term (a2) in (C.14) as

term (a3) 
$$\leq 4L \cdot \mathbb{E}[K]$$
.

As the switching condition in Algorithm 6 implies  $H_{t_k} - H_{t_{k+1}} \ge \log 2$ , we apply Lemma C.3 to have

term (a2) 
$$\leq 4L + \frac{4L \cdot \mathbb{E}[H_0 - H_T]}{\log 2}$$
. (C.56)

Plugging (C.55) and (C.56) into (C.52), we have

$$(1 - \gamma) \cdot \text{term } (a) \le 4L \cdot \sqrt{\mathbb{E}[H_0 - H_T]} \cdot \sqrt{T} + \epsilon \cdot T + 4L + \frac{4L \cdot \mathbb{E}[H_0 - H_T]}{4L \cdot \mathbb{E}[H_0 - H_T]}$$
 (C.57)

Combining (C.49) and (C.57), we obtain

$$\Re(T) \leq \frac{4\sqrt{2}L \cdot \sqrt{\mathbb{E}[H_0 - H_T]}}{1 - \gamma} \cdot \sqrt{T} + \frac{2\epsilon}{1 - \gamma} \cdot T + \frac{4L}{1 - \gamma} + \frac{4L \cdot \mathbb{E}[H_0 - H_T]}{(1 - \gamma)\log 2}$$
$$= \mathcal{O}\left(\frac{L \cdot \sqrt{\mathbb{E}[H_0 - H_T]}}{1 - \gamma} \cdot \sqrt{T} + \frac{\epsilon}{1 - \gamma} \cdot T + \frac{L \cdot \mathbb{E}[H_0 - H_T]}{1 - \gamma}\right).$$

Thus, we finish the proof of Theorem 5.10.

### C.7 Relaxing Assumption 5.3 for Theorem 5.7

In this section, we show that Assumption 5.3 (Posterior Alignment) can be relaxed for Theorem 5.7 to accommodate a generalization error. We remove the dependency on Assumption 5.3 (Posterior Alignment) by introducing the following assumptions.

First, we characterize the data generation process for the pretraining dataset in the following assumption.

**Assumption C.4** (Pretraining Dataset Generation). We assume that the pretraining dataset  $\mathcal{D}_{pre}$  consists N i.i.d. tuples of  $(s', r, \mathcal{D}, s, a)$  generated by Algorithm 9.

For the pretraining dataset  $\mathcal{D}_{pre}$ , we denote by  $\mathbb{P}_{pre}$  the conditional distribution of (s', r) given  $(\mathcal{D}, s, a)$ . We show that  $\mathbb{P}_{pre}$  is equivalent to  $\mathbb{P}_{post}$  as follows,

$$\mathbb{P}_{\text{pre}}(s', r \mid \mathcal{D}, s, a) = \int_{\theta} \mathbb{P}_{\text{pre}}(s', r \mid s, a, \theta) \mathbb{P}_{\text{pre}}(d\theta \mid \mathcal{D}, s, a) 
= \int_{\Theta} P_{\theta}(s' \mid s, a) \mathbf{1}(r = r_{\theta}(s, a)) \mathbb{P}_{\text{pre}}(d\theta \mid \mathcal{D}, s, a) 
= \int_{\Theta} P_{\theta}(s' \mid s, a) \mathbf{1}(r = r_{\theta}(s, a)) \mathbb{P}_{\text{post}}(d\theta \mid \mathcal{D}) 
= \mathbb{P}_{\text{post}}(\xi_{(s, a)} \mid \mathcal{D}, s, a),$$
(C.58)

#### **Algorithm 9** The data collection process for the pretraining dataset.

- 1: **input**: Some (mixed) data collection policy  $\pi_{\text{collect}}$ .
- 2: **initialization**: Initialize the pretraining dataset  $\mathcal{D}_{pre} = \emptyset$ .
- 3: **for** n = 1, ..., N **do**
- 4: Reset the environment such that  $\theta^* \sim \mathbb{P}_0$ .
- 5: Receives  $s_0$  from the environment.
- 6: Initialize the memory buffer  $\mathcal{D}_0 = \emptyset$ .
- 7: **for** t = 0, ..., T **do**
- 8: Execute action  $a_t \sim \pi_{\text{collect}}(s_t)$  to receive reward  $r_t = r_{\theta^*}(s_t, a_t)$  and state  $s_{t+1} \sim P_{\theta^*}(\cdot \mid s_t, a_t)$  from the environment.
- 9: Update memory buffer  $\mathcal{D}_{t+1} \leftarrow \mathcal{D}_t \cup \{(s_t, a_t, s_{t+1}, r_t)\}.$
- 10: end for
- 11: Uniformly sample  $t_0 \in \{0, \dots, T\}$  and update the pretraining dataset  $\mathcal{D}_{pre} = \mathcal{D}_{pre} \cup \{(s_{t_0+1}, r_{t_0}, \mathcal{D}_{t_0}, s_{t_0}, a_{t_0})\}$
- 12: end for
- 13: **output**: The pretraining dataset  $\mathcal{D}_{pre}$ .

where the first equality uses Line 8 in Algorithm 9, the second equality uses the definition of the posterior of  $\theta^*$  in (2.3) and the fact that  $\theta$  and (s, a) are conditionally independent given  $\mathcal{D}$ , and the last equality uses (2.4). Denote by  $\mathcal{F}_{\text{LLM}}$  the function class of LLMs. In the next assumption, we assume that the function class  $\mathcal{F}_{\text{LLM}}$  contains the posterior of  $\xi_{(s,a)}$  in the underlying MDP, which is the conditional distribution of (s', r) given  $(\mathcal{D}, s, a)$  from  $\mathcal{D}_{\text{pre}}$ .

**Assumption C.5** (Realizability). We assume that there exists a LLM LLMPA with a posterior alignment, that is, there exists  $P^{\text{LLMPA}} \in \mathcal{F}_{LLM}$ , such that  $P^{\text{LLMPA}}(\xi_{(s,a)} | \mathcal{D}, s, a) = \mathbb{P}_{post}(\xi_{(s,a)} | \mathcal{D}, s, a)$ , for any query state-action pair (s, a) and in-context dataset  $\mathcal{D}$ .

We introduce the following assumption to require that LLMs are MLEs in the pretraining dataset with uniform coverage.

**Assumption C.6.** We assume that LLMs used in RAFA are Maximum Likelihood Estimators (MLEs) in the pretraining dataset  $\mathcal{D}_{pre}$  satisfying Assumption C.4, that is,

$$P^{\text{LLM}} = \underset{\widehat{P} \in \mathcal{F}_{LLM}}{\operatorname{argmax}} \sum_{(s', r, \mathcal{D}, s, a) \in \mathcal{D}_{pre}} \log \widehat{P}(s', r \mid \mathcal{D}, s, a).$$

Denote by  $\rho_{pre}$  the marginal population distribution of  $(\mathcal{D}, s, a)$  from  $\mathcal{D}_{pre}$ . We also assume that the pretraining dataset satisfies the following coverage condition:

$$\zeta = \sup_{t < T} \left\{ \left\| \frac{\mu_t}{\rho_{pre}} \right\|_{\infty} + \left\| \frac{\mu_t^*}{\rho_{pre}} \right\|_{\infty} \right\} < \infty.$$
 (C.59)

Here,  $\mu_t$  is the marginal distribution of  $(\mathcal{D}_{t_k}, s_t, \pi^k(s_t))$  and  $\mu_t^*$  is the marginal distribution of  $(\mathcal{D}_{t_k}, s_t, \pi^*(s_t))$  with  $s \sim \nu^*(\cdot \mid s_t)$  and  $(s_t, \mathcal{D}_{t_k})$  following the trajectory distribution of RAFA (Algorithm 4), where  $\nu^*$  is defined in (5.1).

We provide the generalization of Theorem 5.7 in the following corollary, which removes the dependency on Assumption 5.3 (Posterior Alignment).

Corollary C.7 (Generalization of Theorem 5.7). Under Assumptions C.5, C.6, 5.5, and 5.6, the Bayesian regret of RAFA (Algorithm 4) satisfies

$$\mathfrak{R}(T) = \mathcal{O}\left(\frac{(\kappa + 1)L \cdot \sqrt{\mathbb{E}[H_0 - H_T]}}{1 - \gamma} \cdot \sqrt{T} + \frac{\epsilon}{1 - \gamma} \cdot T + \frac{L \cdot \mathbb{E}[H_0 - H_T]}{1 - \gamma} + \frac{\zeta L \cdot \sqrt{\frac{\log(|\mathcal{F}_{\text{LLM}}|/\delta)}{N}} \cdot T}{N}\right),$$
Additional Regret Compared with Theorem 5.7

with probability at least  $1 - \delta$ . Here,  $\zeta$  is defined in (C.59),  $|\mathcal{F}_{LLM}|$  is the cardinality of the function class for LLMs, and N is the size of the pretraining dataset.

Comparing Corollary C.7 with Theorem 5.7, we remark that the additional regret decays to zero if N tends to infinity. Hence, we can recover the regret bound based on Assumption 5.3 (posterior alignment) approximately if the pretraining dataset has uniform coverage and is large enough.

*Proof of Corollary C.7.* We start with a standard concentration result for the maximum-likelihood estimator (MLE).

**Lemma C.8** (MLE Concentration). Let  $\mathcal{F}$  be a finite function class used to model a conditional distribution  $\mathbb{P}_{Y|X}(y|x)$  for  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ . Assume there is  $f^* \in \mathcal{F}$  such that  $\mathbb{P}(y|x) = f^*(y|x)$  (realizablility condition). Let  $\{(x_i, y_i)\}_{i=1}^N$  denote a dataset of i.i.d. samples where  $x_i \sim \mathbb{P}_X(x)$  and  $y_i \sim \mathbb{P}_{Y|X}(\cdot|x_i)$ . Let  $\widehat{f}$  be the MLE, which satisfies

$$\widehat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmax}} \sum_{i=1}^{N} \log f(y_i | x_i).$$

Then, it holds that

$$\mathbb{E}_{x \sim \mathbb{P}_X} \left[ d_{\text{TV}} \left( \widehat{f}(\cdot | x), p_{Y|X}(\cdot | x) \right) \right] \le \frac{8 \log(|\mathcal{F}|/\delta)}{N},$$

with probability at least  $1 - \delta$ .

Proof of Lemma C.8. See the proof of Theorem 21 of Agarwal et al. (2020).

Under Assumptions C.5 and C.6, we apply Lemma C.8 to show that

$$\mathbb{E}_{(\mathcal{D}, s, a) \sim \rho_{\text{pre}}} \left[ d_{\text{TV}}(P^{\text{LLMPA}}(\cdot | \mathcal{D}, s, a) \| P^{\text{LLM}}(\cdot | \mathcal{D}, s, a)) \right] \leq \sqrt{\frac{8 \log(|\mathcal{F}_{\text{LLM}}|/\delta)}{N}}$$

holds with probability at least  $1 - \delta$ .

For any fixed distribution  $\mu$  of  $(\mathcal{D}, s, a)$  satisfying  $\|\mu/\rho_{\text{pre}}\|_{\infty} < \infty$ , we use Hölder's inequality to know that

$$\mathbb{E}_{(\mathcal{D}, s, a) \sim \mu} \left[ d_{\text{TV}}(P^{\text{LLMPA}}(\cdot | \mathcal{D}, s, a) \| P^{\text{LLM}}(\cdot | \mathcal{D}, s, a)) \right]$$

$$\leq \left\| \frac{\mu}{\rho_{\text{pre}}} \right\|_{\infty} \cdot \mathbb{E}_{(\mathcal{D}, s, a) \sim \rho_{\text{pre}}} \left[ d_{\text{TV}}(P^{\text{LLMPA}}(\cdot | \mathcal{D}, s, a) \| P^{\text{LLM}}(\cdot | \mathcal{D}, s, a)) \right]$$

$$\leq \left\| \frac{\mu}{\rho_{\text{pre}}} \right\|_{\infty} \cdot \sqrt{\frac{8 \log(|\mathcal{F}_{\text{LLM}}|/\delta)}{N}}$$
(C.61)

holds with probability at least  $1 - \delta$ . Here,  $\|\cdot\|_{\infty}$  denotes the infinity norm. We denote the Bellman operator induced by LLMPA (the LLM with a posterior alignment) and  $\mathcal{D}_{t_k}$  as  $\widetilde{B}_k$ , which is defined as  $(\widetilde{B}_k V)(s,a) = r_{\text{LLMPA}(\mathcal{D}_{t_k})}(s,a) + (P_{\text{LLMPA}(\mathcal{D}_{t_k})}V)(s,a)$  for any s,a, and value function V. Then, by the definition of  $B_k$ , we have

$$\left| \left( (\widetilde{B}_k - B_k) V_t \right) (s, a) \right| = \left| \mathbb{E}_{P^{\text{LLMPA}}} [r + \gamma \cdot V(s')] - \mathbb{E}_{P^{\text{LLM}}} [r + \gamma \cdot V(s')] \right|$$

$$\leq 2L \cdot d_{\text{TV}} (P^{\text{LLMPA}} (\cdot | \mathcal{D}_{t_k}, s, a) || P^{\text{LLM}} (\cdot | \mathcal{D}_{t_k}, s, a)),$$
(C.62)

where the first inequality uses the definition of L (recall that L is the bound of |r + V(s)| for any reward r, state s, and value V) and Hölder's inequality. In the proof of Theorem 5.7 (the analysis of the regret of RAFA), we need to modify (C.14) and (C.21) with the following inequality

$$\begin{aligned} \left| \left( (B_k - B_{\theta^*}) V_t \right) (s, a) \right| &= \left| \left( (\widetilde{B}_k - B_{\theta^*}) V_t \right) (s, a) \right. \\ &+ \left. \left( (\widetilde{B}_k - B_k) V_t \right) (s, a) \right| \\ &\leq \left| \left( (\widetilde{B}_k - B_{\theta^*}) V_t \right) (s, a) \right| \\ &+ 2L \cdot d_{\text{TV}} (P^{\text{LLMPA}} (\cdot \mid \mathcal{D}_{t_k}, s, a) || P^{\text{LLM}} (\cdot \mid \mathcal{D}_{t_k}, s, a)), \end{aligned}$$

which holds for any state  $s \in \mathcal{S}$  and action  $a \in \mathcal{A}$ . By Proposition 5.4 (LLMs with posterior alignments perform BMA) and the fact that  $\widetilde{B}_k$  is the Bellman operator induced by LLMPA (the LLM with a posterior alignment) and  $\mathcal{D}_{t_k}$ , we can analyze  $[(\widetilde{B}_k - B_{\theta^*})V_t)(s, a)]$  in the same way as in the previous proof of Theorem 5.7 (the analysis of the regret of RAFA). It is clear that the additional regret by relaxing Assumption 5.3 (Posterior Alignment) is less than

$$\frac{1}{1-\gamma} \cdot \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} 2L \cdot d_{\text{TV}}(P^{\text{LLMPA}}(\cdot | \mathcal{D}_{t_{k}}, s, a) || P^{\text{LLM}}(\cdot | \mathcal{D}_{t_{k}}, s, a))\right]\right] \\
+ \frac{1}{1-\gamma} \cdot \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{\star}(\cdot | s_{t})} \left[2L \cdot d_{\text{TV}}(P^{\text{LLMPA}}(\cdot | \mathcal{D}_{t_{k}}, s, a) || P^{\text{LLM}}(\cdot | \mathcal{D}_{t_{k}}, s, a))\right]\right]\right] \\
\leq 8L \cdot \sqrt{\frac{2\log(|\mathcal{F}_{\text{LLM}}|/\delta)}{N}} \cdot \zeta \cdot T, \tag{C.63}$$

with probability at least  $1 - \delta$ , where the inequality uses (C.61) and the definition of  $\zeta$  in (C.59). Combining (C.63) and Theorem 5.7, we conclude the proof of Corollary C.7.

# D Missing Proofs in Appendix C

#### D.1 Proof of Lemma C.2

Proof of Lemma C.2. We prove the first part as follows. The Bellman equation (Sutton and Barto, 2018) connects  $Q_{\theta}^{\pi}(s, a)$  and  $V_{\theta}^{\pi}(s)$  by

$$Q_{\theta}^{\pi}(s, a) = r_{\theta}(s, a) + \gamma \left(P_{\theta} V_{\theta}^{\pi}\right)(s, a), \quad V_{\theta}^{\pi} = Q_{\theta}^{\pi}(s, a). \tag{D.1}$$

By the definition of  $B_{\theta}$ , we rewrite (D.1) as  $Q_{\theta}^{\pi}(s, a) = (B_{\theta}V_{\theta}^{\pi})(s, a)$ . For the left-hand side of (C.10) in the first part of Lemma C.2, we have

$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{t}(s_{t}) - V_{\theta^{*}}^{\pi^{k}}(s_{t})\right]\right] \\
= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} Q_{t}(s_{t}, a_{t}) - (B_{\theta^{*}}V_{\theta^{*}}^{\pi^{k}})(s_{t}, a_{t})\right]\right] \\
= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} Q_{t}(s_{t}, a_{t}) - (B_{\theta^{*}}V_{t})(s_{t}, a_{t})\right]\right] \\
+ \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} (B_{\theta^{*}}V_{t})(s_{t}, a_{t}) - r_{\theta^{*}}(s_{t}, a_{t}) - \gamma \cdot V_{t}(s_{t+1})\right]\right] \\
+ \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} r_{\theta^{*}}(s_{t}, a_{t}) + \gamma \cdot V_{\theta^{*}}^{\pi^{k}}(s_{t+1}) - (B_{\theta^{*}}V_{\theta^{*}}^{\pi^{k}})(s_{t}, a_{t})\right]\right] \\
+ \gamma \cdot \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{t}(s_{t+1}) - V_{\theta^{*}}^{\pi^{k}}(s_{t+1})\right]\right], \tag{D.2}$$

where the first equality uses  $a_t = \pi^k(s_t)$ , the condition  $Q(s, \pi^k(s)) = V_t(s)$  for any  $t_k \le t < t_{k+1}$  and k < K in Lemma C.2, and (D.1). Since we have

$$(B_{\theta^*}V)(s_t, a_t) = r_{\theta^*}(s, a) + \gamma \cdot \mathbb{E}_{s_{t+1} \sim P_{\theta^*}(\cdot \mid s_t, a_t)} [V(s_{t+1})],$$

terms (C1) and (C2) in (D.2) are zero. Meanwhile, term (D) in (D.2) satisfies

term (D) = 
$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\sum_{t=t_k}^{t_{k+1}-1} \left(V_t(s_t) - V_{\theta^*}^{\pi^k}(s_t)\right)\right]\right]$$
  
+  $\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^k} \left[\left(V_t(s_{t_{k+1}}) - V_{\theta^*}^{\pi^k}(s_{t_{k+1}})\right) - \left(V_t(s_{t_k}) - V_{\theta^*}^{\pi^k}(s_{t_k})\right)\right]\right],$  (D.3)

where term (B) is defined in the first part of Lemma C.2. Rearranging (D.2) and (D.3), we prove the first part of Lemma C.2.

Next, we show the proof of the second part of Lemma C.2, we. For the left-hand side of (C.11), we have

$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{\theta^{\star}}^{\pi^{\star}}(s_{t}) - V_{t}(s_{t})\right]\right] \\
= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} (B_{\theta^{\star}} V_{\theta^{\star}}^{\pi^{\star}})(s_{t}, \pi^{\star}(s_{t})) - Q_{t}(s_{t}, \pi^{k}(s_{t}))\right]\right] \\
= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} (B_{\theta^{\star}} V_{\theta^{\star}}^{\pi^{\star}})(s_{t}, \pi^{\star}(s_{t})) - (B_{\theta^{\star}} V_{t})(s_{t}, \pi^{\star}(s_{t}))\right]\right] \\
+ \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} (B_{\theta^{\star}} V_{t})(s_{t}, \pi^{\star}(s_{t})) - Q_{t}(s_{t}, \pi^{k}(s_{t}))\right]\right] \\
= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \gamma \cdot \mathbb{E}_{s' \sim P_{\theta^{\star}}(\cdot \mid s_{t}, \pi^{\star}(s_{t}))} \left[V_{\theta^{\star}}^{\pi^{\star}}(s') - V_{t}(s')\right]\right]\right] \\
+ \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} (B_{\theta^{\star}} V_{t})(s_{t}, \pi^{\star}(s_{t})) - Q_{t}(s_{t}, \pi^{k}(s_{t}))\right]\right], \tag{D.5}$$

where the first equality uses the Bellman optimality equation in (2.2), the condition  $Q_t(s, \pi^k(s)) = V_t(s)$ , and  $Q_{\theta^*}^*(s, \pi^*(s)) = V_{\theta^*}^*(s)$  for any state  $s, t_k \leq t < t_{k+1}$ , and k < K. Here, the last equality uses the definition of  $B_{\theta^*}$ . For the simplicity of discussions, we define functions  $F_t, M_t \in \{S \mapsto \mathbb{R}\}$  and the linear operator  $T \in \{\{S \mapsto \mathbb{R}\} \mapsto \{S \mapsto \mathbb{R}\}\}$  as

$$F_t(s) = V_{\theta^*}^{\pi^*}(s) - V_t(s),$$

$$M_t(s) = (B_{\theta^*}V_t)(s, \pi^*(s)) - Q_t(s_t, \pi^k(s)),$$

$$(\mathcal{T}f)(s) = \mathbb{E}_{s' \sim P_{\theta^*}(\cdot \mid s, \pi^*(s))}[f(s')],$$
(D.6)

for any state s and function  $f \in \{S \mapsto \mathbb{R}\}$ . Here, we denote by  $\{S \mapsto \mathbb{R}\}$  the class of all the functions defined on S. By the definitions of  $F_t$ ,  $M_t$ , and T in (D.6), it is clear that

$$F_t(s) = M_t(s) + \gamma \cdot (\mathcal{T}F_t)(s), \tag{D.7}$$

for any state  $s \in \mathcal{S}$ . Then, we introduce the following lemma to bound  $F_t$  by (D.7).

**Lemma D.1.** For the operator  $\mathcal{T}$  defined in (D.6), two arbitrary bounded functions f, m defined on the state space  $\mathcal{S}$ , and any  $\gamma \in [0, 1)$ , if

$$f(s) = m(s) + \gamma \cdot (\mathcal{T}f)(s)$$
 (D.8)

holds for any state  $s \in \mathcal{S}$ , then it holds that for any state  $s \in \mathcal{S}$ ,

$$f(s) = \sum_{\tau=0}^{\infty} \gamma^{\tau} \cdot \left( \underbrace{(\mathcal{T} \circ \dots \circ \mathcal{T})}_{\tau \text{ times}} m \right) (s). \tag{D.9}$$

Proof of Lemma D.1. By the condition in (D.8), we have

$$f(s) = m(s) + \gamma \cdot \left( \mathcal{T}(m + \gamma \cdot (\mathcal{T}f)) \right)(s)$$
  
=  $m(s) + \gamma \cdot (\mathcal{T}m)(s) + \gamma^2 \cdot \left( (\mathcal{T} \circ \mathcal{T})f \right)(s),$  (D.10)

where the last equality relies on the linearity of the operator  $\mathcal{T}$ . Repeating the process in (D.10) for  $N \in \mathbb{N}$  times, we have that

$$f(s) = \gamma^{N+1} \cdot \underbrace{\left( (\mathcal{T} \circ \dots \circ \mathcal{T}) f \right)(s)}_{(N+1) \text{ times}} f(s) + \sum_{\tau=0}^{N} \gamma^{\tau} \cdot \underbrace{\left( (\mathcal{T} \circ \dots \circ \mathcal{T}) m \right)(s)}_{\tau \text{ times}} m(s). \tag{D.11}$$

Since both f and m are bounded functions, we use (D.6) to know that  $\underbrace{(\mathcal{T} \circ \ldots \circ \mathcal{T})}_{\tau \text{ times}} f$  and  $\underbrace{(\mathcal{T} \circ \ldots \circ \mathcal{T})}_{\tau \text{ times}} m$  are also bounded for any  $\tau \in \mathbb{N}$ . As  $\gamma \in [0,1)$ , we let N tend to the infinity to transform (D.11) to

$$f(s) = \sum_{\tau=0}^{\infty} \gamma^{\tau} \cdot (\underbrace{(\mathcal{T} \circ \dots \circ \mathcal{T})}_{\tau \text{ times}} m)(s),$$

for any state  $s \in \mathcal{S}$ . Then, we conclude the proof for Lemma D.1.

By (D.7) and Lemma D.1, we have

$$F_t(s) = \sum_{\tau}^{\infty} \gamma^{\tau} \cdot (\underbrace{(\mathcal{T} \circ \dots \circ \mathcal{T})}_{\tau \text{ times}} M_t)(s)$$

Recalling the definition of the optimal  $\gamma$ -discounted visitation measure in (5.1), we further have

$$F_t(s) = \frac{1}{1 - \gamma} \cdot \mathbb{E}_{s' \sim \nu^*(\cdot \mid s)}[M_t(s')]. \tag{D.12}$$

Plugging (D.12) and (D.6) into (D.5), we have

$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} V_{\theta^{*}}^{\pi^{*}}(s_{t}) - V_{t}(s_{t})\right]\right] \\
= \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \frac{1}{1-\gamma} \cdot \mathbb{E}_{s' \sim \nu^{*}(\cdot \mid s)} \left[ (B_{\theta^{*}} V_{t})(s_{t}, \pi^{*}(s_{t})) - Q_{t}(s_{t}, \pi^{k}(s_{t})) \right] \right]\right] \\
= \frac{1}{1-\gamma} \cdot \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s)} \left[ (B_{\theta^{*}} V_{t})(s, \pi^{*}(s)) - Q_{t}(s, \pi^{*}(s)) \right] \right] \right] \\
+ \frac{1}{1-\gamma} \cdot \mathbb{E}\left[\sum_{k=0}^{K-1} \mathbb{E}_{\pi^{k}} \left[\sum_{t=t_{k}}^{t_{k+1}-1} \mathbb{E}_{s \sim \nu^{*}(\cdot \mid s)} \left[ (Q_{t}(s, \pi^{*}(s)) - Q_{t}(s, \pi^{k}(s))) \right] \right] \right]. \tag{D.13}$$

Multiplying  $(1-\gamma)$  on the two sides of (D.13), we prove the second part of Lemma C.2.

# E Linear Special Case

We specialize RAFA to a linear setting and characterize the Bayesian regret. In particular, we define a Bayesian variant of linear kernel MDPs (Yang and Wang, 2020, 2019; Cai et al., 2020; Zhou et al., 2021b). Here,  $\mathbb{E}_{s'\sim P_{\theta}(\cdot|s,a)}[V(s')]$  is linear in a feature  $\psi_V(s,a)\in\mathbb{R}^d$  for an arbitrary parameter  $\theta\in\mathbb{R}^d$ , while the prior and posterior distributions of the data-generating parameter  $\theta^*\in\mathbb{R}^d$  are Gaussian. Specifically,  $\psi_V(s,a)$  maps the value function V and the state-action pair (s,a) to a d-dimensional vector. Recall that  $\rho$  is the initial distribution of states, t is the step index, and T is the total number of steps. Also,  $\mathbb{P}_t$  is the posterior distribution at the t-th step.

**Definition E.1** (Bayesian Linear Kernel MDP (Ghavamzadeh et al., 2015; Yang and Wang, 2020, 2019; Cai et al., 2020; Zhou et al., 2021b)). A Bayesian linear kernel MDP M satisfies

$$V(s_{t+1}) \mid s_t, a_t \sim \mathcal{N}(\psi_V(s_t, a_t)^\top \theta, 1)$$

for all  $t \geq 0$ ,  $(s_t, a_t) \in \mathcal{S} \times \mathcal{A}$ ,  $s_{t+1} \sim P_{\theta}(\cdot \mid s_t, a_t)$ ,  $\theta \in \mathbb{R}^d$ , as well as all value function V. Here,  $\psi_V(s, a)$  maps V and (s, a) to a d-dimensional vector, which satisfies  $\|\psi_V(s, a)\|_2 \leq R$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$  and all V. Also, M also satisfies  $|\mathbb{E}_{s_0 \sim \rho} V(s_0)| \leq R$  for all V. Here, R is a positive constant that is independent of t and T. The prior distribution of the datagenerating parameter  $\theta^* \in \mathbb{R}^d$  is  $\mathcal{N}(0, \lambda I_d)$ , where  $\lambda$  is a positive constant. Here,  $\psi_V$  is known and  $\theta^*$  is unknown. Without loss of generality, we assume that the reward function is deterministic and known, i.e.,  $(B_{\theta}V)(\cdot, \cdot) = r(\cdot, \cdot) + \gamma \cdot (P_{\theta}V)(\cdot, \cdot)$  for a known reward function r and any  $\theta$ .

By Definition E.1, we obtain the closed form of the posterior  $\mathbb{P}_t$  as follows,

$$\theta \mid \mathcal{D}_t \sim \mathcal{N}(\widehat{\theta}_t; \Sigma_t^{-1}),$$

where

$$\widehat{\theta}_t = \left(\lambda I_d + \sum_{i=0}^{t-1} \psi_{V_i}(s_i, a_i) \psi_{V_i}(s_i, a_i)^\top \right)^{-1} \left(\sum_{i=0}^{t-1} \psi_{V_i}(s_i, a_i) V_i(s_{i+1})\right)$$
(E.1)

and

$$\Sigma_{t} = \lambda I_{d} + \sum_{i=0}^{t-1} \psi_{V_{i}}(s_{i}, a_{i}) \psi_{V_{i}}(s_{i}, a_{i})^{\top}.$$
 (E.2)

Hence, the posterior entropy is

$$H_t = H(\mathbb{P}_t) = 1/2 \cdot \log(\det(\Sigma_t)) + d/2 \cdot (1 + \log(2\pi)).$$
 (E.3)

We specialize the switching condition in Algorithm 6 as follows,

$$H_{t_k} - H_t = 1/2 \cdot \log(\det(\Sigma_{t_k})) - 1/2 \cdot \log(\det(\Sigma_t)) > \log 2, \tag{E.4}$$

which is equivalent to  $\det(\Sigma_{t_k}) > 4 \cdot \det(\Sigma_t)$ . This switching condition is also similarly adopted in work for RL (Zhou et al., 2021b; Abbasi-Yadkori and Szepesvári, 2015). As a result, we have

$$\det(\Sigma_{t_k}) \le 4 \cdot \det(\Sigma_t) \tag{E.5}$$

for all  $t_k \leq t < t_{k+1}$  and k < K.

**Verification of Assumption 5.5** We verify the regularity assumption (Assumption 5.5.) on Bayesian linear kernel MDPs as follows. By (E.4), the condition

$$H_{t_1} - H_{t_2} = 1/2 \cdot \log(\det(\Sigma_{t_1})) - 1/2 \cdot \log(\det(\Sigma_{t_2})) \le \log 2$$

is equivalent to  $\det(\Sigma_{t_1}) \leq 4 \det(\Sigma_{t_2})$ . Since  $t_1 < t_2$  and the posterior variance matrix is positive definite, we have  $\Sigma_{t_1}^{-1} \succeq \Sigma_{t_1}^{-1}$  and  $\det(\Sigma_{t_2}^{-1}) \leq 4 \det(\Sigma_{t_1}^{-1})$ . By the definition of the information gain and (E.4), we have

$$I(\theta; \xi_{(s,a)} | \mathcal{D}_t) = H(\theta | \mathcal{D}_t) - H(\theta | \xi_{(s,a)}, \mathcal{D}_t)$$

$$= 1/2 \cdot \log \left( \frac{\det(\psi_{V_{t_2}}(s, a)\psi_{V_{t_2}}^{\top}(s, a) + \Sigma_t)}{\det(\Sigma_t)} \right)$$

$$= 1/2 \cdot \log \left( 1 + \psi_{V_{t_2}}(s, a)^{\top} \Sigma_t^{-1} \psi_{V_{t_2}}(s, a) \right), \tag{E.6}$$

for  $t = t_1$  and  $t_2$ . Here, the last equality uses the matrix determinant lemma.

Plugging (E.9) into (E.6), we have

$$I(\theta; \xi_{(s,a)} | \mathcal{D}_{t_2}) = 1/2 \cdot \log(1 + \psi_{V_{t_2}}(s, a)^{\top} \Sigma_t^{-1} \psi_{V_{t_2}}(s, a))$$

$$= 1/2 \cdot \log(1 + \psi_{V_{t_2}}(s, a)^{\top} \Sigma_t^{-1} \psi_{V_{t_2}}(s, a))$$

$$\geq \log(1 + d)/(2d) \cdot \|\psi_{V_{t_2}}(s, a)\|_{\Sigma_{t_2}^{-1}}^2, \tag{E.7}$$

where the second equality uses the matrix determinant lemma and the first inequality uses the fact that  $\log(1+x)/x$  is an increasing function for  $x \ge 0$  and

$$0 \leq \psi_{V_{t_2}}(s, a)^{\top} \Sigma_t^{-1} \psi_{V_{t_2}}(s, a)$$

$$\leq \psi_{V_{t_2}}(s, a)^{\top} (\psi_{V_{t_2}}(s, a) \psi_{V_{t_2}}(s, a)^{\top})^{-1} \psi_{V_{t_2}}(s, a)$$

$$= \operatorname{tr} (\psi_{V_{t_2}}(s, a) \psi_{V_{t_2}}(s, a)^{\top} (\psi_{V_{t_2}}(s, a) \psi_{V_{t_2}}(s, a)^{\top})^{-1})$$

$$= d. \tag{E.8}$$

Here, the first inequality uses the nonnegativity of a quadratic form, the first equality uses  $\operatorname{tr}(a^{\top}b) = \operatorname{tr}(ba^{\top})$  for two arbitrary vectors a and b, and the second inequality uses (E.2). By (E.2), we know that

$$\|\psi_{V_{t_2}}(s,a)\|_{\Sigma_{t_1}^{-1}}^2 \le 4 \cdot \|\psi_{V_{t_2}}(s,a)\|_{\Sigma_{t_2}^{-1}}^2, \tag{E.9}$$

where the inequality invokes the following lemma (Lemma E.2).

**Lemma E.2** (Lemma 12 in Abbasi-Yadkori et al. (2011)). Suppose  $A, D \in \mathbb{R}^{d \times d}$  are two positive definite matrices satisfying that  $A \succeq D$ , then for any  $\mathbf{x} \in \mathbb{R}^d$ ,  $\|\mathbf{x}\|_A \leq \|\mathbf{x}\|_D \cdot \sqrt{\det(A)/\det(D)}$ .

Rearranging (E.7), we have

$$\frac{8d}{\log(1+d)} \cdot I(\theta; \xi_{(s,a)} \mid \mathcal{D}_{t_2}) \ge 4 \cdot \|\psi_{V_{t_2}}(s,a)\|_{\Sigma_{t_2}^{-1}}^{2}$$

$$\ge \|\psi_{V_{t_2}}(s,a)\|_{\Sigma_{t_1}^{-1}}^{2}$$

$$\ge \log(1 + \psi_{V_{t_2}}(s,a)^{\top} \Sigma_{t_1}^{-1} \psi_{V_{t_2}}(s,a))$$

$$= 2 \cdot I(\theta; \xi_{(s,a)} \mid \mathcal{D}_{t_1}), \tag{E.10}$$

where the second inequality uses (E.9), the last inequality uses the fact that  $x \ge \log(1+x)$  for any  $x \ge 0$ , and the last equality use (E.6). By (E.10), we know that Bayesian linear kernel MDPs (Definition E.1) satisfy Assumption 5.5 with the coefficient  $\eta = d/\log(1+d)$ .

Analysis of the Cumulative Posterior Entropy  $H_0 - H_T$ . Next, we study the upper bound of the cumulative Information gain  $H_0 - H_T$  in Bayesian linear kernel MDPs. By the definition of  $\Sigma_t$  in (E.2), we have  $\log \det(\Sigma_0) = d \cdot \log \lambda$  and

$$\log \det(\Sigma_T) = \log \det\left(\lambda I_d + \sum_{t=0}^{T-1} \psi_{V_t}(s_t, a_t) \psi_{V_t}^{\top}(s_t, a_t)\right)$$

$$\leq d \cdot \log\left(1/d \cdot \operatorname{tr}\left(\lambda I_d + \sum_{t=0}^{T-1} \psi_{V_t}(s_t, a_t) \psi_{V_t}^{\top}(s_t, a_t)\right)\right)$$

$$= d \cdot \log\left(1/d \cdot \left(\lambda d + \sum_{t=0}^{T-1} \|\psi_{V_t}(s_t, a_t)\|_2^2\right)\right)$$

$$\leq d \cdot \log(\lambda + TR^2/d) \tag{E.11}$$

almost surely. Here, the first inequality uses the relationship between the trace and the determinant of a square matrix, the second equality uses  $\operatorname{tr}(a^{\top}b) = \operatorname{tr}(ba^{\top})$  for two arbitrary vectors a and b, and the last inequality uses the fact that  $\|\psi_V(s,a)\|_2$  is upper bounded by R for all  $(s,a) \in \mathcal{S} \times \mathcal{A}$  and V. Hence, we have

$$H_0 - H_T = \mathcal{O}\left(d \cdot \log(1 + TR^2/(d\lambda))\right) \tag{E.12}$$

almost surely.

**Regret Bounds.** With the verification of Assumption 5.5 and the analysis of the upper bound of the cumulative Information gain  $H_0 - H_T$  in Bayesian linear kernel MDPs, we are ready to specialize the theorems in Appendix 5.3 in Bayesian linear kernel MDPs if we

determine an appropriate L (the bound of the value). We analyze the bound of V(s') for  $s' \sim P_{\theta}(s, a)$  and  $\theta \sim \mathcal{N}(\mu, \lambda I_d)$ . Define that  $\tilde{\epsilon} \sim \mathcal{N}(0, 1)$ . By Definition E.1, we have

$$|V(s')| = |\widetilde{\epsilon} + \psi_V(s, a)^{\top} \theta|$$

$$\leq |\widetilde{\epsilon}| + |\psi_V(s, a)^{\top} (\theta - \mu)| + |\psi_V(s, a)^{\top} \mu|$$

$$\leq |\widetilde{\epsilon}| + ||\psi_V(s, a)||_2 \cdot ||\theta - \mu||_2 + ||\psi_V(s, a)||_2 \cdot ||\mu||_2$$

$$\leq |\widetilde{\epsilon}| + R \cdot ||\theta||_2 + R \cdot ||\mu||_2,$$
(E.13)

where the first inequality uses the triangle inequality, the second inequality uses the Cauchy-Schwartz inequality, and the last inequality uses the definition of R in Definition E.1. By Definition E.1, (E.13), and the tail behavior of the Gaussian distribution (Ghosh, 2021), we have

$$|V(s)| \le \sqrt{2 \cdot \log(2/\delta)} + R \cdot ||\mu||_2 + R \cdot \sqrt{2\lambda d \cdot \log(2d\delta)}$$

for any  $s \in \mathcal{S}$  and value function V with probability at least  $1-\delta$ . Since the prior distribution of  $\theta$  is  $\mathcal{N}(0, \lambda I_d)$ , it is natural to restrict  $\mu$  such that  $\|\mu\|_2 \leq cd \cdot \log(2d)$  for some absolute constant c. Then, we apply the union bound of all T value functions in RAFA and the variants (Algorithms 4, 5, and 6) to have

$$|V_t(s)| \le \sqrt{2 \cdot \log(2T/\delta)} + R \cdot ||\mu||_2 + R \cdot \sqrt{2\lambda d \cdot \log(2dT\delta)}$$
  

$$\le (c+1)R \cdot \sqrt{2\lambda d \log(2dT/\delta)}$$
(E.14)

for any t < T,  $s \in \mathcal{S}$ , and value function V with probability at least  $1 - \delta$ . Hence, we can select  $L = (c+1)R \cdot \sqrt{2\lambda d \log(2dT/\delta)}$  in Theorems 5.7, 5.8, and 5.10. By specializing Theorems 5.7, 5.8, and 5.10, we summarize the corresponding regret bounds in Table 5 for Algorithms 4, 5, and 6, respectively. Here, we choose the planning suboptimality of  $PL^{\epsilon}$  to be  $\epsilon = \mathcal{O}(1/\sqrt{T})$  and all the bounds hold with probability at least  $1 - \delta$ .

Algorithm	Bayesian Regret
RAFA (Algorithm 4)	$\mathcal{O}((1-\gamma)^{-1}(\kappa+1)\sqrt{d^3T}\cdot\log(dT/\delta))$
RAFA with Optimistic Bonus (Algorithm 5)	$((1-\gamma)^{-1}\sqrt{d^3T}\cdot\log(dT/\delta))$
RAFA with Posterior Sampling (Algorithm 6)	$\mathcal{O}((1-\gamma)^{-1}\sqrt{d^3T}\cdot\log(dT/\delta))$

Table 5: Bayesian regret of variants of RAFA in Bayesian linear kernel MDPs (see Definition E.1). Here, we choose the planning suboptimality of  $PL^{\epsilon}$  to be  $\epsilon = \mathcal{O}(1/\sqrt{T})$  and all the bounds hold with probability at least  $1 - \delta$ .

## F More Experiments

In what follows, we provide the detailed setups and additional results of our experiments.

#### F.1 Game of 24

**Task Setup.** Figure 12 gives an illustrative example for Game of 24.

[Illustrative example for Game of 24]
 Numbers: [2, 5, 8, 11]
 Arithmetic Operations: [+, -, ×, /, (, )]
 Solution: (11 - 5) × 8/2 = 24

Figure 12: An illustrative example of the Game of 24. The player uses combinations of basic arithmetic operations with four given numbers to get 24.

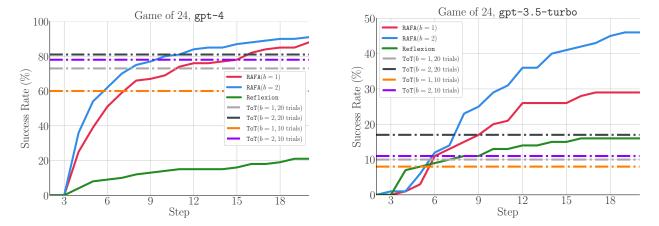


Figure 13: Sample efficiency on Game of 24. RAFA agent achieves strong performance due to an orchestration of reasoning and acting. The success rate at a given step is the number of tasks that is solved within the given step.

Following Yao et al. (2023a), we use the same subset indexed 901-1,000 from a total of 1,362 tasks collected from 4nums.com. The index is arranged from easy to hard by human solving time so the subset is relatively challenging. The agent receives a reward of 1 if the proposed formula is correct and the proposed formula is accepted and concatenated into the state; if the final result is exactly 24, the agent receives a reward of 10, and the episode terminates. Otherwise, the agent receives a reward of 0, and the proposed formula is not accepted. We limit the maximum trials for each task to 20 to avoid meaningless retries. The task is successful if the agent receives a return larger than 10 <sup>1</sup> (i.e., find a valid solution

<sup>&</sup>lt;sup>1</sup>For gpt-3.5-turbo, we report the success rate when the agent receives a return no less than 3 (i.e., find all sub-steps to get 24 but not necessarily generate a whole correct formula). This is because ToT with gpt-3.5-turbo is known to suffer from correctly get a whole formula due to limited reasoning ability and non-perfect prompts. See <a href="https://github.com/princeton-nlp/tree-of-thought-llm/issues/24">https://github.com/princeton-nlp/tree-of-thought-llm/issues/24</a> for more details.

within 20 steps). We report the final success rate and sample efficiency for each method on the subset of 100 tasks. Notably, a task is considered successful if the RAFA agent returns one and only one correct formula, which is more strictly evaluated than Tree of Thoughts (ToT, Yao et al. (2023a)): we allow open-loop agents like ToT to retry 20 times and consider them successful if they generate a valid solution in any of the 20 trials. For CoT (Wei et al., 2022) and Reflexion (Shinn et al., 2023) agents, we allow them to reflect on the environment's feedback but require them to generate a plan immediately without sophisticated reasoning.

**RAFA Setup.** In the Game of 24, the RAFA agent uses ToT as the planner, regenerates a plan when the agent receives a zero reward and continues acting according to the previous plan when the agent receives a positive reward. We set the base ToT planner with beam search width b = 1, 2 and use both gpt-3.5-turbo and gpt-4 to test the RAFA's boost-up over LLM agents with different reasoning abilities. We set the temperature t = 0.2 by default to favor rigorous reasoning and t = 0.7 for majority voting.

Reduced Hallucination Through Interaction. A comprehensive review of various method proposals revealed significant hallucination, especially with gpt-3.5-turbo. A common hallucination is that the agent believes she can reuse the same number (e.g. using the number 2 twice as illustrated in Figure 4). RAFA efficiently mitigates such hallucination by actively interacting with the environment, displaying exceptional hallucination resistance and improved performance.

Enhanced Efficiency Through Planning. Evidenced in Figure 5, the RAFA agent substantially surpasses the Reflexion baseline, reflecting heightened efficiency and minimized regret by negating careless trials. For example, without carefully planning, agent may give negative answers, e.g., "Impossible to obtain 24 with the given numbers, or unchecked answers, e.g., "Answer: 6 \* 9 / (3 - 2) = 24". This reduction of careless trails is especially achieved when a strong backbone LLMs (e.g., gpt-4) is used, even with a basic planning method, such as BFS with B = 1.

Ablation Study. The RAFA agent's performance is dissected by individually examining its components: (1) Planning modules or model/elite LLM, (2) Reflection modules or critic LLM, and (3) Different LLMs. Results, displayed in Table 2 and Figure 5, affirm the substantial contribution of each segment to the aggregate performance. Compared to absent or rudimentary zero-shot planning, a basic planner markedly enhances overall performance. However, augmenting planner strength only offers marginal performance enhancements. Both critic LLM and robust LLM usage emerge as pivotal for optimal performance.

### F.2 ALFWorld

**Task Setup.** The action space of ALFWorld consists of high-level actions such as "heat a potato with a microwave", which is executed in the underlying embodied simulator through

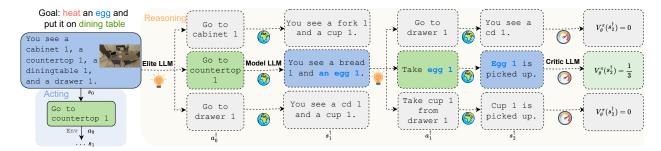


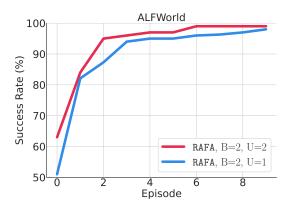
Figure 14: An illustration of RAFA in the ALFWorld environment.

low-level action primitives. The egocentric visual observations of the simulator are translated into natural language before being provided to the agent. The state is the history of the observations. If a task goal can be precisely achieved by the agent, it will be counted as a success.

RAFA Setup. In the ALFWorld environment, the RAFA planner is instantiated as Breadth First Search (BFS). Specifically, B and U are both set to 2, and we use gpt-3 (text-davinci-003) for the Critic, Model, and Elite modules. Besides, since it is challenging to prompt the LLM with the stored full trajectories in the memory buffer due to the token limit, we make the following modifications: the Model LLM instance uses only the partial trajectory executed so far in the current episode, and the Elite LLM instance uses the same partial executed trajectory with additional model-generated state-action pairs during the planning subroutine. When switching is triggered after 20 failed timesteps (i.e., an episode), a summary from the failure trajectory is generated by gpt-4 and added to the Critic prompt.

Reduced Hallucination Through Interaction. The baselines are more likely to hallucinate when the target object is not found after exploring many locations. On the other hand, the critic LLM used in RAFA is able to probe the hallucination by generating the summary "In this environment, my critic assigned a 1/3 value after taking a knife. However, the task is to take and cool a tomato." and avoid it in the next episode. Therefore, RAFA is more sample-efficient due to an orchestration of reasoning and acting and the ability to mitigate hallucination through interaction.

Ablation Study. To better understand the role that the planning subroutine plays in the RAFA algorithm, we conduct ablation studies on the search depth U and search breadth B. The results are shown in Figure 15 and 16, respectively. We observe that when setting the search depth to B = U = 2, the success rate is higher than when setting the search depth to U = 1 or setting the search breadth B = 1, especially at the initial episode. This indicates that the reasoning ability of RAFA is enhanced through the planning subroutine. Besides, the algorithm is also more sample-efficient when setting B = U = 2, indicating a better capacity for learning and planning through interaction and reasoning.



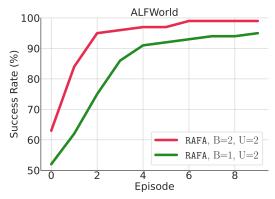


Figure 15: Ablation on the search depth U in the ALFWorld environment.

Figure 16: Ablation on the search breadth B in the ALFWorld environment.

#### F.3 BlocksWorld

**Task Setup.** The reported success rates are averaged in tasks that require different minimum steps. Specifically, the evaluation is conducted in 57 4-step tasks and 114 6-step tasks. We set the state as the current arrangement of the blocks and the actions contain Stack, Unstack, Put, and Pickup, coupled with a block being operated.

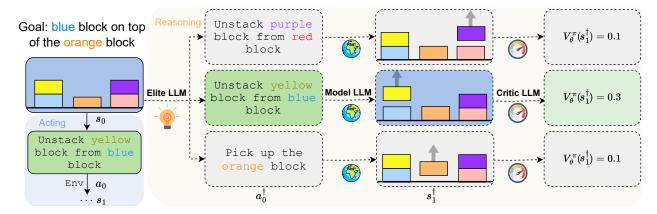


Figure 17: Illustration of RAFA in the BlocksWorld environment.

RAFA Setup. The search space is up to 5<sup>4</sup> for a 4-step task and is up to 5<sup>6</sup> for a 6-step task. For 4-step tasks, RAFA can achieve over 50% success rate within 8 learning steps with Vicuna-13B(v1.3) and achieve over 80% success rate within 8 learning steps with Vicuna-33B(v1.3). For 6-step tasks, RAFA can achieve over 40% success rate within 20 learning steps with Vicuna-13B(v1.3) and achieve over 50% success rate within 20 learning steps with Vicuna-33B(v1.3). Empirical results show that Vicuna could produce wrong state transition in the planning phase. RAFA can mitigate hallucination with feedback from failure trajectories and active exploration. One can draw such a conclusion by comparing RAFA with RAP as RAP does not receive feedback from the real environment.

### F.4 Tic-Tac-Toe

**Task Setup.** Tic-Tac-Toe (Beck, 2008) is a competitive game in which two players take turns to mark a three-by-three grid with X or O, and a player succeeds when their marks occupy a diagonal, horizontal, or vertical line.

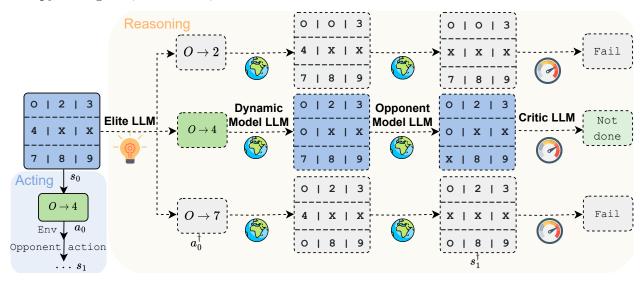


Figure 18: Illustration of RAFA (playing O) in the Tic-Tac-Toe game. States are represented by a numbered 3×3 grid and actions are represented by a number between 1-9. The opponent is considered part of the environment.

We adopt the convention that X plays first. As illustrated below in Figure 18, we use a numbered  $3 \times 3$  grid to represent a state and a number between 1 and 9 to represent an action, which also illustrates the transition and reward function. Although Tic-Tac-Toe is a solved game with a forced draw assuming the best play from both players, it remains a challenge for LLMs to accomplish this task even when prompted to play only the optimal moves. We collected the battle outcomes between different LLM models in Table 6, where we notice that gpt-4 performs worse when playing as "O". Thus, in our experiments, we let RAFA play as "O" and let baseline LLM models play as "X".

X wins : Tie : O wins		O	
		gpt-3.5	gpt-4
X	gpt-3.5	55%:35%:10%	90%:0%:10%
	gpt-4	65%:15%:20%	90%:0%:10%

Table 6: Probability of "X wins," "Tie," and "O wins" in Tic-Tac-Toe. The results are obtained by averaging over 20 simulated games.

**RAFA Setup.** For implementation, we set B=3 and adopt MCTS to evaluate the proposed actions. We set U=4 which is the maximum game depth. We set a prediction-

based switching condition triggered when the prediction does not agree with the observation. Specifically, policy switches when one of the following events occurs:

- The RAFA agent takes an action and predicts the next state, which is different from the observed next state.
- Before the opponent takes an action, the RAFA agent tries to predict such an action, which is different from the actual action that the opponent takes.
- After the opponent takes an action, RAFA agent predicts the next state, which is different from the observed next state.
- The RAFA agent predicts the current game status (X wins, O wins, Tie, Not finished), which is different from the environment's feedback.

Besides, we use the ground truth of those predictions to update the agent's belief of the world, which also implicitly affects the agent's policy.

We define a discrete reward function with r = -1, 0, 1 corresponding to lose, tie, and win. The agent only gets rewards when the current episode is completed. We define the score of an agent as its expected reward which can be approximated by simulation. The empirical results are shown in figure 19. We conduct experiments using both gpt-4 as the backend. The score of RAFA (B=4) increases as it interacts more with the environment. By analyzing the generated trajectories, we also notice that although RAFA agent is not perfect, it exploits the weakness of the baseline model well, which is why it almost never loses after 7 episodes.

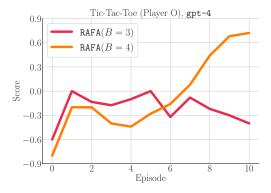


Figure 19: Score curves in the Tic-Tac-Toe game. We use gpt-4 as backend. Results are averaged across 10 simulations and smoothed with a window size of 5.

## G Prompts

In this section, we give details of the prompts used for each task.

#### G.1 Game of 24

Critic LLM. For the LLM instance of the Critic, we prompt it with the current action (formula) with reward and feedback from the environment. The critic is required to determine whether each formula is valid or not and give a "sure" or "impossible" label for the formula. We use critic prompts to generate reflection for formula proposal and evaluation, respectively.

Critic prompt (for formula proposal)

Now we would like to play a game of 24. That is, given 4 numbers, try to use them with arithmetic operations (+ - \* /) to get 24. Now we consider the following puzzle: {input}.

Here is an attempt answer:

{answer}

And we have the following feedback:

 $\{feedback\}$ 

Now using the above feedback, give 'sure' or 'impossible' labels for each formula with left numbers from each step. Give 'sure' if the formula is correct and can lead to 24 and give 'impossible' if the formula is incorrect or illegal. First repeat the formula with left numbers from each step above and then give the label, with the following form: {{formula}} (left: {{left numbers}}): {{label}}.

Critic prompt (for formula evaluation)

Now we would like to play a game of 24. That is, given 4 numbers, try to use them with arithmetic operations (+ - \* /) to get 24. Now we consider the following puzzle: {input}.

Here is an attempt answer:

{answer}

And we have the following feedback:

{feedback}

Now using the above feedback, give 'sure' or 'impossible' labels for left numbers from each step. Give 'sure' if the formula is correct and left numbers can lead to 24 and give 'impossible' if the formula is incorrect or illegal. First repeat the left numbers from each step above and then give the label, with the following form: {{left numbers}}: {{label}}.

Elite LLM. We adopt the same prompts used in Tree-of-Thoughts (Yao et al., 2023a) to propose and evaluate formulas, except that we concatenate the reflections from each step to avoid making repeated mistakes.

Elite prompt (for formula proposal)

Now we would like to play a game of 24. That is, given 4 numbers, try to use them with arithmetic operations (+ - \* /) to get 24.

Evaluate if given numbers can reach 24 and choose labels from 'sure', 'likely' and 'impossible'.

What you have learned about the puzzle are summarized below. {reflections}

Now use numbers and basic arithmetic operations (+ - \* /) to generate possible next steps. Make sure use steps that is sure to leads to 24 and avoid steps that are impossible to generate 24. Note that it is possible that we are considering intermediate steps so the numbers of the input may be less than 4. Example: Input: 2 8 8 14 Possible next steps: 2 + 8 = 10 (left: 8 10 14) 8 / 2 = 4 (left: 4 8 14) 14 + 2 = 16 (left: 8 8 16) 2 \* 8 = 16 (left: 8 14 16) 8 - 2 = 6 (left:  $6 \ 8 \ 14$ ) 14 - 8 = 6 (left: 2 6 8)14 / 2 = 7 (left: 7 8 8)14 - 2 = 12 (left: 8 8 12)Example: Input: 2 5 8 5 - 2 = 3 (left: 3 8) 5 \* 2 = 10 (left: 10 8)8 / 2 = 4 (left: 4 5)Now try with the following input: Input: {input} Possible next steps: {input}

#### Elite prompt (for formula evaluation)

Now we would like to play a game of 24. That is, given 4 numbers, try to use them with arithmetic operations (+ - \* /) to get 24.

Evaluate if given numbers can reach 24 and choose labels from 'sure', 'likely' and 'impossible'.

What you have learned about the puzzle are summarized below.

{reflections}

If the given numbers are already in the feedback above, just give the answer. Otherwise enumerate possible steps and try to give an approximate answer. Give the

final answer in a separated line.

{input}

#### Elite prompt (for last step formula evaluation)

Now we would like to play a game of 24. That is, given 4 numbers, try to use them with arithmetic operations (+ - \* /) to get 24.

Evaluate if given numbers can reach 24 and choose labels from 'sure', 'likely' and 'impossible'.

What you have learned about the puzzle are summarized below.

{reflections}

Use numbers and basic arithmetic operations (+ - \* /) to obtain 24. Given an input and an answer, give a judgement (sure/impossible) if the answer is correct, i.e. it uses each input exactly once and no other numbers, and reach 24.

```
Input: 4 4 6 8
Answer: (4 + 8) * (6 - 4) = 24
Judge:
sure
Input: 2 9 10 12
Answer: 2 * 12 * (10 - 9) = 24
Judge:
sure
Input: 4 9 10 13
Answer: (13 - 9) * (10 - 4) = 24
Judge:
sure
Input: 4 4 6 8
Answer: (4 + 8) * (6 - 4) + 1 = 25
Judge:
impossible
Input: 2 9 10 12
Answer: 2 * (12 - 10) = 24
Judge:
impossible
Input: 4 9 10 13
Answer: (13 - 4) * (10 - 9) = 24
Judge:
impossible
Input: {input}
Answer: {answer}
Judge:
```

For Chain-of-Thought baselines, we adopt the same methodology, and keep the original prompts except for adding reflections as below.

```
Elite prompt (for chain-of-thought proposals)
Now we would like to play a game of 24. That is, given 4 numbers, try to use them
with arithmetic operations (+ - * /) to get 24.
Evaluate if given numbers can reach 24 and choose labels from 'sure', 'likely' and
'impossible'.
What you have learned about the puzzle are summarized below.
{reflections}
Now just remember the tips from before (if any) and focus on the new task. Use
numbers and basic arithmetic operations (+ - * /) to obtain 24. Each step, you are
only allowed to choose two of the remaining numbers to obtain a new number.
Input: 4 4 6 8
Steps:
4 + 8 = 12 (left: 4 6 12)
6 - 4 = 2 (left: 2 12)
2 * 12 = 24 (left: 24)
Answer: (6 - 4) * (4 + 8) = 24
Input: 2 9 10 12
Steps:
```

```
12 * 2 = 24 (left: 9 10 24)
10 - 9 = 1 (left: 1 24)
24 * 1 = 24 (left: 24)
Answer: (12 * 2) * (10 - 9) = 24
Input: 4 9 10 13
Steps:
13 - 10 = 3 (left: 3 4 9)
9 - 3 = 6 (left: 4 6)
4 * 6 = 24 (left: 24)
Answer: 4 * (9 - (13 - 10)) = 24
Input: 1 4 8 8
Steps:
8 / 4 = 2 (left: 1 2 8)
1 + 2 = 3 (left: 3 8)
3 * 8 = 24 (left: 24)
Answer: (1 + 8 / 4) * 8 = 24
Input: 5 5 5 9
Steps:
5 + 5 = 10 (left: 5 9 10)
10 + 5 = 15 (left: 9 15)
15 + 9 = 24 \text{ (left: } 24)
Answer: ((5 + 5) + 5) + 9 = 24
Input: {input}
```

# G.2 ALFWorld

Critic LLM. For the LLM instance of the Critic, we prompt it with successful examples for each of the six task types to generate the value within each episode. Specifically, the critic LLM first determines if the value of the current state changes. If it changes, the critic LLM then generates the scalar value with the corresponding reason. We provide the following two examples in the put and cool tasks, respectively.

```
Critic prompt (for value generation), "put" task
You are in the middle of a room. Looking quickly around you, you see a cabinet 4, a
cabinet 3, a cabinet 2, a cabinet 1, a countertop 1, a garbagecan 1, a
handtowelholder 2, a handtowelholder 1, a sinkbasin 2, a sinkbasin 1, a toilet 1, a
toiletpaperhanger 1, and a towelholder 1. Your task is to: put some spraybottle on
toilet.
> critic: My task requires two sub-goals in order: take a spraybottle and put the
spraybottle on the toilet. My current state satisfies zero of the two sub-goals. The
value is 0/2=0.
> OK.
On the cabinet 1, you see a cloth 1, a soapbar 1, a soapbottle 1.
> OK.
The cabinet 2 is closed.
> OK.
You open the cabinet 2. The cabinet 2 is open. In it, you see a candle 1, and a
spraybottle 2.
```

```
> OK.
You pick up the spraybottle 2 from the cabinet 2.
> critic: Now I take a spraybottle. My current state satisfies the first of the two sub-goals: take a spraybottle. The value is 1/2=0.5.
> OK.
On the toilet 1, you see a soapbottle 2.
> OK.
You put the spraybottle 2 in/on the toilet 1.
> critic: Now I put the spraybottle on the toilet. My current state satisfies all the two sub-goals. The value is 2/2=1.
```

```
Critic prompt (for value generation), "cool" task
You are in the middle of a room. Looking quickly around you, you see a cabinet 16, a
cabinet 15, a cabinet 14, a cabinet 13, a cabinet 12, a cabinet 11, a cabinet 10, a
cabinet 9, a cabinet 8, a cabinet 7, a cabinet 6, a cabinet 5, a cabinet 4, a cabinet
3, a cabinet 2, a cabinet 1, a coffeemachine 1, a countertop 2, a countertop 1, a
diningtable 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer 1, a fridge
1, a garbagecan 1, a microwave 1, a safe 1, a sinkbasin 1, a stoveburner 4, a
stoveburner 3, a stoveburner 2, a stoveburner 1, and a toaster 1.
Your task is to: cool some pan and put it in stoveburner.
> critic: My task requires three sub-goals in order: take a pan, cool the pan with
fridge, and put the pan in the stoveburner. My current state satisfies zero of the
three sub-goals. The value is 0/3=0.
> OK.
On the stoveburner 1, you see nothing.
On the stoveburner 2, you see a pot 1.
> OK.
On the stoveburner 3, you see a pan 1.
> OK.
You pick up the pan 1 from the stoveburner 3.
> critic: Now I take a pan. My current state satisfies the first of the three
sub-goals: take a pan. The value is 1/3=0.333.
> OK.
The fridge 1 is closed.
> OK.
You cool the pan 1 using the fridge 1.
> critic: Now I cool the pan. My current state satisfies the first two of the three
sub-goals: take a pan and cool the pan with fridge. The value is 2/3=0.667.
> OK.
On the stoveburner 1, you see nothing.
You put the pan 1 in/on the stoveburner 1.
> critic: Now I put the pan in the stoveburner. My current state satisfies all the
three sub-goals. The value is 3/3=1.
```

In addition to the prompt for value generation within each episode, we also prompt the Critic LLM to generate summary feedback according to the past failure trajectory. We use

the same prompt for all six task types. Specifically, it contains the full failure trajectory (observation only), as well as the failure feedback at the end of each trajectory.

```
Critic prompt (for feedback generation)
You are in the middle of a room. Looking quickly around you, you see a cabinet 6, a
cabinet 5, a cabinet 4, a cabinet 3, a cabinet 2, a cabinet 1, a coffeemachine 1, a
countertop 3, a countertop 2, a countertop 1, a drawer 3, a drawer 2, a drawer 1, a
fridge 1, a garbagecan 1, a microwave 1, a shelf 3, a shelf 2, a shelf 1, a sinkbasin
1, a stoveburner 4, a stoveburner 3, a stoveburner 2, a stoveburner 1, and a toaster
Your task is to: put a clean plate in countertop.
> critic: My task requires three sub-goals in order: take a plate, clean the plate
with sinkbasin, and put the plate on the countertop. My current state satisfies zero
of the three sub-goals. The value is 0/3=0.
> OK.
The drawer 1 is closed.
> OK.
You open the drawer 1. The drawer 1 is open. In it, you see a creditcard 2, a knife
1, and a spatula 2.
> OK.
On the stoveburner 2, you see a pan 1.
You pick up the pan 1 from the stoveburner 2.
> critic: Now I take a pan. My current state satisfies the first of the three
sub-goals: take a plate. The value is 1/3=0.167.
> OK.
On the sinkbasin 1, you see nothing.
> OK.
You put the pan 1 in/on the sinkbasin 1.
> OK.
The microwave 1 is closed.
> OK.
You open the microwave 1. The microwave 1 is open. In it, you see a potato 1.
> OK.
On the garbagecan 1, you see a soapbottle 2.
> OK.
The microwave 1 is open. In it, you see a potato 1.
> OK.
On the coffeemachine 1, you see nothing.
> OK.
On the countertop 2, you see a bread 1, a cellphone 2, a cellphone 1, a
papertowelroll 1, a plate 2, and a soapbottle 1.
> OK.
The drawer 2 is closed.
> OK.
You open the drawer 2. The drawer 2 is open. In it, you see a spatula 1.
> OK.
On the sinkbasin 1, you see a pan 1.
> OK.
On the cabinet 3, you see a cup 1.
```

> NK

On the countertop 1, you see a apple 2, a dishsponge 2, a potato 3, and a potato 2. STATUS: FAIL

Failure feedback: In this environment, my critic assigned a 1/3 value after taking a pan. However, the task is to take and clean a plate. I noticed that the plate was found on countertop 2. In the next trial, I will go to countertop 2 to take the plate, then go to a sinkbasin to clean the plate.

You are in the middle of a room. Looking quickly around you, you see a cabinet 20, a cabinet 19, a cabinet 18, a cabinet 17, a cabinet 16, a cabinet 15, a cabinet 14, a cabinet 13, a cabinet 12, a cabinet 11, a cabinet 10, a cabinet 9, a cabinet 8, a cabinet 7, a cabinet 6, a cabinet 5, a cabinet 4, a cabinet 3, a cabinet 2, a cabinet 1, a coffeemachine 1, a countertop 3, a countertop 2, a countertop 1, a drawer 3, a drawer 2, a drawer 1, a fridge 1, a garbagecan 1, a microwave 1, a shelf 3, a shelf 2, a shelf 1, a sinkbasin 1, a stoveburner 4, a stoveburner 3, a stoveburner 2, a stoveburner 1, and a toaster 1.

Your task is to: put two saltshaker in drawer.

> critic: My task requires four sub-goals in order: take the first saltshaker, put the first saltshaker in the drawer, take the second saltshaker, and put the second saltshaker in the drawer. My current state satisfies zero of the four sub-goals. The value is 0/4=0.

> OK.

On the coffeemachine 1, you see a cup 1, a fork 4, and a potato 2.

> UK

On the toaster 1, you see a pan 1 and a knife 1.

> OK.

On the countertop 3, you see a bowl 2, a houseplant 1, a knife 2, a peppershaker 1, a saltshaker 4, a soapbottle 1, a spatula 2, and a tomato 3.

> OK.

You pick up the saltshaker 4 from the countertop 3.

> critic: Now I take the first saltshaker. My current state satisfies the first of the four sub-goals: take the first saltshaker. The value is 1/4=0.25.

> OK.

The drawer 1 is open. In it, you see nothing.

> OK.

You put the saltshaker 4 in/on the drawer 1.

> critic: Now I put the first creditcard in the dresser. My current state satisfies the first two of the four sub-goals: take the first saltshaker and put the first saltshaker in the drawer. The value is 2/4=0.5.

> OK.

On the stoveburner 1, you see a bread 1.

> OK.

On the stoveburner 2, you see a potato 1, a pot 2, and a knife 3.

> OK.

On the stoveburner 3, you see a fork 1, a fork 2, and a cup 2.

> OK.

On the stoveburner 4, you see a peppershaker 1 and a pan 2.

> OK.

You pick up the peppershaker 1 from the stoveburner 4.

> critic: Now I take the second saltshaker. My current state satisfies the first three of the four sub-goals: take the first saltshaker, put the first saltshaker in the drawer, and take the second saltshaker. The value is 3/4=0.75.

The drawer 1 is open. In it, you see a saltshaker 4.

> OK.

You put the peppershaker 1 in/on the drawer 1.

> critic: Now I put the second creditcard in the dresser. My current state satisfies all the four sub-goals. The value is 4/4=1.

STATUS: FAIL

Failure feedback: In this environment, my critic assigned a 1/4 value after taking a saltshaker and a 3/4 value after taking a peppershaker. However, the task is to put two saltshaker in the drawer. I noticed that one saltshaker was put on the countertop 3. I didn't find the second saltshaker in this trial. In the next trial, I will first go to countertop 3 to take the first saltshaker and then explore other locations to find a second saltshaker, including drawer 2-3, shelf 1-3, and cabinet 1-20.

You are in the middle of a room. Looking quickly around you, you see a cabinet 5, a cabinet 4, a cabinet 3, a cabinet 2, a cabinet 1, a coffeemachine 1, a countertop 3, a countertop 2, a countertop 1, a drawer 3, a drawer 2, a drawer 1, a fridge 1, a garbagecan 1, a microwave 1, a shelf 3, a shelf 2, a shelf 1, a sinkbasin 1, a stoveburner 4, a stoveburner 3, a stoveburner 2, a stoveburner 1, and a toaster 1. Your task is to: put a cool tomato in microwave.

> critic: My task requires three sub-goals in order: take a tomato, cool the tomato with fridge, and put the tomato in the microwave. My current state satisfies zero of the three sub-goals. The value is 0/3=0.

> OK.

On the countertop 1, you see a mug 1, and a saltshaker 1.

> OK.

On the sinkbasin 1, you see a fork 1, a glassbottle 2, a spatula 1, and a spoon 3. > OK.

You pick up the fork 1 from the sinkbasin 1.

> critic: Now I take a fork. My current state satisfies the first of the six sub-goals: take a fork. The value is 1/3=0.333.

> OK.

On the countertop 2, you see a butterknife 2, a butterknife 1, a cellphone 1, a glassbottle 1, a knife 1, a lettuce 3, a peppershaker 1, a statue 2, a statue 1, and a vase 1.

> OK.

On the countertop 3, you see a bread 1, a butterknife 3, a creditcard 2, a houseplant 1, a knife 3, a knife 2, a mug 2, a peppershaker 2, and a spatula 2.

> OK.

On the stoveburner 4, you see a pan 1.

> OK.

The drawer 3 is closed.

> OK.

You open the drawer 3. The drawer 3 is open. In it, you see a saltshaker 3.

The fridge 1 is closed.

> OK.

```
On the countertop 3, you see a bread 1, a butterknife 3, a creditcard 2, a houseplant
1, a knife 3, a knife 2, a mug 2, a peppershaker 2, and a spatula 2.
> OK.
On the cabinet 1, you see a dishsponge 1.
> OK.
The cabinet 2 is closed.
> OK.
You open the cabinet 2. The cabinet 2 is open. In it, you see nothing.
> OK.
On the cabinet 3, you see a cd 1 and a plate 4.
STATUS: FAIL
Failure feedback: In this environment, my critic assigned a 1/3 value after taking a
fork. However, the task is to take and cool a tomato. I didn't find the tomato in
this trial. In the next trial, I will explore other locations to find a tomato,
including cabinet 4, cabinet 5, coffeemachine 1, microwave 1, shelf 1-3, stoveburner
1-4 and toaster 1, etc.
```

Model LLM and Elite LLM. We use the same prompt for both the Model LLM and the Elite LLM. Specifically, we simply prompt these two instances with the successful past trajectories (observations and actions) for each task type and expect the Model LLM and Elite LLM to generate the possible next observation and the potential actions, respectively. Below, we provide two prompt examples in the put and cool tasks, respectively.

```
Model & Elite prompt, "put" task
You are in the middle of a room. Looking quickly around you, you see a cabinet 4, a
cabinet 3, a cabinet 2, a cabinet 1, a countertop 1, a garbagecan 1, a
handtowelholder 2, a handtowelholder 1, a sinkbasin 2, a sinkbasin 1, a toilet 1, a
toiletpaperhanger 1, and a towelholder 1.
Your task is to: put some spraybottle on toilet.
> go to cabinet 1
On the cabinet 1, you see a cloth 1, a soapbar 1, a soapbottle 1.
> go to cabinet 2
The cabinet 2 is closed.
> open cabinet 2
You open the cabinet 2. The cabinet 2 is open. In it, you see a candle 1, and a
spraybottle 2.
> take spraybottle 2 from cabinet 2
You pick up the spraybottle 2 from the cabinet 2.
> go to toilet 1
On the toilet 1, you see a soapbottle 2.
> put spraybottle 2 in/on toilet 1
You put the spraybottle 2 in/on the toilet 1.
```

```
Model & Elite prompt, "cool" task
```

```
You are in the middle of a room. Looking quickly around you, you see a cabinet 16, a
cabinet 15, a cabinet 14, a cabinet 13, a cabinet 12, a cabinet 11, a cabinet 10, a
cabinet 9, a cabinet 8, a cabinet 7, a cabinet 6, a cabinet 5, a cabinet 4, a cabinet
3, a cabinet 2, a cabinet 1, a coffeemachine 1, a countertop 2, a countertop 1, a
diningtable 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer 1, a fridge
1, a garbagecan 1, a microwave 1, a safe 1, a sinkbasin 1, a stoveburner 4, a
stoveburner 3, a stoveburner 2, a stoveburner 1, and a toaster 1.
Your task is to: cool some pan and put it in stoveburner.
> go to stoveburner 1
On the stoveburner 1, you see nothing.
> go to stoveburner 2
On the stoveburner 2, you see a pot 1.
> go to stoveburner 3
On the stoveburner 3, you see a pan 1.
> take pan 1 from stoveburner 3
You pick up the pan 1 from the stoveburner 3.
> go to fridge 1
The fridge 1 is closed.
> cool pan 1 with fridge 1
You cool the pan 1 using the fridge 1.
> go to stoveburner 1
On the stoveburner 1, you see nothing.
> put pan 1 in/on stoveburner 1
You put the pan 1 in/on the stoveburner 1.
```

# G.3 Blocksworld

Critic LLM. We evaluate RAFA and RAP with the reward scheme proposed by Hao et al. (2023). We prompt the language model with the previous state-action trajectory and calculate the log probabilities of taking each feasible action. Given the action taken in the current state, the Model LLM predicts the next state and we calculate the percentage of subgoals completed in the next state. We adopt the prompt examples from Hao et al. (2023) to ensure fairness in comparison.

```
Critic prompt example (for log probability), "step-4" task

I am playing with a set of blocks where I need to arrange the blocks into stacks.

Here are the actions I can do

Pick up a block
Unstack a block from on top of another block

Put down a block
Stack a block on top of another block

I have the following restrictions on my actions:
I can only pick up or unstack one block at a time.
I can only pick up or unstack a block if my hand is empty.
```

I can only pick up a block if the block is on the table and the block is clear. A block is clear if the block has no other blocks on top of it and if the block is not picked up.

I can only unstack a block from on top of another block if the block I am unstacking was really on top of the other block.

I can only unstack a block from on top of another block if the block I am unstacking is clear.

Once I pick up or unstack a block, I am holding the block.

I can only put down a block that I am holding.

I can only stack a block on top of another block if I am holding the block being stacked.

I can only stack a block on top of another block if the block onto which I am stacking the block is clear.

Once I put down or stack a block, my hand becomes empty.

## [STATEMENT]

As initial conditions I have that, the red block is clear, the yellow block is clear, the hand is empty, the red block is on top of the blue block, the yellow block is on top of the orange block, the blue block is on the table and the orange block is on the table.

My goal is to have that the orange block is on top of the red block.

My plan is as follows:

#### [PLAN]

unstack the yellow block from on top of the orange block put down the yellow block pick up the orange block stack the orange block on top of the red block [PLAN END]

## [STATEMENT]

As initial conditions I have that, the orange block is clear, the yellow block is clear, the hand is empty, the blue block is on top of the red block, the orange block is on top of the blue block, the red block is on the table and the yellow block is on the table.

My goal is to have that the blue block is on top of the red block and the yellow block is on top of the orange block.

My plan is as follows:

## [PLAN]

pick up the yellow block
stack the yellow block on top of the orange block
[PLAN END]

## [STATEMENT]

As initial conditions I have that, the red block is clear, the blue block is clear, the orange block is clear, the hand is empty, the blue block is on top of the yellow block, the red block is on the table, the orange block is on the table and the yellow block is on the table.

My goal is to have that the blue block is on top of the orange block and the yellow block is on top of the red block.

My plan is as follows:

### [PLAN]

unstack the blue block from on top of the yellow block stack the blue block on top of the orange block pick up the yellow block stack the yellow block on top of the red block [PLAN END]

#### [STATEMENT]

As initial conditions I have that, the red block is clear, the blue block is clear, the yellow block is clear, the hand is empty, the yellow block is on top of the orange block, the red block is on the table, the blue block is on the table and the orange block is on the table.

My goal is to have that the orange block is on top of the blue block and the yellow block is on top of the red block.

My plan is as follows:

## [PLAN]

unstack the yellow block from on top of the orange block stack the yellow block on top of the red block pick up the orange block stack the orange block on top of the blue block [PLAN END]

**Model LLM.** we prompt the Model LLM with few-shot examples and the current state and action. The Model LLM generates the predicted next state description. We adopt the prompt examples from Hao et al. (2023) to ensure fairness in comparison.

Model prompt template, "Pick up" action

I am playing with a set of blocks where I need to arrange the blocks into stacks. Here are the actions I can do  $\,$ 

Pick up a block

 ${\tt Unstack} \ {\tt a} \ {\tt block} \ {\tt from} \ {\tt on} \ {\tt top} \ {\tt of} \ {\tt another} \ {\tt block}$ 

Put down a block

Stack a block on top of another block

I have the following restrictions on my actions:

I can only pick up or unstack one block at a time.

I can only pick up or unstack a block if my hand is empty.

I can only pick up a block if the block is on the table and the block is clear. A block is clear if the block has no other blocks on top of it and if the block is not picked up.

I can only unstack a block from on top of another block if the block I am unstacking was really on top of the other block.

I can only unstack a block from on top of another block if the block I am unstacking is clear. Once I pick up or unstack a block, I am holding the block.

I can only put down a block that I am holding.

I can only stack a block on top of another block if I am holding the block being stacked.

I can only stack a block on top of another block if the block onto which I am stacking the block is clear. Once I put down or stack a block, my hand becomes empty.

After being given an initial state and an action, give the new state after performing the action.

#### [SCENARIO 1]

[STATE 0] I have that, the white block is clear, the cyan block is clear, the brown block is clear, the hand is empty, the white block is on top of the purple block, the purple block is on the table, the cyan block is on the table and the brown block is on the table.

[ACTION] Pick up the brown block.

[CHANGE] The hand was empty and is now holding the brown block, the brown block was on the table and is now in the hand, and the brown block is no longer clear.

[STATE 1] I have that, the white block is clear, the cyan block is clear, the brown block is in the hand, the hand is holding the brown block, the white block is on top of the purple block, the purple block is on the table and the cyan block is on the table.

#### [SCENARIO 2]

[STATE 0] I have that, the purple block is clear, the cyan block is clear, the white block is clear, the hand is empty, the white block is on top of the brown block, the purple block is on the table, the cyan block is on the table and the brown block is on the table.

[ACTION] Pick up the cyan block.

[CHANGE] The hand was empty and is now holding the cyan block, the cyan block was on the table and is now in the hand, and the cyan block is no longer clear.

[STATE 1] I have that, the cyan block is in the hand, the white block is clear, the purple block is clear, the hand is holding the cyan block, the white block is on top of the brown block, the purple block is on the table and the brown block is on the table.

Model prompt template, "Unstack" action

I am playing with a set of blocks where I need to arrange the blocks into stacks. Here are the actions I can do  $\,$ 

Pick up a block

Unstack a block from on top of another block Put down a block

Stack a block on top of another block

- I have the following restrictions on my actions:
- I can only pick up or unstack one block at a time.
- I can only pick up or unstack a block if my hand is empty.
- I can only pick up a block if the block is on the table and the block is clear. A block is clear if the block has no other blocks on top of it and if the block is not picked up.
- I can only unstack a block from on top of another block if the block I am unstacking was really on top of the other block.
- I can only unstack a block from on top of another block if the block I am unstacking is clear. Once I pick up or unstack a block, I am holding the block.
- I can only put down a block that I am holding.
- I can only stack a block on top of another block if I am holding the block being stacked.
- I can only stack a block on top of another block if the block onto which I am stacking the block is clear. Once I put down or stack a block, my hand becomes empty.

After being given an initial state and an action, give the new state after performing the action.

## [SCENARIO 1]

[STATE 0] I have that, the white block is clear, the cyan block is clear, the brown block is clear, the hand is empty, the white block is on top of the purple block, the purple block is on the table, the cyan block is on the table and the brown block is on the table.

[ACTION] Unstack the white block from on top of the purple block.

[CHANGE] The hand was empty and is now holding the white block, the white block was on top of the purple block and is now in the hand, the white block is no longer clear, and the purple block is now clear.

[STATE 1] I have that, the purple block is clear, the cyan block is clear, the brown block is clear, the hand is holding the white block, the white block is in the hand, the purple block is on the table, the cyan block is on the table and the brown block is on the table.

## [SCENARIO 2]

[STATE 0] I have that, the purple block is clear, the cyan block is clear, the white block is clear, the hand is empty, the cyan block is on top of the brown block, the purple block is on the table, the white block is on the table and the brown block is on the table.

[ACTION] Unstack the cyan block from on top of the brown block.

[CHANGE] The hand was empty and is now holding the cyan block, the cyan block was on top of the brown block and is now in the hand, the cyan block is no longer clear, and the brown block is now clear.

[STATE 1] I have that, the purple block is clear, the brown block is clear, the cyan block is in the hand, the white block is clear, the hand is holding the cyan block, the purple block is on the table, the white block is on the table and the brown block is on the table.

Model prompt template, "Put down" action

I am playing with a set of blocks where I need to arrange the blocks into stacks. Here are the actions I can do  $\,$ 

Pick up a block

Unstack a block from on top of another block

Put down a block

Stack a block on top of another block

- I have the following restrictions on my actions:
- I can only pick up or unstack one block at a time.
- I can only pick up or unstack a block if my hand is empty.
- I can only pick up a block if the block is on the table and the block is clear. A block is clear if the block has no other blocks on top of it and if the block is not picked up.
- I can only unstack a block from on top of another block if the block I am unstacking was really on top of the other block.
- I can only unstack a block from on top of another block if the block I am unstacking is clear. Once I pick up or unstack a block, I am holding the block.
- I can only put down a block that I am holding.
- I can only stack a block on top of another block if I am holding the block being stacked.
- I can only stack a block on top of another block if the block onto which I am stacking the block is clear. Once I put down or stack a block, my hand becomes empty.

After being given an initial state and an action, give the new state after performing the action.

## [SCENARIO 1]

[STATE 0] I have that, the white block is clear, the purple block is clear, the cyan block is in the hand, the brown block is clear, the hand is holding the cyan block, the white block is on the table, the purple block is on the table, and the brown block is on the table.

[ACTION] Put down the cyan block.

[CHANGE] The hand was holding the cyan block and is now empty, the cyan block was in the hand and is now on the table, and the cyan block is now clear.

[STATE 1] I have that, the cyan block is clear, the purple block is clear, the white block is clear, the brown block is clear, the hand is empty, the white block is on the table, the purple block is on the table, the cyan block is on the table, and the brown block is on the table.

## [SCENARIO 2]

[STATE 0] I have that, the purple block is clear, the black block is in the hand, the white block is clear, the hand is holding the black block, the white block is on top of the brown block, the purple block is on the table, and the brown block is on the table.

[ACTION] Put down the black block.

[CHANGE] The hand was holding the black block and is now empty, the black block was in the hand and is now on the table, and the black block is now clear.

[STATE 1] I have that, the black block is clear, the purple block is clear, the white block is clear, the hand is empty, the white block is on top of the brown block, the purple block is on the table, the brown block is on the table, and the black block is on the table.

## Model prompt template, "Stack" action

I am playing with a set of blocks where I need to arrange the blocks into stacks. Here are the actions I can do  $\,$ 

Pick up a block

Unstack a block from on top of another block

Put down a block

Stack a block on top of another block

- I have the following restrictions on my actions:
- I can only pick up or unstack one block at a time.
- I can only pick up or unstack a block if my hand is empty.
- I can only pick up a block if the block is on the table and the block is clear. A block is clear if the block has no other blocks on top of it and if the block is not picked up.
- I can only unstack a block from on top of another block if the block I am unstacking was really on top of the other block.
- I can only unstack a block from on top of another block if the block I am unstacking is clear. Once I pick up or unstack a block, I am holding the block.
- I can only put down a block that I am holding.
- I can only stack a block on top of another block if I am holding the block being stacked.
- I can only stack a block on top of another block if the block onto which I am stacking the block is clear. Once I put down or stack a block, my hand becomes empty.

After being given an initial state and an action, give the new state after performing the action.

#### [SCENARIO 1]

[STATE 0] I have that, the white block is clear, the purple block is clear, the cyan block is in the hand, the brown block is clear, the hand is holding the cyan block, the white block is on the table, the purple block is on the table, and the brown block is on the table.

[ACTION] Stack the cyan block on top of the brown block.

[CHANGE] The hand was holding the cyan block and is now empty, the cyan block was in the hand and is now on top of the brown block, the brown block is no longer clear, and the cyan block is now clear.

[STATE 1] I have that, the cyan block is clear, the purple block is clear, the white block is clear, the hand is empty, the cyan block is on top of the brown block, the brown block is on the table, the purple block is on the table, and the white block is on the table.

## [SCENARIO 2]

[STATE 0] I have that, the purple block is clear, the black block is in the hand, the white block is clear, the hand is holding the black block, the white block is on top of the brown block, the purple block is on the table, and the brown block is on the table.

[ACTION] Stack the black block on top of the purple block.

[CHANGE] The hand was holding the black block and is now empty, the black block was in the hand and is now on top of the purple block, the purple block is no longer clear, and the black block is now clear.

[STATE 1] I have that, the black block is clear, the white block is clear, the hand is empty, the black block is on top of the purple block, the white block is on top of the brown block, the brown block is on the table, and the purple block is on the table.

# G.4 Tic-Tac-Toe

#### Elite LLM

Elite prompt, propose n actions

In the game of Tic-Tac-Toe, two players, "X" and "O," alternate placing their symbols on a 3x3 grid. The objective is to be the first to get three of their symbols in a row, either horizontally, vertically, or diagonally. We use numbers to indicate empty positions, and then replace them with "X" or "O" as moves are made. For example, an empty board is denoted by

Your task is to identify the optimal position for the next move based on the current board state. Assume that it's your turn and you're playing as "{role}". Please make sure the optimal position is EMPTY. For example, in the following Tic-Tac-Toe Board:

1 | 2 | 3 ------4 | X | 6 ------7 | 8 | 9

Position 5 is occupied by "X". Thus, position 5 is not an optimal position. Provide only the optimal position in the first line. In the second line, give a brief explanation for this choice.

```
Current Tic-Tac-Toe Board:
{state}
Role: {role}
Optimal Position:
```

## Model LLM

Model prompt, predict next state

Predict the Next State of the Tic-Tac-Toe Board

In a game of Tic-Tac-Toe, two players, "X" and "O," take turns to place their symbols on a 3x3 grid. Your task is to predict what the board will look like after a specified move has been made.

 ${\tt Examples}$ 

{examples}

Now, Predict the Next State of the Following Tic-Tac-Toe Board:

Initial Tic-Tac-Toe Board:

{state}

Move: Player puts "{role}" in position {action}.

Updated Board:

Model prompt, predict opponent's action

In Tic-Tac-Toe, each player takes turns placing their respective symbols ("X" or "O") on a 3x3 board. Your task is to predict where the opponent will place their symbol based on their past moves and the current board state.

Example

Tic-Tac-Toe Board:

0 | X | 0

x | 0 | x

\_\_\_\_\_

7 | 8 | X

Opponent's Move: "O" in position 7

{examples}

Here's how the Tic-Tac-Toe board currently looks:

Tic-Tac-Toe Board:

## {state}

Given the history and current board state, where do you think the opponent will place their "{role}" next? Please make sure the output is an empty position without "X" or "O".

Opponent's Move: "{role}" in position

## Critic LLM

Critic prompt, evaluate winner

Determine the Winner in a Tic-Tac-Toe Game

In Tic-Tac-Toe, two players, "X" and "O" take turns to place their respective symbols on a 3x3 board. The first player to get three of their symbols in a row, either horizontally, vertically, or diagonally, wins the game. Your task is to evaluate the board state and determine if there is a winner.

Examples

Example

Tic-Tac-Toe Board:

0 | X | 0

x | x | x

A | A | A

0 | 0 | X

Question: Is there a winner?

Answer: Let's think step by step.

First row: 0 X O, no winner Second row: X X X, X wins

Therefore, "X" wins

Example

Tic-Tac-Toe Board:

 $X \mid 2 \mid 0$ 

-----4 | 0 | X

-----

0 | X | 9

Question: Is there a winner?

Answer: Let's think step by step. First row: X 2 0, no winner Second row: 4 0 X, no winner Third row: 0 X 9, no winner First column: X 4 0, no winner Second column: 2 0 X, no winner Thrid column: 0 X 9, no winner Main diagonal: X 0 9, no winner Anti-diagonal: 0 0 0, 0 wins Therefore, "O" wins. {examples} Now, for the Current Tic-Tac-Toe Board: Tic-Tac-Toe Board: {state} Question: Is there a winner? Answer: Let's think step by step.

Critic prompt, evaluate tie (when there is no winner)

In the game of Tic-Tac-Toe, two players alternate turns to fill a 3x3 grid with their respective symbols: "X" and "O". A board is considered "completely filled" when all nine cells of the grid contain either an 'X' or an 'O', with no empty spaces or other characters.

Examples:

 $\{\tt examples\}$ 

Now for the Current Tic-Tac-Toe Board:

Tic-Tac-Toe Board:

{state}

Is the board completely filled?

Answer: