# How CAN LLM GUIDE RL? A VALUE-BASED AP-PROACH

Anonymous authors

Paper under double-blind review

### ABSTRACT

Reinforcement learning (RL) has become the de facto standard practice for sequential decision-making problems by improving future acting policies with feedback. However, RL algorithms may require extensive trial-and-error interactions to collect useful feedback for improvement. On the other hand, recent developments in large language models (LLMs) have showcased impressive capabilities in language understanding and generation, yet they fall short in exploration and selfimprovement capabilities for planning tasks, lacking the ability to autonomously refine their responses based on feedback. Therefore, in this paper, we study how the policy prior provided by the LLM can enhance the sample efficiency of RL algorithms. Specifically, we develop an algorithm named LINVIT that incorporates LLM guidance as a regularization factor in value-based RL, leading to significant reductions in the amount of data needed for learning, particularly when the difference between the ideal policy and the LLM-informed policy is small, which suggests that the initial policy is close to optimal, reducing the need for further exploration. Additionally, we present a practical algorithm SLINVIT that simplifies the construction of the value function and employs sub-goals to reduce the search complexity. Our experiments across three interactive environments-ALFWorld, InterCode, and BlocksWorld-demonstrate that the proposed method achieves state-of-the-art success rates and also surpasses previous RL and LLM approaches in terms of sample efficiency.

029 030 031

032

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

### 1 INTRODUCTION

Trained on the web-scale corpora, Large Language Models (LLMs) have exhibited emergent capabilities and seen tremendous success across various fields, such as code development (Chen et al., 2021; Roziere et al., 2023; Li et al., 2023) and theorem proving (Yang et al., 2023b; Romera-Paredes et al., 2023). The recent advances in robotics (Huang et al., 2023b; Liang et al., 2023) and games (Wang et al., 2023a; Yuan et al., 2023; Wang et al., 2023c; Liu et al., 2023b) further highlight the potential of LLMs to build effective agents in well-designed interactive environments.

However, the reasoning and planning abilities of LLMs, which are important for intelligent agents, 040 have been found to be inconsistent and often unreliable (Valmeekam et al., 2023b; Mahowald et al., 041 2023; Huang et al., 2023a; Pallagani et al., 2023). Besides, agents powered by LLMs tend to have 042 limited abilities to explore different strategies, frequently defaulting to repeating established poli-043 cies. This limitation becomes particularly pronounced in complex decision-making scenarios that 044 LLMs are not specifically attuned to, resulting in significant difficulties in refining their strategies based on environmental feedback by reasoning the environment feedback based solely on its inherent capabilities (Valmeekam et al., 2023a; Shinn et al., 2023; Ivanova, 2023; Zhang et al., 2024). 046 On the contrary, Reinforcement Learning (RL) is a well-studied methodology for improving future 047 acting policies with feedback. Unfortunately, improving from scratch without the guidance of prior 048 knowledge, such as common sense, requires the RL agents to take a huge amount of random interactions to collect useful feedback, leading to poor sample efficiency and even failure in sparse-reward environments. 051

052

2 Hence, in this paper, we aim to tackle these issues and answer the following question:

053

Can we improve the sample efficiency of Reinforcement Learning with Large Language Models?





088

090

099 100

103

Our primary objective is to develop an algorithm that is both theoretically robust and empirically ef-073 fective, utilizing LLMs to enhance sample efficiency. Our pivotal insight is the utilization of LLMs 074 to define a regularizer, as opposed to directly employing them in decision-making. Leveraging 075 the properties of regularized-MDPs, we find that sample complexity can be significantly reduced 076 when the LLM-provided policy closely aligns with the optimal policy. Moreover, our approach 077 retains the capability to identify the optimal policy even in scenarios where the LLM policy falls 078 short. An illustration comparing the standard RL algorithms, LLM agents, and the RL framework 079 with LLM as a prior is shown in Figure 1. To demonstrate this concept, we introduce an algorithm 080 named Language-INtegrated Value Iteration (LINVIT), which shows a marked improvement in sam-081 ple complexity, particularly when the Kullback-Leibler (KL) divergence between the optimal policy 082 and the LLM policy is minimal.

We further present a practical algorithm called SLINVIT and empirically validate it in various benchmarks, including ALFWorld (Shridhar et al., 2020), the interactive coding environment InterCode (Yang et al., 2023a), and the planning benchmark BlocksWorld (Valmeekam et al., 2023b). Experimental results show that the proposed algorithm outperforms previous RL and LLM algorithms by a large margin, achieving higher success rates with fewer numbers of samples.

2 BACKGROUND

**Reinforcement Learning.** Consider the problem of learning to optimize a finite *H*-horizon Markov Decision Process (MDP) over repeated episodes of interaction. We denote by S and A the state and action space, respectively. When taking an action  $a \in A$  at a state  $s \in S$  at timestep *h*, the agent receives a reward  $r_h(s, a)$  and the MDP transits to a new state s' according to  $s' \sim P_h^*(\cdot | s, a)$ .

We aim to find a policy  $\pi$  that maps a state to an action distribution to maximize the expected cumulative reward. We denote by  $Q_h^* : S \times A \to \mathbb{R}$  and  $V_h^* : S \to \mathbb{R}$  the state-action value function and the state value function associated with  $\pi$ , respectively, which are defined as follows,

$$Q_h^*(s,a) = r_h(s,a) + \sum_{s'} P_h^*(s' \mid s,a) V_{h+1}^t(s'), \qquad V_h^*(s) = \sum_a \pi(a \mid s) Q_h^*(s,a),$$

where  $s \in S$ ,  $a \in A$ . The objective of the decision-making problem is to maximize the state value at the initial timestep  $V_1^*(s_1)$ , where  $s_1 \in S$  is the initial state.

**Decision-Making with LLMs.** To solve decision-making problems, one can prompt the LLM agent to generate action responses  $\pi_h^{\text{LLM}}(a \mid s)$  based on its state *s*, which consists of the observation-action history up to the current timestep *h* in partially observable MDPs, or contains an additional reasoning step (Yao et al., 2022) in a chain-of-thought (Wei et al., 2022) manner. Unfortunately, without the domain knowledge of a specific task or environment, the LLM policy is hard to be optimal, especially when the planning problems have long horizons and the LLM lacks
 the necessary reasoning abilities. In this work, we explore another manner of leveraging the LLM
 policy as the priorin RL.

111 112

113 114

## 3 USING LANGUAGE MODEL AS A POLICY PRIOR

In this section, we present an algorithm leveraging LLM to enhance sample efficiency. We begin by discussing the algorithm's motivation, followed by a detailed explanation of its procedure.

117 **Motivation.** A simplistic approach to using Large Language Models (LLMs) in decision-making 118 is directly applying the LLM-generated policy to target tasks. However, pretrained and fine-tuned 119 on the static datasets, LLMs are not inherently attuned to the specific interactive environments of 120 concern and are unable to modify their policies based on environmental feedback. Consequently, 121 effectively leveraging LLM information for decision-making remains an unresolved challenge. In-122 spired by the property of entropy-regularized MDP (Neu et al., 2017), our approach employs the 123 large language model as a supplemental regularizer within the original algorithm, rather than using 124 it as the primary decision-making tool. This methodology significantly improves sample efficiency 125 when the LLM's policy is closely aligned with the optimal policy. Moreover, we can still identify the optimal policy for the original MDP even if the LLM's policy is suboptimal. 126

127 128

Algorithm 1 Language-INtegrated Value Iteration(LINVIT)

129 **Input:** Target precision  $\epsilon$ , target probability  $\delta$ , bonus function  $b^0$ ,  $b^{0,\text{KL}}$ . 130 1: for t = 0, ..., T do 131 Construct the model estimator  $P_h^t$  and  $u_h^t$  as (3.1). 2: 132 Compute the optimistic and pessimistic value  $\overline{V}_{h}^{t}$  and  $\underline{V}_{h}^{t}$  as (3.2) and (3.3). 3: 133 Compute  $\pi^t$  as (3.4). 4: 134 5: for  $h = 1, \ldots, H$  do 135 6: Sample  $a_h^t \sim \pi_h^t(\cdot | s_h^t)$ , and observe  $s_{h+1}^t$  from the environment. 136 7: end for 137 8: end for 138 9: Return  $\hat{\pi}$ , which the uniform mixture of  $\{\bar{\pi}^t\}_{t=1}^T$ . 139 140

With the above motivation, we propose a novel algorithm Language-INtegrated Value Iteration 141 (LINVIT), which is an iterative algorithm that outputs a policy after T iterations. In each iteration 142  $t \in [T]$ , we first utilize the gathered data to estimate the transition model and calculate the uncer-143 tainty associated with our estimation. This estimator is employed to formulate both the optimistic 144 and pessimistic regularized value functions. The final step of each iteration involves leveraging these 145 value functions to develop an exploration policy, which is then used to acquire additional data from 146 the environment. We summarize our algorithm in Algorithm 1. In the following part, we elaborate 147 each of the above steps in detail. 148

Model and Uncertainty Estimation. We estimate the transition model as follows. Let  $n_h^t(s, a)$ denote the number of times the state-action pair (s, a) has been visited at step h during the first tepisodes, and let  $n_h^t(s, a, s')$  denotes the number of times the state-action-next-state triplet (s, a, s')at the same step and episode count. Our dynamics estimator  $P_h^t(s'|s, a)$  is defined as  $P_h^t(s'|s, a) =$  $n_h^t(s, a, s')/n_h^t(s, a)$  if  $n_h^t(s, a) > 0$  and  $P_h^t(s'|s, a) \triangleq 1/S$  for all  $s' \in S$  else. We then define the uncertainty quantifier  $u_h^t$  by

155 156

157 158

$$u_h^t(s,a) \triangleq \max\left\{2H, \sqrt{\frac{\log(4HTS^2 A/\delta)}{n_h^t(s,a)}}\right\}.$$
(3.1)

Intuitively, the uncertainty quantifier  $u_h^t$  is inversely related to the frequency of visits to a state-action pair; the less frequently a state-action pair is visited, the greater the value of  $u_h^t$ . This relationship means that  $u_h^t$  effectively measures our uncertainty regarding each state-action pair, capturing the uncertainty of our estimation. 162 **Regularized Value Functions.** After estimating the transition model and the uncertainty, we com-163 pute the optimistic and the pessimistic regularized value function by 164

$$\overline{Q}_h^t(s,a) = \operatorname{clip}\left\{r_h(s,a) + \sum_{s'} P_h^t(s' \mid s,a)\overline{V}_{h+1}^t(s') + u_h^t(s,a)\right\},$$

$$\overline{V}_{h}^{t}(s) = \max_{\pi \in \Delta_{A}} \left\{ \sum_{a} \pi(a|s) \overline{Q}_{h}^{t}(s) - \lambda \mathrm{KL} \left( \pi(\cdot \mid s) \| \pi_{h}^{\mathrm{LLM}}(\cdot \mid s) \right) \right\},$$
(3.2)

with  $\overline{V}_{H+1}^t = \underline{V}_{H+1}^t = 0$  by convention, and  $\operatorname{clip}(x) = \min\left\{\max\{x,0\},H\right\}$ . The definition of  $\overline{Q}$  and  $\overline{V}$  comprise three components: the expected reward, the uncertainty estimator, and the regularization defined via  $\pi^{\text{LLM}}$ . It can be viewed as an optimistic estimation of the regularized 172 value function. We similarly define Q and  $\underline{V}$  as

174 175 176

170

171

173

$$\underline{Q}_h^t(s,a) = \operatorname{clip}\left(r_h(s,a) + \sum_{s'} P_h^t(s' \mid s,a)\overline{V}_{h+1}^t(s') - u_h^t(s,a)\right),$$

177 178 179

$$\underline{V}_{h}^{t}(s) = \max_{\pi \in \Delta_{A}} \left\{ \sum_{a} \pi(a|s) \underline{Q}_{h}^{t}(s,a) - \lambda \mathrm{KL} \left( \pi(\cdot \mid s) \| \pi_{h}^{\mathrm{LLM}}(\cdot \mid s) \right) \right\}.$$
(3.3)

Similar to  $\overline{Q}$  and  $\overline{V}$ , Q and  $\underline{V}$  can be regarded as pessimistic estimations of the regularized value 181 function. The primary distinction between Q and  $\overline{Q}$  lies in the sign of the uncertainty estimator. 182

**Sampling Policy.** We explore the environment and collect data with  $\pi^{t+1} = {\{\pi_h^{t+1}\}}_{h=1}^H$ , which is defined as

188

189

183

$$\pi_{h}^{t}(\cdot \mid s) = \frac{1}{H} \cdot \mathbb{1}\left\{a = \operatorname{argmax} \overline{Q}_{h}^{t}(s, a) - \underline{Q}_{h}^{t}(s, a)\right\} + \frac{H - 1}{H} \cdot \bar{\pi}_{h}^{t}(\cdot \mid s),$$

$$\text{where } \bar{\pi}_{h}^{t}(\cdot \mid s) = \operatorname{argmax}_{\pi \in \Delta_{A}} \left\{\sum_{a} \pi(a \mid s) \overline{Q}_{h}^{t}(s) - \lambda \mathrm{KL}\left(\pi(\cdot \mid s) \| \pi_{h}^{\mathrm{LLM}}(\cdot \mid s)\right)\right\}.$$

$$(3.4)$$

190 Intuitively, the difference  $\overline{Q}_h^t(s,a) - \underline{Q}_h^t(s,a)$  captures the uncertainty of the estimation of the 191 regularized value function. As a result, the policy  $\pi_h^t$  is designed to act predominantly as a greedy 192 policy concerning the optimistic regularized value function, doing so with a probability of 1 - 1/H. 193 Conversely, with a probability of 1/H, it opts for the action associated with the greatest uncertainty. 194 This approach ensures a balance between exploiting known rewards and exploring actions with 195 higher uncertainty to refine the value function estimation. 196

197 198

### **RELATED WORK** 4

199 **Reinforcement Learning with Language.** Language offers a particularly effective medium for 200 tackling decision-making challenges due to its succinct and structured format. This quality has 201 made it a valuable tool for numerous reinforcement learning (RL) algorithms, enabling them to learn 202 from high-level specifications of goals (Jiang et al., 2019; Lynch & Sermanet, 2020b; Hejna et al., 203 2023) or to benefit from the step-by-step instructions provided by large language models (LLMs) 204 (Ahn et al., 2022; Huang et al., 2022). Additionally, research has ventured into harnessing more 205 expansive language applications to model the dynamics and reward mechanisms of environments, 206 employing planning algorithms to guide decision-making processes (Bialystok, 1978; Liu et al., 207 2023b). Different from these approaches, our work introduces the novel concept of applying LLMs as regularizing agents within value-based RL frameworks. 208

209 Decision-Making with Language Models. The strong capabilities language models exhibit have 210 opened a new avenue for LLM agents to interact with the real world autonomously for decision-211 making tasks. Inspired by classical planning literature (Bonet & Geffner, 2001; Hoffmann & Nebel, 212 2001; Chitnis et al., 2016; Gehring et al., 2022) that uses heuristic functions as dense reward genera-213 tors to perform informed search, recent works (Lin et al., 2023; Hao et al., 2023) proposed to use the LLM as the heuristic function. The remarkable programming abilities exhibited by the LLM have 214 also enabled converting natural language instructions into planning languages and then adopting the 215 classical planner (Liu et al., 2023a; Liang et al., 2023; Silver et al., 2023; Xie et al., 2023), which, 216 however, are constrained in narrowed domains and predefined environments. Moreover, a recent line 217 of work (Yao et al., 2023a;b; Sel et al., 2023; Zhang et al., 2023) has developed various prompting 218 schemes to enhance LLM reasoning, though these approaches generally do not integrate feedback 219 from the environment into the decision-making process.

220 A large body of previous works focused on prompt engineering by providing the LLM with addi-221 tional contextual information and templates to complete the task. Among them, the ReAct (Yao 222 et al., 2022) agent generates both reasoning traces and task-specific actions in an interleaved manner, Plan-and-Solve (Wang et al., 2023b) improves the Chain-of-Thought (Wei et al., 2022) prompt 224 to devise a fixed high-level plan before taking actions in the environment, and Huang et al. (2022) 225 prompt the LLM to extract temporally extended plans in a zero-shot manner. Besides, other works 226 directly train the model on embodied decision-making data (Suglia et al., 2021; Sharma et al., 2021; Mezghani et al., 2023; Driess et al., 2023) or multi-modal data (Lu et al., 2019; Li et al., 2019; 227 Radford et al., 2021; Zellers et al., 2021b) by learning additional downstream networks on top of the 228 pre-trained LLM (Lynch & Sermanet, 2020a; Akakzia et al., 2020; Zellers et al., 2021a) or finetun-229 ing in the environment (Reid et al., 2022; Li et al., 2022; Chen et al., 2023). Similar to our work, Ahn 230 et al. (2022); Hu & Sadigh (2023); Lin et al. (2023); Hao et al. (2023) also use value functions to 231 ground the LLM agent, but the sub-problem horizon is set to 1 and the executed actions are one-step 232 greedy without backtracking, which still suffers from the curse of the long horizon. Works building 233 on the self-reflection abilities of LLMs (Shinn et al., 2023; Sun et al., 2023; Ma et al., 2023) also 234 demonstrate limitations in refining initial strategies based on feedback, as shown in our experiments. 235

Moreover, there are various approaches that fine-tune the LLM policies to maximize the cumulative 236 reward using RL. Examples include GLAM (Carta et al., 2023) and ETPO (Wen et al., 2024), which 237 employs a token-level entropy-augmented RL method. Additionally, ILQL (Snell et al., 2022) uses 238 the LLM as a perturbation for Q-values in offline RL. In contrast, our method integrates value 239 iterations into the LLM for effective policy improvement via feedback without fine-tuning.

240 241 242

243

244

245

246

247

248

249

251

### 5 **EXPERIMENTS**

In this section, we conduct empirical studies in several text-based benchmarks, including the embodied environment ALFWorld (Shridhar et al., 2020), the interactive coding environment InterCode (Yang et al., 2023a), and the standard planning benchmark BlocksWorld (Valmeekam et al., 2023b). Across these three benchmarks, we measure the algorithm's effectiveness by its success rate, which we define as the ratio of the number of task instances the algorithm successfully completes. More precisely, each task instance is defined by a target state  $s_g \in S$ , and a task is deemed successfully completed if the algorithm reaches this target state  $s_H = s_q$  at the end.

250 We will first delve into a detailed discussion of our implementation approach. Following this, we will present and analyze the results of our experiments. 252

npı	<b>it:</b> Sub-problem horizon $N$ , BFS breadth $k$ .
1: (	Construct the value estimator $\widehat{V}$ (rule-based or Monte-Carlo)
2: <b>f</b>	for $i=1,\ldots,H/N$ do
3:	Solve $(5.1)$ with breadth-k BFS
4:	Execute the resulting $a_{(i-1)N+1:iN}$
5: e	end for

In our experiments, we introduce two primary simplifications to the LINVIT algorithm to enhance its 264 efficiency and practicality. These modifications lead to a streamlined variant we denote as SLINVIT, 265 detailed in Algorithm 2. Below, we elaborate on each simplification: 266

**Construction of Value Function and Exploration Policy.** In Algorithm 1, we use a bonus in the 267 construction of the value function and combine the uniform policy with the optimal policy in the 268 regularized MDP in the exploration policy. This approach, while sample-efficient, introduces sig-269 nificant computational complexity. To optimize computation in our experiments, we simplify these 270 processes. Specifically, we directly combine the original value estimator  $\widehat{V}$  with the log probability 271 of the LLM policy  $\mathbb{P}_{LLM}$ , which is the key component in (3.2), to construct the value estimator in the 272 experiment. Furthermore, we adopt a greedy policy with respect to this adjusted, regularized value 273 estimator. This simplification enables a simplier computation by focusing on the core elements that 274 drive the decision-making process.

275 Using Sub-Goals to Reduce Searching Complexity. Since the complexity of directly searching 276 for the maximizor of the regularized value function is exponential in the horizon H, we leverage 277 sub-goal states to reduce the searching complexity. In the experiment, our algorithm works by 278 decomposing the H-horizon planning problem into  $H/N^1$  sub-problems, each of which has a sub-279 goal and is of horizon N. More specifically, for each sub-problem  $i \in [1, H/N]$ , the corresponding 280 sub-goal  $s_{iN+1}$  is determined by solving

281 282 283

$$Q^{\text{LLM}}(s_{(i-1)N}, a_{(i-1)N+1:iN}) := \widehat{V}(s_{iN+1}) + \sum_{h=(i-1)N+1}^{iN} \lambda \pi_h^{\text{LLM}}(a_h|s_h),$$
$$a_{(i-1)N+1:iN} := \underset{a_{(i-1)N+1:iN}}{\operatorname{argmax}} Q^{\text{LLM}}(s_{(i-1)N}, a_{(i-1)N+1:iN}).$$
(5.1)

284 285

287

288

289

290

291

where  $\hat{V}$  is the estimator of the true value  $V^{\pi}$  and  $\lambda > 0$  is a hyperparameter for the regularization. Compared with the regularized value function  $\overline{V}$  in Section 3, we remove the logarithm term before  $\mathbb{P}_{\text{LLM}}$  for stability, such that it has a similar scale with  $r \in [0, 1]$ . Here,  $\widehat{V}$  can take various forms depending on both the true policy value it estimates and its own approximators, which we will discuss in more detail in Section 5.2.

In practice, we implement a breadth-k Breadth First Search (BFS) to approximate the actions in 292 (5.1). Specifically, the following procedure is repeated N times: making k copies of the agent that 293 execute the top-k outputs of the LLM by querying it "What is the potential next-step action?". This will generate  $k^N$  lookahead action sequences and the one with the highest  $Q^{\text{LLM}}$ 295 is selected as  $a_{(i-1)N+1:iN}$  and executed in the environment.

### 296 297

298 299

300

302

307

## 5.2 INSTANTIATIONS OF VALUE ESTIMATOR

In this section, we describe two instantiations of the value estimator  $\widehat{V}$  in (5.1), named rule-based and Monte-Carlo value estimators. An illustration is provided in Figure 2.

301 **Rule-Based Value Estimation.** The rule-based value estimator is designed for scenarios where achieving the goal  $s_q$  requires fulfilling multiple preconditions, such as in ALFWorld and 303 BlocksWorld. It outputs the ratio of preconditions currently met by the state s. To achieve 304 this, we prompt the LLM with "Estimate the value of the task by generating 305 Python functions" as well as the task description ahead of evaluation. To avoid uncontrol-306 lable mistakes of the LLM, its response undergoes a one-time human review. This step is necessary only once because, although  $s_q$  is changing during evaluation (e.g., "put a cup on table" and "put a pen in drawer"), the nature of the task and the structure of the preconditions 308 (e.g., "pick A" and "place A on B" for any "put" task) do not vary. We provide how we 309 implement this in Appendix C. 310



Carlo estimator. 323

<sup>&</sup>lt;sup>1</sup>For notation convenience, we assume w.l.o.g. that h can be divided by N.

324 Monte-Carlo Value Estimation. At the state  $s_h$  to be evaluated, by sampling actions from the 325 LLM policy  $\pi^{\text{LLM}}(\cdot | s_h)$  until the planning horizon H is reached, we obtain a partial trajectory. 326 The Monte-Carlo value estimation is then given by averaging the cumulative reward received in M327 such rollouts to approximate the value of the LLM policy, i.e.,

$$\widehat{V} = \frac{1}{M} \sum_{m=1}^{M/(H-h)} \sum_{n=h}^{H} r_h(s_n^m, a_n^m), \qquad \text{where } a_n^m \sim \pi_n^{\text{LLM}}(\cdot \mid s_n^m) \text{ and } s_{n+1}^m \sim P_n^*(\cdot \mid s_n^m, a_n^m).$$

### 5.3 ALFWORLD

ALFWorld (Shridhar et al., 2020) is an interactive text-based environment with aligned embodied simulators. The benchmark encompasses 134 virtual household task instances with predefined and fixed goals, each of which can be categorized into one of the six task types as shown in Table 1.



Figure 3: Demonstration of the SLINVIT algorithm in the ALFWorld environment when N = 2and the tree breadth of BFS is set to k = 3. The task is to "clean a cloth and put it on countertop". The hallucination that LLM faces, i.e., the towel should be taken (instead of cloth), is addressed by the inherent exploration mechanism in our RL framework.

We provide a visualization of how SLINVIT works in the ALFWorld environment in Figure 3. Specifically, for each type of the six tasks, we use the rule-based value estimator to generate Python code that determines the preconditions of the task goal and the current state. The code outputs the portion of the satisfied preconditions as the estimated value. We set N = 2 in our implementation for the ALFWorld benchmark.

	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
SLINVIT	100.00	100.00	91.30	90.48	100.00	100.00	97.01

Table 1: Success rate (%) comparison of SLINVIT and baselines, including LLM agents and RL algorithms, in the ALFWorld environment.

We compare the success rate of SLINVIT and baselines in the ALFWorld environment in Table 1. We use GPT-3 (text-davinci-003) in our implementation. All the LLM agent baselines, including ReAct (Yao et al., 2022), AdaPlanner (Sun et al., 2023), and Reflexion (Shinn et al., 2023), use GPT-3 that is the same as ours, while BUTLER (Shridhar et al., 2020) is a RL-style imitation learning



Figure 4: Success rates with different numbers of samples.

389 390

391

392

393

403

404

	SQL S	ample N	umber	Bash S	Sample Number			
	10	20	30	10	20	30		
TryAgain	48.45	50.97	50.97	34.67	40.20	50.25		
ReAct	52.61	52.71	53.67	20.50	21.10	21.61		
SLINVIT	52.80	58.51	64.02	46.23	50.25	54.27		

Table 2: Success rate (%) under different maximum numbers of samples per episode in the InterCode-SQL and InterCode-Bash environments.

algorithm. Notably, the inferior performance of BUTLER indicates that RL algorithms may have difficulties understanding the task and generalize beyond. On the contrary, SLINVIT achieves the highest success rate in most categories of the task types and outperforms the baselines in terms of the overall success rate.

394 We also report the changes in success rate when the numbers of samples are different. The results are shown in Figure 4. The data points of SLINVIT are obtained by changing the tree breadth when 396 performing BFS. Specifically, we set the tree breadth to  $k = 2, 3, \dots, 10$  and report the overall 397 success rate corresponding to different numbers of samples in all the 134 tasks. For the Reflexion 398 and ReAct baselines, the points in the plot are the results at the end of each trial. The sample number 399 is calculated as the number of samples taken in the successful tasks in the current trial, plus all the 400 samples in the previous trials. We observe that the proposed algorithm is able to achieve a higher 401 success rate with fewer numbers of samples compared to methods that incorporate the environment feedback summary into the LLM as additional context. 402

## 5.4 INTERCODE

InterCode (Yang et al., 2023a) is an interactive coding benchmark with code or command as actions and the environment feedback after executing an action as observations. It provides two benchmarks to evaluate the planning abilities of the large language models, namely InterCode-SQL and InterCode-Bash which use SQL and Bash commands as action spaces, respectively. For each benchmark, there are hundreds of tasks with predefined goals, such as "*Find the name of airports which do not have any flight in and out.*" and "*Find all text files in the testbed directory and subdirectories and concatenate them into a single file.*".

Unlike the ALFWorld environment where the task goals can be described by preconditions, there is no straightforward way to directly measure the value of the current state with explicit and simple rules in the InterCode environment. Therefore, we implement SLINVIT using the Monte-Carlo value estimator. Besides, we use the original dense reward as proposed in (Yang et al., 2023a). We set the sub-problem horizon N = 1 and the Monte-Carlo sampling number M = 1. For our method and all the baselines, we use GPT-3.5 (gpt-3.5-turbo).

	Int				
	Easy	Medium	Hard	Extra	Total
M = 1	90.73	71.08	60.34	50.00	70.60
M=2	88.31	77.13	65.52	52.42	73.89

Table 3: Ablation study on	SLINVIT with different $M$ .
----------------------------	------------------------------

The success rates in the InterCode-SQL and InterCode-Bash environments are reported in Table 4.
The baselines we compare include ReAct (Yao et al., 2022), Plan & Solve (Wang et al., 2023b), and
Try Again (Yang et al., 2023a), which is a vanilla LLM-based planning algorithm. We observe that
SLINVIT achieves the highest success rate in all the hardness modes of the InterCode-SQL benchmark, all the file systems of the InterCode-Bash benchmark, and has the best overall performance.
In Table 2, we investigate the sample efficiency of the proposed algorithm in the InterCode environment. Specifically, we set the maximum number of samples in each episode to be 10, 20, and 30

		InterCode-SQL Hardness			InterCode-Bash File System				n		
	Ea	isy	Med.	Hard	Extra	Total	Sys 1	Sys 2	Sys 3	Sys 4	Total
TryAgain	75	5.81	48.65	49.43	21.69	50.97	45.00	49.06	45.00	48.15	46.50
Plan & Sol	ve 77	.42	49.78	32.18	22.89	49.13	0.00	45.28	43.33	22.22	28.00
ReAct	80	).24	65.47	47.13	20.48	58.70	21.67	5.66	30.00	25.93	20.50
SLINVIT	90	.73	71.08	60.34	50.00	70.60	55.00	60.38	64.41	66.67	60.80

Table 4: Success rate (%) comparison in the InterCode-SQL and InterCode-Bash environment.

and report the corresponding success rates. For SLINVIT, both the samples taken to maximize (5.1) and the samples for Monte-Carlo value estimation are counted. The results indicate that SLINVIT consistently outperforms the baselines with the same sample size and is thus more sample-efficient.

Ablation. We conduct an ablation study on the number of rollouts M and the results are shown in Table 3. We observe that a more accurate value estimation corresponding to a larger M leads to higher success rates for harder problems. Therefore, a trade-off can be taken between the sample number and the performance.

5.5 BLOCKSWORLD

BlocksWorld (Valmeekam et al., 2023b; Liu et al., 2023a) is another planning benchmark that contains various tasks to arrange blocks in specific configurations. The state is the current configuration of the blocks and the action is an instruction that moves blocks. Specifically, an action is composed of one of the four verbs (STACK, UNSTACK, PUT, and PICKUP) and the operated block. Similar to the implementation in ALFWorld, we also adopt the rule-based value estimator. Specifically, the goal state of each task is defined as the combination of several block arrangements (e.g., "block 1 is on top of block 2"). With the current state as input, the value estimator then returns the proportion of the satisfied arrangements.

Following RAP (Hao et al., 2023), we group the task instances in Valmeekam et al. (2023b) by the minimum number of actions required, resulting in 57 cases that are solvable within 4 steps, and 114 cases that are solvable within 6 steps. We evaluate our method and the RAP baseline by comparing the success rates under different numbers of samples. We evaluate the Vicuna-13b (v1.3) model and the results are shown in Figure 5. Our method consistently outperforms RAP and achieves a higher success rate with fewer samples.



Figure 5: Success rate (%) of SLINVIT and baselines on the 4-step and 6-step BlocksWorld tasks.

## 6 THEORY



486 **Definition 6.1** (KL-divergence between two policy). For two policies  $\pi_1 = {\{\pi_1^h\}_{h=1}^H}$  and  $\pi_2 = {\{\pi_2^h\}_{h=1}^H}$ , we define

489 490

491 492

493

494

495

496 497 498

500 501

502

$$\mathrm{KL}(\pi^{1} \| \pi^{2}) = \sum_{h=1}^{H} \mathbb{E}_{\pi^{1}} \Big[ \mathrm{KL} \big( \pi_{h}^{1}(\cdot | s_{h}) \| \pi_{h}^{2}(\cdot | s_{h}) \big) \Big].$$

Definition 6.1 provides a quantitative measure of the similarity between two policies. More specifically, the divergence is small when two policies are similar. Building on this foundational understanding, we present the following theorem.

**Theorem 6.2.** We assume that  $KL(\pi^* || \pi^{LLM}) \leq \epsilon_{LLM}$ , and set the tuning parameter

$$\lambda = \epsilon / (2\epsilon_{\rm LLM}), \ T = CH^6 SA^4 \log^2(HSA/\delta)\epsilon_{\rm LLM}/\epsilon^2$$

for some absolute constant C. We then have  $V_1^*(s_1) - V_1^{\hat{\pi}}(s_1) \le \epsilon$  with probability as least  $1 - \delta$ .

*Proof.* See Appendix §A for a detailed proof.

503 Theorem 6.2 demonstrates that the number of samples required to achieve  $\epsilon$ -optimality is proportional to the KL divergence,  $\text{KL}(\pi^* || \pi^{\text{LLM}})$ , given that  $\lambda$  is suitably chosen. This relationship implies a reduced sample necessity when  $\pi^{\text{LLM}}$  closely aligns with the optimal policy  $\pi^*$ . The in-504 505 tuitive rationale behind this is that the demand for exploration diminishes when an initial policy is 506 nearly optimal. Consequently, this theorem underscores the efficacy of our algorithm in capitaliz-507 ing on the policy information provided by the Large Language Model (LLM), thereby validating its 508 practical utility in decision-making scenarios. Theorem 6.2 further demonstrates that Algorithm 1 509 is capable of achieving  $\epsilon$ -optimality, even in cases where  $\epsilon \leq \epsilon_{\rm LLM}$ , contingent upon the collection 510 of a sufficient number of samples. This finding underscores the algorithm's robustness in attaining 511 a specified level of optimality. 512

In Theorem 6.2, the regularization coefficient  $\lambda$  became bigger as the KL divergence become smaller. This trend aligns intuitively with the principle of relying more heavily on the information provided by the Large Language Model (LLM) when there is evidence that the policy it offers is effective. Essentially, a smaller KL divergence indicates a closer alignment between the LLM's policy and the optimal policy, justifying increased reliance on the LLM's guidance in these scenarios.

Connection Between the Theoretical Analysis and SLINVIT. Although SLINVIT differs from the algorithm described in our theoretical analysis in several respects, it is crucial to emphasize that the core objective of our theoretical analysis is to justify the use of log-probability as regularization. Despite the simplifications made for practical implementation, SLINVIT retains log-probability regularization. This adherence ensures that the fundamental element of our analysis is preserved.

523

### 7 CONCLUSION

524 525

Large language models (LLMs) have shown remarkable capabilities in quickly generating viable 526 initial strategies for decision-making tasks, even with minimal or no prior examples. However, these 527 LLM-driven agents struggle to iteratively refine their strategies based on feedback from their envi-528 ronment, mainly because they lack the ability to effectively explore and reason from the feedback. In 529 contrast, reinforcement learning (RL) excels at adapting and improving through feedback. However, 530 it often requires an extensive amount of trial-and-error to gather improvable feedback, hindered by 531 its inability to leverage common sense reasoning. To improve the sample efficiency of RL algo-532 rithms, in this work, we propose a novel RL framework with LLM as a policy prior. We prove that 533 the number of samples required by our algorithm is proportional to the KL divergence between the 534 LLM and the optimal policy. This result is further evidenced through experiments in interactive en-535 vironments such as ALFWorld, InterCode, and BlocksWorld, underscoring our method's improved 536 sample efficiency. For future work, we would like to extend our experiments to more complex and 537 diverse environments would test the scalability and robustness of our framework. Moreover, exploring how various model architectures and pre-training tasks in LLMs influence reinforcement 538 learning performance could provide valuable insights, helping to identify key factors that enhance or hinder learning efficiency and generalization.

## 540 REFERENCES

548

553

560

561

562

542	Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea
543	Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say:
544	Grounding language in robotic affordances. arXiv preprint arXiv:2204.01691, 2022.

- Ahmed Akakzia, Cédric Colas, Pierre-Yves Oudeyer, Mohamed Chetouani, and Olivier Sigaud.
   Grounding language to autonomously-acquired skills via goal generation. arXiv preprint arXiv:2006.07185, 2020.
- Ellen Bialystok. A theoretical model of second language learning 1. Language learning, 28(1):
   69–83, 1978.
- Blai Bonet and Héctor Geffner. Planning as heuristic search. Artificial Intelligence, 129(1-2):5–33, 2001.
- Thomas Carta, Clément Romac, Thomas Wolf, Sylvain Lamprier, Olivier Sigaud, and Pierre-Yves
   Oudeyer. Grounding large language models in interactive environments with online reinforcement
   learning. In *International Conference on Machine Learning*, pp. 3676–3713. PMLR, 2023.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
  - Xiaoyu Chen, Shenao Zhang, Pushi Zhang, Li Zhao, and Jianyu Chen. Asking before action: Gather information in embodied decision making with language models. *arXiv preprint arXiv:2305.15695*, 2023.
- Rohan Chitnis, Dylan Hadfield-Menell, Abhishek Gupta, Siddharth Srivastava, Edward Groshev, Christopher Lin, and Pieter Abbeel. Guided search for task and motion plans using learned heuristics. In 2016 IEEE International Conference on Robotics and Automation (ICRA), pp. 447–454.
  IEEE, 2016.
- Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multi-modal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- 572 Clement Gehring, Masataro Asai, Rohan Chitnis, Tom Silver, Leslie Kaelbling, Shirin Sohrabi,
  573 and Michael Katz. Reinforcement learning for classical planning: Viewing heuristics as dense
  574 reward generators. In *Proceedings of the International Conference on Automated Planning and*575 *Scheduling*, volume 32, pp. 588–596, 2022.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*, 2023.
- Joey Hejna, Pieter Abbeel, and Lerrel Pinto. Improving long-horizon imitation through instruction
   prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 7857–7865, 2023.
- Jörg Hoffmann and Bernhard Nebel. The ff planning system: Fast plan generation through heuristic search. *Journal of Artificial Intelligence Research*, 14:253–302, 2001.
- Hengyuan Hu and Dorsa Sadigh. Language instructed reinforcement learning for human-ai coordination. *arXiv preprint arXiv:2304.07297*, 2023.
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. Large language models cannot self-correct reasoning yet. arXiv preprint arXiv:2310.01798, 2023a.
- Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot
   planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*, pp. 9118–9147. PMLR, 2022.

594 595 596	Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer: Composable 3d value maps for robotic manipulation with language models. <i>arXiv preprint</i>
597	arXiv:2307.03973, 20230.
598	Anna A Ivanova. Running cognitive evaluations on large language models: The do's and the don'ts.
599	arXiv preprint arXiv:2312.01276, 2023.
600	
601	Yiding Jiang, Shixiang Shane Gu, Kevin P Murphy, and Chelsea Finn. Language as an abstrac-
602	tion for hierarchical deep reinforcement learning. Advances in Neural Information Processing
603	<i>Systems</i> , <i>32</i> , <i>2019</i> .
604	Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang, Visualbert: A simple
605	and performant baseline for vision and language. arXiv preprint arXiv:1908.03557, 2019.
606	
607	Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou,
608	Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with
609	you! arxiv preprint arxiv:2305.06161, 2023.
610	Shuang Li, Xavier Puig, Chris Paxton, Yilun Du, Clinton Wang, Linxi Fan, Tao Chen, De-An Huang,
611	Ekin Akyürek, Anima Anandkumar, et al. Pre-trained language models for interactive decision-
612	making. Advances in Neural Information Processing Systems, 35:31199–31212, 2022.
613	Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and
614	Andy Zeng. Code as policies: Language model programs for embodied control. In 2023 IEEE
615	International Conference on Robotics and Automation (ICRA), pp. 9493–9500. IEEE, 2023.
616	
017	Kevin Lin, Christopher Agia, Toki Migimatsu, Marco Pavone, and Jeannette Bong. Text2motion:
610	From natural language instructions to reasible plans. arXiv preprint arXiv:2505.12155, 2025.
620	Bo Liu, Yuqian Jiang, Xiaohan Zhang, Qiang Liu, Shiqi Zhang, Joydeep Biswas, and Peter Stone.
621	Llm+ p: Empowering large language models with optimal planning proficiency. arXiv preprint
622	<i>arXiv:2304.11477</i> , 2023a.
623	
624	Zninan Liu, Hao Hu, Snenao Znang, Hongyi Guo, Snuqi Ke, Boyi Liu, and Znaoran wang. Reason
625	efficiency. arXiv preprint arXiv:2309.17382, 2023b.
626	The Discount of the Discount o
627	Jiasen Lu, Diruv Batra, Devi Parikin, and Stelan Lee. Vilbert: Pretraining task-agnostic visionin-
628 629	systems, 32, 2019.
630	Corey Lynch and Pierre Sermanet, Grounding language in play, arXiv preprint arXiv:2005.07648
631	3. 2020a.
632	-,
633	Corey Lynch and Pierre Sermanet. Language conditioned imitation learning over unstructured data.
634	arXiv preprint arXiv:2005.07648, 2020b.
635	Vashang Jasan Ma, William Liang, Guanzhi Wang, Da An Huang, Oshart Pastani, Dinash Jayara
636	man Yuke Zhu Linyi Fan and Anima Anandkumar. Fureka: Human-level reward design via
637	coding large language models. arXiv preprint arXiv:2310 12931 2023
638	
639	Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and
640	Evelina Fedorenko. Dissociating language and thought in large language models: a cognitive
641	perspective. arXiv preprint arXiv:2301.06627, 2023.
642	Pierre Ménard Omar Darwiche Domingues Anders Jonsson Emilie Kaufmann Edouard Laurent
643	and Michal Valko. Fast active learning for pure exploration in reinforcement learning. In Inter-
644	national Conference on Machine Learning, pp. 7599–7608. PMLR. 2021.
645	J
646	Lina Mezghani, Piotr Bojanowski, Karteek Alahari, and Sainbayar Sukhbaatar. Think before

you act: Unified policy for interleaving language reasoning with actions. *arXiv preprint arXiv:2304.11063*, 2023.

685

686

687

688

689

- 648
   649
   649
   650
   650
   651
   652
   653
   654
   654
   655
   655
   656
   656
   657
   658
   658
   659
   659
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
   650
- Vishal Pallagani, Bharath Muppasani, Keerthiram Murugesan, Francesca Rossi, Biplav Srivastava,
   Lior Horesh, Francesco Fabiano, and Andrea Loreggia. Understanding the capabilities of large
   language models for automated planning. *arXiv preprint arXiv:2305.16151*, 2023.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Machel Reid, Yutaro Yamada, and Shixiang Shane Gu. Can wikipedia help offline reinforcement learning? *arXiv preprint arXiv:2201.12122*, 2022.
- Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog,
   M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang,
   Omar Fawzi, et al. Mathematical discoveries from program search with large language models.
   *Nature*, pp. 1–3, 2023.
- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi
   Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code.
   *arXiv preprint arXiv:2308.12950*, 2023.
- Bilgehan Sel, Ahmad Al-Tawaha, Vanshaj Khattar, Lu Wang, Ruoxi Jia, and Ming Jin. Algorithm of thoughts: Enhancing exploration of ideas in large language models. *arXiv preprint* arXiv:2308.10379, 2023.
- Pratyusha Sharma, Antonio Torralba, and Jacob Andreas. Skill induction and planning with latent language. *arXiv preprint arXiv:2110.01517*, 2021.
- Noah Shinn, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and
   Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *arXiv preprint arXiv:2303.11366*, 2023.
- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. *arXiv preprint arXiv:2010.03768*, 2020.
- Tom Silver, Soham Dan, Kavitha Srinivas, Joshua B Tenenbaum, Leslie Pack Kaelbling, and
   Michael Katz. Generalized planning in pddl domains with pretrained large language models.
   *arXiv preprint arXiv:2305.11014*, 2023.
  - Charlie Snell, Ilya Kostrikov, Yi Su, Mengjiao Yang, and Sergey Levine. Offline rl for natural language generation with implicit language q learning. *arXiv preprint arXiv:2206.11871*, 2022.
  - Alessandro Suglia, Qiaozi Gao, Jesse Thomason, Govind Thattai, and Gaurav Sukhatme. Embodied bert: A transformer model for embodied, language-guided visual task completion. *arXiv preprint arXiv:2108.04927*, 2021.
- Haotian Sun, Yuchen Zhuang, Lingkai Kong, Bo Dai, and Chao Zhang. Adaplanner: Adaptive planning from feedback with language models. *arXiv preprint arXiv:2305.16653*, 2023.
- Karthik Valmeekam, Matthew Marquez, and Subbarao Kambhampati. Can large language models
   really improve by self-critiquing their own plans? *arXiv preprint arXiv:2310.08118*, 2023a.
- Karthik Valmeekam, Sarath Sreedharan, Matthew Marquez, Alberto Olmo, and Subbarao Kambhampati. On the planning abilities of large language models (a critical investigation with a proposed benchmark). *arXiv preprint arXiv:2302.06706*, 2023b.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023a.

702 703 704	Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. <i>arXiv preprint arXiv:2305.04091</i> , 2023b.
706 707 708 709	Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. <i>arXiv preprint arXiv:2302.01560</i> , 2023c.
710 711 712 713	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837, 2022.
714 715 716	Muning Wen, Cheng Deng, Jun Wang, Weinan Zhang, and Ying Wen. Entropy-regularized token- level policy optimization for large language models. <i>arXiv preprint arXiv:2402.06700</i> , 2024.
717 718 719	Yaqi Xie, Chen Yu, Tongyao Zhu, Jinbin Bai, Ze Gong, and Harold Soh. Translating natural lan- guage to planning goals with large-language models. <i>arXiv preprint arXiv:2302.05128</i> , 2023.
720 721 722 722	John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. Intercode: Standardizing and benchmarking interactive coding with execution feedback. <i>arXiv preprint arXiv:2306.14898</i> , 2023a.
724 725 726	Kaiyu Yang, Aidan M Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan Prenger, and Anima Anandkumar. Leandojo: Theorem proving with retrieval-augmented language models. <i>arXiv preprint arXiv:2306.15626</i> , 2023b.
728 729 730	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. <i>arXiv preprint arXiv:2210.03629</i> , 2022.
731 732 733 734	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. <i>arXiv</i> preprint arXiv:2305.10601, 2023a.
735 736 737	Yao Yao, Zuchao Li, and Hai Zhao. Beyond chain-of-thought, effective graph-of-thought reasoning in large language models. <i>arXiv preprint arXiv:2305.16582</i> , 2023b.
738 739 740 741	Haoqi Yuan, Chi Zhang, Hongcheng Wang, Feiyang Xie, Penglin Cai, Hao Dong, and Zongqing Lu. Plan4mc: Skill reinforcement learning and planning for open-world minecraft tasks. <i>arXiv</i> preprint arXiv:2303.16563, 2023.
742 743 744 745	Rowan Zellers, Ari Holtzman, Matthew Peters, Roozbeh Mottaghi, Aniruddha Kembhavi, Ali Farhadi, and Yejin Choi. Piglet: Language grounding through neuro-symbolic interaction in a 3d world. <i>arXiv preprint arXiv:2106.00188</i> , 2021a.
746 747 748 749	Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and Yejin Choi. Merlot: Multimodal neural script knowledge models. <i>Advances in Neural Information Processing Systems</i> , 34:23634–23651, 2021b.
750 751 752 753	Wenqi Zhang, Yongliang Shen, Linjuan Wu, Qiuying Peng, Jun Wang, Yueting Zhuang, and Weim- ing Lu. Self-contrast: Better reflection through inconsistent solving perspectives. <i>arXiv preprint</i> <i>arXiv:2401.02009</i> , 2024.
754 755	Yifan Zhang, Jingqin Yang, Yang Yuan, and Andrew Chi-Chih Yao. Cumulative reasoning with large language models. <i>arXiv preprint arXiv:2308.04371</i> , 2023.

### **PROOF OF THEOREM 6.2** А

*Proof.* To analyze the property of LINVIT, we introduce the definition of the prior-regularized value function. The prior-regularized value function  $Q^{\pi}_{\text{LLM},\lambda,h}$ ,  $V^{\pi}_{\text{LLM},\lambda,h}(s_h)$  and  $V^*_{\text{LLM},\lambda,h}$  are defined as

$$Q_{\mathrm{LLM},\lambda,h}^{\pi}(s_h, a_h) = r_h(s_h, a_h) + \mathbb{E}_{s_{h+1} \sim P_h} \left[ V_{\mathrm{LLM},\lambda,h+1}^{\pi}(s_{h+1}) \right],$$
  

$$V_{\mathrm{LLM},\lambda,h}^{\pi}(s_h) = \mathbb{E}_{a_h \sim \pi_h(\cdot|s_h)} Q_{\mathrm{LLM},\lambda,h}^{\pi}(s_h, a_h) - \lambda \mathrm{KL} \left( \pi_h(\cdot|s_h) \| \pi_h^{\mathrm{LLM}}(\cdot|s_h) \right).$$
  

$$V_{\mathrm{LLM},\lambda,h}^{\pi}(s_h) = \max V_{\mathrm{LLM},\lambda,h}^{\pi}(s_h).$$

LINVIT can be viewed as an algorithm that maximizes the prior-regularized value function by in-teracting with the environment. The prior-regularized value function can be viewed as a regularized version of the original value functions that favors the policies that similar with  $\pi^{\text{LLM}}$ . The following lemma connects the KL-regularized value with the original value. 

**Lemma A.1.** When 
$$KL(\pi^* || \pi^{LLM}) \leq \epsilon_{LLM}$$
, we have

$$V_1^*(s_1) - V_1^{\overline{\pi}}(s_1) \le V_{\mathrm{LLM},\lambda,h}^*(s_1) - V_{\mathrm{LLM},\lambda,h}^{\overline{\pi}}(s_1) + \lambda \epsilon_{\mathrm{LLM}}.$$

Here  $\operatorname{KL}(\pi^1 \| \pi^2) = \sum_{h=1}^{H} \mathbb{E}_{\pi^1} [\operatorname{KL}(\pi_h^1(\cdot | s_h) \| \pi_h^2(\cdot | s_h))].$ 

*Proof.* See §B.1 for a detailed proof.

The following lemma show that, our algorithm provably find the optimal policy with respect to the prior-regularized value function with high probability. 

**Lemma A.2.** With probability  $1 - \delta$ , we have

$$V_{\text{LLM},\lambda,h}^{\star}(s_1) - V_{\text{LLM},\lambda,h}^{\widehat{\pi}}(s_1) \le \epsilon/2.$$

*Proof.* See §B.2 for a detailed proof.

We denote by  $\mathcal{E}_1$  the event in Lemma A.2. We then have  $P(\mathcal{E}_1) \geq 1 - \delta$ . Since we set  $\lambda =$  $\epsilon/(2\epsilon_{\rm LLM})$ , we have

$$V_1^*(s_1) - V_1^{\widehat{\pi}}(s_1) \le V_{\mathrm{LLM},\lambda,h}^*(s_1) - V_{\mathrm{LLM},\lambda,h}^{\widehat{\pi}}(s_1) + \epsilon/2 \le \epsilon$$

by Lemma A.1 when we condition on Event  $\mathcal{E}_1$ . Therefore, we conclude the proof of Theorem 6.2.

### **PROOF OF AUXILIARY LEMMAS** В

B.1 PROOF OF LEMMA A.1 

*Proof.* By the definitions of  $V_{\pi^{\text{LLM}},\lambda,h}^{\pi^*}$  and  $V_{\pi^{\text{LLM}},\lambda,h}^{\star}(s_1)$ , we have

$$V_1^*(s_1) = V_{\pi^{\mathrm{LLM}},\lambda,h}^{\pi^*}(s_1) + \lambda \mathrm{KL}(\pi^*, \pi^{\mathrm{LLM}})$$
$$= V_{\pi^{\mathrm{LLM}},\lambda,h}^*(s_1) + \lambda \mathrm{KL}(\pi^*, \pi^{\mathrm{LLM}}).$$

Therefore, when  $\mathrm{KL}(\pi^*,\pi^{\mathrm{LLM}}) \leq \epsilon_{\mathrm{LLM}}$ , we have

$$V_1^*(s_1) - V_1^{\widehat{\pi}}(s_1) \leq V_{\pi^{\text{LLM}},\lambda,h}^*(s_1) - V_{\pi^{\text{LLM}},\lambda,h}^{\widehat{\pi}}(s_1) + \lambda \text{KL}(\pi^*, \pi^{\text{LLM}})$$
$$\leq V_{\pi^{\text{LLM}},\lambda,h}^*(s_1) - V_{\pi^{\text{LLM}},\lambda,h}^{\widehat{\pi}}(s_1) + \lambda \epsilon_{\text{LLM}},$$

which concludes the proof of Lemma A.1.

### 810 B.2 PROOF OF LEMMA A.2

*Proof.* The proof of Lemma A.2 consists of two parts. We decompose the regret in the first part,and upper bound the decomposed regret in the second part.

By the definition of  $\hat{\pi}$  in Algorithm LINVIT, we have

$$V_{\pi^{\text{LLM}},\lambda,1}^{\star}(s) - V_{\pi^{\text{LLM}},\lambda,1}^{\overline{\pi}^{t+1}}(s) = \frac{1}{T} \sum_{t=1}^{T} \left[ V_{\pi^{\text{LLM}},\lambda,1}^{\star}(s) - V_{\pi^{\text{LLM}},\lambda,1}^{\overline{\pi}^{t+1}}(s) \right].$$
(B.1)

819 We also have the following lemma.

Lemma B.1. We have

$$V_{\pi^{\mathrm{LLM}},\lambda,1}^{\star}(s) - V_{\pi^{\mathrm{LLM}},\lambda,1}^{\overline{\pi}^{t+1}}(s) \le \frac{A}{2\lambda} \sum_{h=1}^{H} \mathbb{E}_{\overline{\pi}^{t+1}} \Big[ \max_{a \in \mathcal{A}} \big( \overline{Q}_{h}^{t}(s_{h},a) - \underline{Q}_{h}^{t}(s_{h},a) \big)^{2} \Big]$$

holds for all  $(t, s) \in [T] \times S$  with probability  $1 - \delta/2$ .

*Proof.* See §B.4 for a detailed proof.

Therefore, we have

$$V_{\pi^{\text{LLM}},\lambda,1}^{\star}(s) - V_{\pi^{\text{LLM}},\lambda,1}^{\overline{\pi}^{t+1}}(s) \le \frac{A}{2T\lambda} \sum_{t=1}^{T} \sum_{h=1}^{H} \mathbb{E}_{\overline{\pi}^{t+1}} \Big[ \max_{a \in \mathcal{A}} \big( \overline{Q}_{h}^{t}(s_{h},a) - \underline{Q}_{h}^{t}(s_{h},a) \big)^{2} \Big]$$
(B.2)

The following lemma upper bounds the decomposed regret.

Lemma B.2. We have

$$\sum_{t=1}^{T} \sum_{h=1}^{H} \mathbb{E}_{\overline{\pi}^{t+1}} \left[ \max_{a \in \mathcal{A}} \left( \overline{Q}_h^t(s_h, a) - \underline{Q}_h^t(s_h, a) \right)^2 \right] \le 920SA^3H^6 \log^2(12HTS^2A/\delta)$$

holds with probability at least  $1 - \delta/4$ .

*Proof.* See §B.5 for a detailed proof.

Therefore, we have

$$V_{\pi^{\text{LLM}},\lambda,1}^{\star}(s) - V_{\pi^{\text{LLM}},\lambda,1}^{\overline{\pi}^{t+1}}(s) \le 460SA^4H^6\log^2(12HTS^2A/\delta)/(T\lambda)$$
(B.3)

when the event in Lemmas B.1 and B.2 hold. Therefore, when  $\lambda = \epsilon/(2\epsilon_{\text{LLM}})$  and  $T = CSA^4H^6\log^2(HTSA/\delta)\epsilon_{\text{LLM}}/\epsilon^2$  for some absolute constant c, we have  $V_{\pi^{\text{LLM}},\lambda,1}^{\star}(s) - V_{\pi^{\text{LLM}},\lambda,1}^{\pi^{t+1}}(s) \leq \epsilon$ , which concludes the proof of Lemma A.1.

### B.3 PROOF OF LEMMA B.6

*Proof.* First, we have the following lemma, which shows that the bonus we define characterizes the uncertainty of the model estimation with high probability.

Lemma B.3. We have

$$\mathbb{1}_{n_h^t(s,a)>0} \operatorname{TV} \left( P_h^*(\cdot \mid s,a), P_h^t(\cdot \mid s,a) \right) \le A \sqrt{\frac{\log(4HTS^2 A/\delta)}{n_h^t(s,a)}}$$

holds for all  $(t, h, s, a) \in [T] \times [H] \times S \times A$  with probability at least  $1 - \delta/2$ .

*Proof.* See §B.6 for a detailed proof.

We denote by  $\mathcal{E}_3$  the event in Lemma B.3. In the following part of the proof, we condition on Event  $\mathcal{E}_3$ . We prove Lemma B.6 using induction on h. By the definition of  $\overline{Q}_{H+1}^{\iota}$  and  $Q_{\text{LLM},\lambda,H+1}^{\star}(s,a)$ , Lemma B.6 holds when h = H + 1. We also have the following lemma, which shows that the prior-regularized value function is bounded. 

**Lemma B.4** (Boundedness of  $Q_{\text{LLM},\lambda,h}^{\star}$ ). We have

$$0 \le Q^{\star}_{\mathrm{LLM},\lambda,h}(s,a) \le H + 1 - h; \qquad 0 \le V^{\star}_{\mathrm{LLM},\lambda,h}(s) \le H + 1 - h;$$

holds for all  $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H + 1]$ .

Proof. See §B.7 for a detailed proof.

By Lemma B.4, Lemma B.6 obviously holds when  $\overline{Q}_h^t \ge H$ . Otherwise, we have  $\overline{Q}_{h}^{t}(s,a) - Q_{\text{LLM }\lambda h}^{\star}(s,a) = \sum_{k} P_{h}^{t}(s' \mid s,a) \overline{V}_{h+1}^{t}(s') - \sum_{k} P_{h}^{\star}(s' \mid s,a) V_{\text{LLM }\lambda h+1}^{\star}(s') + u_{h}^{t}(s,a)$ 

$$\sum_{s' \in \mathcal{S}} [P_h^t(s' \mid s, a) - P_h^*(s' \mid s, a)] V_{\text{LLM},\lambda,h+1}^\star(s') + u_h^t(s, a)$$

$$\geq \sum_{s' \in \mathcal{S}} [P_h^t(s' \mid s, a) - P_h^*(s' \mid s, a)] V_{\text{LLM},\lambda,h+1}^\star(s') + u_h^t(s, a)$$

when the induction hypothesis holds. Since  $|V_{\text{LLM},\lambda,h+1}^{\star}(s')| \leq H$ , we have

$$\overline{Q}_h^t(s,a) - Q_{\text{LLM},\lambda,h}^\star(s,a) \ge u_h^t(s,a) - H \sum_{s' \in \mathcal{S}} \left| P_h^t(s' \mid s,a) - P_h^\star(s' \mid s,a) - u_h^\star(s,a) - H \operatorname{TV}(P_h^t(\cdot \mid s,a), P_h^\star(\cdot \mid s,a)) \right|$$

$$= u_h^t(s, a) - H \operatorname{TV}(P_h^t(\cdot \mid s, a), P_h^*(\cdot \mid$$

When we condition on  $\mathcal{E}_3$ , we have

$$H \operatorname{TV}(P_{h}^{t}(\cdot \mid s, a), P_{h}^{*}(\cdot \mid s, a)) \leq \max\left\{2H, \sqrt{\frac{\log(4HTS^{2}A/\delta)}{n_{h}^{t}(s, a)}}\right\} = u_{h}^{t}(s, a),$$
(B.4)

which implies  $\overline{Q}_{h}^{t}(s, a) \geq Q_{\text{LLM},\lambda,h}^{\star}(s, a)$ . Therefore, we have

$$\begin{split} \overline{V}_{h}^{t}(s) &= \max_{\pi(\cdot \mid s)} \sum_{a \in \mathcal{A}} \pi(a \mid s) \overline{Q}_{h}^{t}(s, a) - \lambda \mathrm{KL} \big( \pi(\cdot \mid s), \pi^{\mathrm{LLM}}(\cdot \mid s) \big) \\ &\geq \max_{\pi(\cdot \mid s)} \sum_{a \in \mathcal{A}} \pi(a \mid s) Q_{\mathrm{LLM},\lambda,h}^{\star}(s, a) - \lambda \mathrm{KL} \big( \pi(\cdot \mid s), \pi^{\mathrm{LLM}}(\cdot \mid s) \big) = V_{\mathrm{LLM},\lambda,h}^{\star}(s), \end{split}$$

which concludes the first part of the proof.

By the definition of  $\underline{Q}_{H+1}^t$  and  $Q_{\text{LLM},\lambda,H+1}^{\star}(s,a)$ , Lemma B.6 holds when h = H + 1. By the boundedness of  $Q_{\text{LLM},\lambda,h}^{\star}$ , Lemma B.6 obviously holds when  $Q_{h}^{t} \leq 0$ . Otherwise, we have

$$\underline{Q}_{h}^{t}(s,a) - Q_{\mathrm{LLM},\lambda,h}^{\star}(s,a) \leq \sum_{s' \in \mathcal{S}} P_{h}^{t}(s' \mid s,a) \underline{V}_{h+1}^{t}(s') - \sum_{s' \in \mathcal{S}} P_{h}^{\star}(s' \mid s,a) V_{\mathrm{LLM},\lambda,h+1}^{\star}(s') - u_{h}^{t}(s,a)$$
$$\leq \sum_{s' \in \mathcal{S}} \left[ P_{h}^{t}(s' \mid s,a) - P_{h}^{\star}(s' \mid s,a) \right] V_{\mathrm{LLM},\lambda,h+1}^{\star}(s') - u_{h}^{t}(s,a)$$

when the induction hypothesis holds. Since  $|V_{\text{LLM},\lambda,h+1}^{\star}(s')| \leq H$ , we have

$$\begin{split} \underline{Q}_h^t(s,a) - Q_{\mathrm{LLM},\lambda,h}^\star(s,a) &\geq H \sum_{s' \in \mathcal{S}} \left| P_h^t(s' \mid s,a) - P_h^*(s' \mid s,a) \right| - u_h^t(s,a) \\ &= H \operatorname{TV}(P_h^t(\cdot \mid s,a), P_h^*(\cdot \mid s,a)) - u_h^t(s,a). \end{split}$$

Therefore, we have  $\underline{Q}_{h}^{t}(s, a) \geq Q_{\text{LLM},\lambda,h}^{\star}(s, a)$  when we condition on  $\mathcal{E}_{3}$  by (B.4). We have

$$\underline{V}_{h}^{t}(s) = \max_{\pi(\cdot|s)} \sum_{a \in \mathcal{A}} \pi(a \mid s) \underline{Q}_{h}^{t}(s, a) - \lambda \mathrm{KL} \big( \pi(\cdot \mid s), \pi^{\mathrm{LLM}}(\cdot \mid s) \big)$$
$$\geq \max_{\pi(\cdot|s)} \sum_{a \in \mathcal{A}} \pi(a \mid s) Q_{\mathrm{LLM},\lambda,h}^{\star}(s, a) - \lambda \mathrm{KL} \big( \pi(\cdot \mid s), \pi^{\mathrm{LLM}}(\cdot \mid s) \big) = V_{\mathrm{LLM},\lambda,h}^{\star}(s)$$

which conclude the proof of Lemma B.6.

### 918 B.4 PROOF OF LEMMA B.1

*Proof.* We have the following lemma.

Lemma B.5. For two vectors 
$$x = (x_1, \dots, x_A) \in \mathbb{R}^A$$
 and  $\bar{x} = (\bar{x}_1, \dots, \bar{x}_A) \in \mathbb{R}^A$ , we have  

$$\max_{\substack{\sum_{i=1}^A \beta_i = 1, \beta_i > 0}} \left[ \sum_{i=1}^A x_i \beta_i - \lambda \sum_{i=1}^A \beta_i \log \frac{\beta_i}{\tilde{\beta}_i} \right] - \max_{\substack{\sum_{i=1}^A \beta_i = 1, \beta_i > 0}} \left[ \sum_{i=1}^A \bar{x}_i \beta_i - \lambda \sum_{i=1}^A \beta_i \log \frac{\beta_i}{\tilde{\beta}_i} \right]$$

$$\leq \sum_{i=1}^A \bar{\beta}_i (x_i - \bar{x}_i) + \frac{A}{2\lambda} \max_i |x_i - \bar{x}_i|^2,$$

where  $\{\bar{\beta}_i\} = \operatorname{argmax}_{\sum_{i=1}^A \beta_i = 1, \beta_i > 0} \left[ \sum_{i=1}^A \bar{x}_i \beta_i - \sum_{i=1}^A \beta_i \log \frac{\beta_i}{\bar{\beta}_i} \right].$ 

*Proof.* See §B.8 for a detailed proof.

We prove Lemma B.1 using induction on h. The statement obviously holds when 
$$h = H + 1$$

$$V_{\pi^{\text{LLM}},\lambda,h}^{\star}(s) - V_{\pi^{\text{LLM}},\lambda,h}^{\overline{\pi}^{t+1}}(s) = V_{\pi^{\text{LLM}},\lambda,h}^{\star}(s) - \overline{V}_{h}^{t}(s) + \overline{V}_{h}^{t}(s) - V_{\pi^{\text{LLM}},\lambda,h}^{\overline{\pi}^{t+1}}(s).$$
(B.5)

By the definition of  $V_{\pi^{\text{LLM}},\lambda,h}^{\star}$ ,  $V_{\pi^{\text{LLM}},\lambda,h}^{\overline{\pi}^{t+1}}$  and Lemma B.5, we have

$$V_{\pi^{\mathrm{LLM}},\lambda,h}^{\star}(s) - \overline{V}_{h}^{t}(s)$$

$$\leq \sum_{a \in \mathcal{A}} \bar{\pi}_{h}^{t+1}(a \mid s) \left[ Q_{\pi^{\mathrm{LLM}},\lambda,h}^{\star}(s,a) - \overline{Q}_{h}^{t}(s,a) \right] + \frac{A}{2\lambda} \max_{a} |Q_{\pi^{\mathrm{LLM}},\lambda,h}^{\star}(s,a) - \overline{Q}_{h}^{t}(s,a)|^{2}.$$
(B.6)

We also have

$$\overline{V}_{h}^{t}(s) - V_{\pi^{\mathrm{LLM}},\lambda,h}^{\overline{\pi}^{t+1}}(s) = \sum_{a \in \mathcal{A}} \overline{\pi}_{h}^{t+1}(a \mid s) \left[ \overline{Q}_{h}^{t}(s,a) - Q_{\pi^{\mathrm{LLM}},\lambda,h}^{\overline{\pi}^{t+1}}(s,a) \right].$$
(B.7)

Combining (B.5), (B.6) and (B.7), we have

$$\begin{split} V_{\pi^{\text{LLM}},\lambda,h}^{\star}(s) &- V_{\pi^{\text{LLM}},\lambda,h}^{\overline{\pi}^{t+1}}(s) \\ &\leq \sum_{a \in \mathcal{A}} \bar{\pi}_{h}^{t+1}(a \mid s) \Big[ Q_{\pi^{\text{LLM}},\lambda,h}^{\star}(s,a) - Q_{\pi^{\text{LLM}},\lambda,h}^{\overline{\pi}^{t+1}}(s,a) \Big] + \frac{A}{2\lambda} \max_{a} |Q_{\pi^{\text{LLM}},\lambda,h}^{\star}(s,a) - \overline{Q}_{h}^{t}(s,a)|^2 \\ &= \mathbb{E}_{\overline{\pi}^{t+1}} \Big[ V_{\pi^{\text{LLM}},\lambda,h+1}^{\star}(s_{h+1},a_{h+1}) - V_{\pi^{\text{LLM}},\lambda,h+1}^{\overline{\pi}^{t+1}}(s_{h+1},a_{h+1}) \mid s_{h} = s \Big] + \frac{A}{2\lambda} \max_{a} |Q_{\pi^{\text{LLM}},\lambda,h}^{\star}(s,a) - \overline{Q}_{h}^{t}(s,a)|^2 \end{split}$$

By induction, we easily have

$$V_{\pi^{\text{LLM}},\lambda,h}^{\star}(s) - V_{\pi^{\text{LLM}},\lambda,h}^{\overline{\pi}^{t+1}}(s) \le \frac{A}{2\lambda} \sum_{h'=h}^{H} \mathbb{E}_{\overline{\pi}^{t+1}} \Big[ \max_{a} |Q_{\pi^{\text{LLM}},\lambda,h'}^{\star}(s_{h'},a_{h'}) - \overline{Q}_{h'}^{t}(s_{h'},a_{h'})|^2 |s_h = s \Big]$$
(B.8)

We also have the following lemma, which shows that  $Q_{\pi^{\text{LLM}},\lambda,h'}^{\star}$  lies between  $\overline{Q}_{h'}^{t}$  and  $\underline{Q}_{h'}^{t}$ . Lemma B.6. We have

$$\underline{Q}_{h}^{t}(s,a) \leq Q_{\pi^{\text{LLM}},\lambda,h}^{\star}(s,a) \leq \overline{Q}_{h}^{t}(s,a), \qquad \underline{V}_{\lambda,h}^{t}(s) \leq V_{\pi^{\text{LLM}},\lambda,h}^{\star}(s) \leq \overline{V}_{h}^{t}(s)$$
  
holds for all  $t \in \mathbb{N}, (h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$  with probability  $1 - \delta/2$ .

*Proof.* See §B.3 for a detailed proof.

Therefore, by  $(\mathbf{B.8})$ , we have

$$V_{\pi^{\text{LLM}},\lambda,h}^{\star}(s) - V_{\pi^{\text{LLM}},\lambda,h}^{\overline{\pi}^{t+1}}(s) \le \frac{A}{2\lambda} \sum_{h'=h}^{H} \mathbb{E}_{\overline{\pi}^{t+1}} \Big[ \max_{a} |\overline{Q}_{h'}^{t}(s_{h'},a_{h'}) - \underline{Q}_{h'}^{t}(s_{h'},a_{h'})|^{2} | s_{h} = s \Big]$$
(B.9)

when the event in Lemma **B**.6 holds, which conclude the proof of Lemma **B**.1.

972 B.5 PROOF OF LEMMA B.2 

*Proof.* By the definition of  $\pi^t$  in (3.4), we have 

$$\mathbb{E}_{\overline{\pi}^{t+1}} \left[ \max_{a \in \mathcal{A}} \left( \overline{Q}_h^t(s_h, a) - \underline{Q}_h^t(s_h, a) \right)^2 \right] \leq H \left( \frac{H}{H-1} \right)^{h-1} \mathbb{E}_{\pi^{t+1}} \left[ \left( \overline{Q}_h^t(s_h, a_h) - \underline{Q}_h^t(s_h, a_h) \right)^2 \right]$$
(B.10)  
$$\leq e H \mathbb{E}_{\pi^{t+1}} \left[ \left( \overline{Q}_h^t(s_h, a_h) - \underline{Q}_h^t(s_h, a_h) \right)^2 \right].$$

Next we analyze the expression under the square. First, we have

$$\overline{Q}_{h}^{t}(s_{h},a_{h}) - \underline{Q}_{h}^{t}(s_{h},a_{h}) \leq 2u_{h}^{t}(s_{h},a_{h}) + \sum_{a \in \mathcal{A}} P_{h}^{t}(s'|s_{h},a_{h}) [\overline{V}_{h+1}^{t}(s') - \underline{V}_{h+1}^{t}(s')] \\
\leq 2u_{h}^{t}(s_{h},a_{h}) + \sum_{a \in \mathcal{A}} P_{h}^{*}(s'|s_{h},a_{h}) [\overline{V}_{h+1}^{t}(s') - \underline{V}_{h+1}^{t}(s')] + H \operatorname{TV}(P_{h}^{t}(\cdot \mid s,a), P_{h}^{*}(\cdot \mid s,a)).$$

Therefore, by  $(\mathbf{B.4})$ , we have

$$\overline{Q}_h^t(s_h, a_h) - \underline{Q}_h^t(s_h, a_h) \le 3u_h^t(s_h, a_h) + \sum_{a \in \mathcal{A}} P_h^*(s'|s_h, a_h) [\overline{V}_{h+1}^t(s') - \underline{V}_{h+1}^t(s')].$$

Therefore, we have

$$\begin{aligned} \overline{Q}_{h}^{t}(s_{h}, a_{h}) - \underline{Q}_{h}^{t}(s_{h}, a_{h}) &\leq 3AH \cdot \mathbb{E}_{\overline{\pi}^{t+1}} \left[ \sum_{h'=h}^{H} \sqrt{\frac{\log(4HTS^{2}A/\delta)}{n_{h}^{t}(s, a)}} \middle| s_{h} \right] \\ &\leq 9AH \cdot \mathbb{E}_{\pi^{t+1}} \left[ \sum_{h'=h}^{H} \sqrt{\frac{\log(4HTS^{2}A/\delta)}{n_{h}^{t}(s, a)}} \middle| s_{h} \right] \end{aligned}$$

by induction and the definition of  $u_h^t$  in (3.1). By Cauchy inequality, we have

$$\overline{Q}_{h}^{t}(s_{h},a_{h}) - \underline{Q}_{h}^{t}(s_{h},a_{h}) \leq 9AH^{3/2} \sqrt{\mathbb{E}_{\pi^{t+1}} \left[ \sum_{h'=h}^{H} \frac{\log(4HTS^{2}A/\delta)}{n_{h}^{t}(s,a)} \middle| s_{h} \right]}.$$
(B.11)

Combining (B.10) and (B.11), we have

$$\mathbb{E}_{\overline{\pi}^{t+1}} \left[ \max_{a \in \mathcal{A}} \left( \overline{Q}_h^t(s_h, a) - \underline{Q}_h^t(s_h, a) \right)^2 \right] \le 230 A^2 H^4 \mathbb{E}_{\pi^{t+1}} \left[ \mathbb{E}_{\pi^{t+1}} \left[ \mathbb{E}_{\pi^{t+1}} \left[ \sum_{h'=h}^H \frac{\log(4HTS^2 A/\delta)}{n_h^t(s, a)} \middle| s_h \right] \right] \right]$$
(B.12)  
$$= 230 A^2 H^4 \mathbb{E}_{\pi^{t+1}} \left[ \sum_{h'=h}^H \frac{\log(4HTS^2 A/\delta)}{n_h^t(s, a)} \middle| s_h \right].$$

1012 We also have the following lemma.

Lemma B.7. We define 1014

$$\mathcal{E}^{\text{cnt}} \triangleq \{n_h^t(s, a) \ge \frac{1}{2}\bar{n}_h^t(s, a) - \log(12SAH/\delta) \text{ for all } (s, a, h, t) \in \mathcal{S} \times \mathcal{A} \times [H] \times [T]\},$$
  
Then we have  $P(\mathcal{E}^{\text{cnt}}) \ge 1 - \delta/4.$ 

1018<br/>1019Proof. This is Lemma 3 of Ménard et al. (2021). See Ménard et al. (2021) for a detailed proof.□

When we condition on  $\mathcal{E}^{\text{cnt}}$ , we have

$$\begin{array}{ll} \begin{array}{l} \mbox{1022} \\ \mbox{1023} \\ \mbox{1023} \\ \mbox{1024} \\ \mbox{1024} \\ \mbox{1025} \end{array} & \begin{array}{l} \frac{\log(4HTS^2A/\delta)}{n_h^t(s,a)} \leq \frac{\log(4HTS^2A/\delta) + \log(12SAH/\delta)}{n_h^t(s,a) + \log(12SAH/\delta)} \leq \frac{2\log(12HTS^2A/\delta)}{n_h^t(s,a) + \log(12SAH/\delta)} \\ \leq \frac{4\log(12HTS^2A/\delta)}{\bar{n}_h^t(s,a)}. \end{array} \\ \end{array}$$

1026 Therefore, by (B.12), we have

$$\mathbb{E}_{\overline{\pi}^{t+1}}\left[\max_{a\in\mathcal{A}}\left(\overline{Q}_{h}^{t}(s_{h},a)-\underline{Q}_{h}^{t}(s_{h},a)\right)^{2}\right] \leq 920A^{2}H^{4}\mathbb{E}_{\pi^{t+1}}\left[\sum_{h'=h}^{H}\frac{\log(12HTS^{2}A/\delta)}{\overline{n}_{h}^{t}(s,a)}\Big|s_{h}\right].$$

<sup>1031</sup> Therefore, we have

$$\sum_{t=1}^{1032} \sum_{h=1}^{T} \sum_{h=1}^{H} \mathbb{E}_{\overline{\pi}^{t+1}} \left[ \max_{a \in \mathcal{A}} \left( \overline{Q}_{h}^{t}(s_{h}, a) - \underline{Q}_{h}^{t}(s_{h}, a) \right)^{2} \right] \leq \sum_{t=1}^{T} \sum_{h=1}^{H} 920A^{2}H^{4} \mathbb{E}_{\pi^{t+1}} \left[ \sum_{h'=h}^{H} \frac{\log(12HTS^{2}A/\delta)}{\overline{n}_{h'}^{t}(s, a)} \middle| s_{h} \right]$$

$$\leq \sum_{t=1}^{T} \sum_{h=1}^{H} 920A^{2}H^{5} \mathbb{E}_{\pi^{t+1}} \left[ \frac{\log(12HTS^{2}A/\delta)}{\overline{n}_{h}^{t}(s, a)} \middle| s_{h} \right]$$

$$\leq 920SA^{3}H^{6} \log(12HTS^{2}A/\delta) \log T,$$

1040 which concludes the proof of Lemma B.2.

### 1042 B.6 PROOF OF LEMMA B.3 1043

1044 *Proof.* We denote by  $n_h^t(s, a)$  the number of times the state action-pair (s, a) was visited in step 1045 *h* in the first *t* episodes,  $n_h^t(s, a, s')$  the number of times the state action next state -pair (s, a, s')1046 was visited in step *h* in the first *t* episodes, and define  $N_h^t(s, a, s') = n_h^t(s, a)P_h(s' | s, a)$ . By 1047 Hoeffding's inequality, we have

1056 1057 1058

1061

1062

1068 1069

1076 1077 1078

1041

1028 1029 1030

 $\mathbb{1}_{n_h^t(s,a)>0} \left| \frac{n_h^t(s,a,s') - N_h^t(s,a,s')}{\sqrt{n_h^t(s,a)}} \right| \le \sqrt{\log(4HTS^2 A/\delta)}$ (B.13)

holds for a fix  $(t, h, s, a, s') \in [T] \times [H] \times S \times A \times S$  with probability at least  $1 - \delta/(2HTS^2A)$ . By taking a union bound over all  $(t, h, s, a, s') \in [T] \times [H] \times S \times A \times S$ , we have (B.13) holds for all  $(t, h, s, a, s') \in [T] \times [H] \times S \times A \times S$  with probability at least  $1 - \delta/2$ . We denote by  $\mathcal{E}_2$  that (B.13) holds for all  $(t, h, s, a, s') \in [T] \times [H] \times S \times A \times S$ . When condition on  $\mathcal{E}_2$ , we have

$$\begin{split} \mathbb{1}_{n_{h}^{t}(s,a)>0} \operatorname{TV}(P_{h}^{*}(\cdot \mid s,a), P_{h}^{t}(\cdot \mid s,a)) &= \sum_{s' \in \mathcal{S}} \mathbb{1}_{n_{h}^{t}(s,a)>0} \left| \frac{n_{h}^{t}(s,a,s') - N_{h}^{t}(s,a,s')}{n_{h}^{t}(s,a)} \right| \\ &\leq A \sqrt{\frac{\log(4HTS^{2}A/\delta)}{n_{h}^{t}(s,a)}}. \end{split}$$

We conclude the proof by noticing that  $P(\mathcal{E}_2) \ge 1 - \delta/2$ .

1064 B.7 PROOF OF LEMMA B.4

Proof. We prove this using induction. Lemma B.4 obviously holds when h = H + 1. If Lemma B.4 hold for h + 1, we have

$$Q^{\star}_{\mathrm{LLM},\lambda,h}(s,a) = r_h(s,a) + \sum_{s' \in \mathcal{S}} P^{\star}_h(s'|s,a) V^{\star}_{\mathrm{LLM},\lambda,h+1}(s,a) \ge 0.$$

1070 1071 We also have

1072 
$$Q_{\text{LLM},\lambda,h}^{\star}(s,a) = r_h(s,a) + \sum_{s' \in \mathcal{S}} P_h^{\star}(s'|s,a) V_{\text{LLM},\lambda,h+1}^{\star}(s,a) \le 1 + \sum_{s' \in \mathcal{S}} P_h^{\star}(s'|s,a) (H-h) = H-h+1$$
1073 1074 Using the set of the set o

Using the property of constraint optimization, we have

$$V^{\star}_{\mathrm{LLM},\lambda,h}(s) = \lambda \log \bigl( \sum_{a \in \mathcal{A}} \pi^{\mathrm{LLM}}_{h}(a|s) \exp(Q^{\star}_{\mathrm{LLM},\lambda,h}(s,a)/\lambda) \bigr).$$

since  $0 \le Q_{\text{LLM},\lambda,h}^{\star}(s,a) \le H + 1 - h$ , we have  $0 \le V_{\text{LLM},\lambda,h}^{\star}(s) \le H + 1 - h$ . Therefore, we conclude the proof of Lemma B.4 using induction.

## B.8 PROOF OF LEMMA B.5

*Proof.* For  $\boldsymbol{x} = (x_1, \dots, x_A)^\top$  and  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_A)^\top$ , we define 

1084  
1085  
1086
$$f(\boldsymbol{x},\boldsymbol{\beta}) = \boldsymbol{x}^{\top}\boldsymbol{\beta} - \lambda \sum_{i=1}^{A} \beta_i \log \frac{\beta_i}{\widetilde{\beta_i}}, \quad f^*(\boldsymbol{x}) = \max_{\sum_{i=1}^{A} \beta_i = 1, \beta_i > 0} f(\boldsymbol{x},\boldsymbol{\beta}), \quad \boldsymbol{\beta}^*(\boldsymbol{x}) = \arg_{\sum_{i=1}^{A} \beta_i = 1, \beta_i > 0} f(\boldsymbol{x},\boldsymbol{\beta})$$
(B.14)

We have the following lemma.

Lemma B.8. We have 

$$f^*(\boldsymbol{x}) = \lambda \log\left[\sum_{i=1}^A \widetilde{\beta}_i \exp(x_i/\lambda)\right], \nabla f^*(\boldsymbol{x}) = \boldsymbol{\beta}^*(\boldsymbol{x}), \beta_i(\boldsymbol{x}) = \frac{\partial}{\partial x_i} f^*(\boldsymbol{x}) = \frac{\widetilde{\beta}_i \exp(x_i/\lambda)}{\sum_{j=1}^A \widetilde{\beta}_j \exp(x_j/\lambda)}$$

Moreover, we have  $\|\nabla f^*(\boldsymbol{x}_1) - \nabla f^*(\boldsymbol{x}_2)\|_1 \leq \frac{A}{\lambda} \|\boldsymbol{x}_1 - \boldsymbol{x}_2\|_{\infty}$ . 

*Proof.* See §B.9 for a detailed proof. 

First, we have 

$$f^*(\boldsymbol{x}) - f^*(\bar{\boldsymbol{x}}) = \nabla f^*(\bar{\boldsymbol{x}})(\boldsymbol{x} - \bar{\boldsymbol{x}}) + + \int_0^1 \left( \nabla f^*(t\boldsymbol{x} + (1-t)\bar{\boldsymbol{x}}) - \nabla f^*(\bar{\boldsymbol{x}}) \right)^\top (\boldsymbol{x} - \bar{\boldsymbol{x}}) dt$$
$$\leq \nabla f^*(\bar{\boldsymbol{x}})(\boldsymbol{x} - \bar{\boldsymbol{x}}) + + \int_0^1 \left\| \nabla f^*(t\boldsymbol{x} + (1-t)\bar{\boldsymbol{x}}) - \nabla f^*(\bar{\boldsymbol{x}}) \right\|_1 \|\boldsymbol{x} - \bar{\boldsymbol{x}}\|_\infty dt.$$

By Lemma B.8, we have

$$f^*(\boldsymbol{x}) - f^*(\bar{\boldsymbol{x}}) \le \nabla f^*(\bar{\boldsymbol{x}})(\boldsymbol{x} - \bar{\boldsymbol{x}}) + \frac{A}{\lambda} \int_0^1 \|t\boldsymbol{x} + (1 - t)\bar{\boldsymbol{x}} - \bar{\boldsymbol{x}}\|_{\infty} \|\boldsymbol{x} - \bar{\boldsymbol{x}}\|_{\infty} dt$$
$$= \nabla f^*(\bar{\boldsymbol{x}})(\boldsymbol{x} - \bar{\boldsymbol{x}}) + \frac{A}{2\lambda} \|\boldsymbol{x} - \bar{\boldsymbol{x}}\|_{\infty}^2$$
(B.15)

We conclude the proof of Lemma B.5 by combining (B.15) with Lemma B.8. 

#### B.9 PROOF OF LEMMA B.8

*Proof.* By the property of constraint optimization, we have  $\frac{\partial}{\partial \beta} f(\boldsymbol{x}, \boldsymbol{\beta}(\boldsymbol{x})) = c(1, \dots, 1)^{\top}$  for some constant c. By direct computation, we have 

$$\frac{\partial}{\partial \beta} f(\boldsymbol{x}, \boldsymbol{\beta}(\boldsymbol{x})) = x - \lambda \log \boldsymbol{\beta} + \lambda \log \widetilde{\boldsymbol{\beta}} - \lambda (1, \dots, 1)^{\top},$$

which implies  $\beta_i(\boldsymbol{x}) \propto \widetilde{\beta}_i \exp(x_i/\lambda)$ . Since  $\sum_{i=1}^A \beta_i(\boldsymbol{x} \propto) = 1$ , we have 

$$\beta_i(\boldsymbol{x}) = \frac{\widetilde{\beta}_i \exp(x_i/\lambda)}{\sum_{j=1}^A \widetilde{\beta}_j \exp(x_j/\lambda)}.$$
(B.16)

By the definition of  $f^*$ , we have 

$$f^{*}(\boldsymbol{x}) = f(\boldsymbol{x}, \boldsymbol{\beta}_{i}(\boldsymbol{x})) = \frac{\sum_{i=1}^{A} x_{i} \widetilde{\beta}_{i} \exp(x_{i}/\lambda)}{\sum_{i=1}^{A} \widetilde{\beta}_{i} \exp(x_{i}/\lambda)} - \frac{\sum_{i=1}^{A} x_{i} \widetilde{\beta}_{i} \exp(x_{i}/\lambda)}{\sum_{i=1}^{A} \widetilde{\beta}_{i} \exp(x_{i}/\lambda)} + \lambda \log(\sum_{i=1}^{A} \widetilde{\beta}_{i} \exp(x_{i}/\lambda))$$

$$= \frac{\sum_{i=1}^{A} x_{i} \widetilde{\beta}_{i} \exp(x_{i}/\lambda)}{A} - \frac{\sum_{i=1}^{A} x_{i} \widetilde{\beta}_{i} \exp(x_{i}/\lambda)}{\sum_{i=1}^{A} \widetilde{\beta}_{i} \exp(x_{i}/\lambda)} + \lambda \log(\sum_{i=1}^{A} \widetilde{\beta}_{i} \exp(x_{i}/\lambda))$$

1130  
1131 
$$= \lambda \log\left(\sum_{i=1}^{A} \widetilde{\beta}_{i} \exp(x_{i}/\lambda)\right).$$
(B.17)

Combining (B.16) and (B.17), we have  $\nabla f^*(x) = \beta(x)$ . We define  $\Delta = \{\beta \mid \sum_{i=1}^A \beta_i = 1, \beta_i > 0\}$ 0}, we have the following lemma.

**Lemma B.9.** For any  $x, \beta \in \Delta$ , we have

$$\lambda \sum_{i=1}^{A} \beta_i \log \frac{\beta_i}{\widetilde{\beta}_i} \geq \lambda \sum_{i=1}^{A} \beta_i(\boldsymbol{x}) \log \frac{\beta_i(\boldsymbol{x})}{\widetilde{\beta}_i} + \boldsymbol{x}^\top \big(\boldsymbol{\beta} - \boldsymbol{\beta}(\boldsymbol{x})\big) + \frac{\lambda}{2A} \|\boldsymbol{\beta} - \boldsymbol{\beta}(\boldsymbol{x})\|_1^2$$

1139 where  $\beta(x)$  is defined in (B.14).

*Proof.* See §B.10 for a detailed proof.

By Lemma **B**.9, we have

$$\lambda \sum_{i=1}^{A} \beta_i(\boldsymbol{x}_2) \log \frac{\beta_i(\boldsymbol{x}_2)}{\widetilde{\beta}_i} \ge \lambda \sum_{i=1}^{A} \beta_i(\boldsymbol{x}_1) \log \frac{\beta_i(\boldsymbol{x}_1)}{\widetilde{\beta}_i} + \boldsymbol{x}_1^{\top} \left( \boldsymbol{\beta}(\boldsymbol{x}_2) - \boldsymbol{\beta}(\boldsymbol{x}_1) \right) + \frac{\lambda}{2A} \| \boldsymbol{\beta}(\boldsymbol{x}_2) - \boldsymbol{\beta}(\boldsymbol{x}_1) \|_1^2$$
$$\lambda \sum_{i=1}^{A} \beta_i(\boldsymbol{x}_1) \log \frac{\beta_i(\boldsymbol{x}_1)}{\widetilde{\beta}_i} \ge \lambda \sum_{i=1}^{A} \beta_i(\boldsymbol{x}_2) \log \frac{\beta_i(\boldsymbol{x}_2)}{\widetilde{\beta}_i} + \boldsymbol{x}_2^{\top} \left( \boldsymbol{\beta}(\boldsymbol{x}_1) - \boldsymbol{\beta}(\boldsymbol{x}_2) \right) + \frac{\lambda}{2A} \| \boldsymbol{\beta}(\boldsymbol{x}_2) - \boldsymbol{\beta}(\boldsymbol{x}_1) \|_1^2$$

which implies

$$rac{\lambda}{A} \| oldsymbol{eta}(oldsymbol{x}_2) - oldsymbol{eta}(oldsymbol{x}_1) \|_1^2 \leq (oldsymbol{x}_1 - oldsymbol{x}_2)^ op ig(oldsymbol{x}_1) - oldsymbol{eta}(oldsymbol{x}_2) ig) \leq \| oldsymbol{eta}(oldsymbol{x}_2) - oldsymbol{eta}(oldsymbol{x}_1) \|_1 \|oldsymbol{x}_1 - oldsymbol{x}_2 \|_\infty.$$

1154 Therefore, we have

$$\frac{\lambda}{A} \| \boldsymbol{\beta}(\boldsymbol{x}_2) - \boldsymbol{\beta}(\boldsymbol{x}_1) \|_1 \le \| \boldsymbol{x}_1 - \boldsymbol{x}_2 \|_{\infty}$$

which concludes the proof of Lemma B.8

### 1160 В.10 Ркооf of Lemma В.9

*Proof.* Let  $\beta_i$  be the *i*-th element of  $\beta$ , and  $\beta_i(\boldsymbol{x})$  be the *i*-th element of  $\beta$ . We define  $g(\beta) = \lambda \sum_{i=1}^{A} \beta_i \log \frac{\beta_i}{\beta_i}$ . By the property of the function g, we have

$$g(\boldsymbol{\beta}) - g(\boldsymbol{\beta}(\boldsymbol{x})) = \nabla g(\boldsymbol{\beta}(\boldsymbol{x}))^{\top} (\boldsymbol{\beta} - \boldsymbol{\beta}(\boldsymbol{x})) + \sum_{i=1}^{A} \int_{\beta_{i}(\boldsymbol{x})}^{\beta_{i}} \int_{\beta_{i}(\boldsymbol{x})}^{u} \frac{\lambda}{v} \mathrm{d}v \mathrm{d}u$$

1167 Since  $\beta_i < 1$ , we have

$$g(\boldsymbol{\beta}) - g(\boldsymbol{\beta}(\boldsymbol{x})) \ge \nabla g(\boldsymbol{\beta}(\boldsymbol{x}))^{\top} (\boldsymbol{\beta} - \boldsymbol{\beta}(\boldsymbol{x})) + \lambda \sum_{i=1}^{A} \int_{\beta_{i}(\boldsymbol{x})}^{\beta_{i}} \int_{\beta_{i}(\boldsymbol{x})}^{u} 1 \mathrm{d}v \mathrm{d}u$$
(B.18)

$$egin{split} &= 
abla gig(oldsymbol{eta}(oldsymbol{x})ig)^{ op}ig(oldsymbol{eta}-oldsymbol{eta}(oldsymbol{x})ig) + rac{\lambda}{2}\sum_{i=1}^A |eta_i(oldsymbol{x})-eta_i|^2 \end{split}$$

$$\geq \nabla g \big( \boldsymbol{\beta}(\boldsymbol{x}) \big)^\top \big( \boldsymbol{\beta} - \boldsymbol{\beta}(\boldsymbol{x}) \big) + \frac{\lambda}{2A} \| \boldsymbol{\beta}(\boldsymbol{x}) - \boldsymbol{\beta} \|_1^2.$$

1177 By the definition of the function g, we have

$$\nabla g(\boldsymbol{\beta}(\boldsymbol{x})) = (\lambda \log \beta_1 + \lambda - \lambda \log \widetilde{\beta}_1, \dots, \lambda \log \beta_A + \lambda - \lambda \log \widetilde{\beta}_A)^\top.$$
(B.19)

1180 By the definition of  $\beta(x)$  in (B.14), we can easily obtain

$$\beta_i(\boldsymbol{x}) = \frac{\widetilde{\beta}_i \exp(x_i/\lambda)}{\sum_{i=1}^A \widetilde{\beta}_i \exp(x_i/\lambda)}.$$
(B.20)

1184 Combining (B.19) with (B.20), we have

1186  
1187 
$$\nabla g(\boldsymbol{\beta}(\boldsymbol{x})) = \boldsymbol{x} + \lambda \Big(1 - \log(\sum_{i=1}^{A} \widetilde{\beta}_{i} \exp(x_{i}/\lambda))\Big)(1, 1, \dots, 1)^{\top}.$$

1188 Therefore, we have 1189

1190 1191

1192 1193

1194

1195

1196 1197

1200

1201

1202

$$(\nabla g(\boldsymbol{\beta}(\boldsymbol{x})) - \boldsymbol{x})^{\top} (\boldsymbol{\beta} - \boldsymbol{\beta}(\boldsymbol{x}))$$
  
=  $\lambda \Big( 1 - \log \Big( \sum_{i=1}^{A} \widetilde{\beta}_{i} \exp(x_{i}/\lambda) \Big) \Big) (1, 1, \dots, 1)^{\top} \boldsymbol{\beta} - \lambda \Big( 1 - \log \Big( \sum_{i=1}^{A} \widetilde{\beta}_{i} \exp(x_{i}/\lambda) \Big) \Big) (1, 1, \dots, 1)^{\top} \boldsymbol{\beta}(\boldsymbol{x})$ 

By the definition of the region  $\Delta$  in Lemma B.8, we have  $(1,1,\ldots,1)^{\top}\beta = (1,1,\ldots,1)^{\top}\beta(x) =$ 1. Therefore, we have

$$\nabla g(\boldsymbol{\beta}(\boldsymbol{x})) - \boldsymbol{x})^{\top} (\boldsymbol{\beta} - \boldsymbol{\beta}(\boldsymbol{x})) = 0.$$
(B.21)

1198 We conclude the proof of Lemma B.9 by combining (B.18) and (B.21). 1199

### **EXPERIMENT DETAILS** С

1203 We first describe how the rule-based value estimator is implemented. For the "put" task, we prompt the LLM with the following prompt: 1204

```
1205
                                          Input
1206
        Given the following successful trajectory, output sub-goals that are
1207
        needed to complete the task ``put A on B" and write code to determine
1208
       how many sub-goals are satisfied at the current state.
1209
1210
        Looking quickly around you, you see a shelf 6, a safe 1, and a drawer
1211
        4. Your task is to: put some vase in safe.
        > go to shelf 6
1212
        On the shelf 6, you see a vase 2.
1213
        > take vase 2 from shelf 6
1214
        You pick up the vase 2 from the shelf 6.
1215
        > go to safe 1
1216
        The safe 1 is closed.
1217
        > open safe 1
        Safe 1 is open. In it, you see a keychain 3.
1218
        > put vase 2 in/on safe 1
1219
        Vase 2 is in/on safe 1.
1220
1221
```

1222 The LLM responds as follows:

```
Output
1224
        Sub-goals needed to complete the task "put A on B" are
1226
        - Locate object A
        - Find B
        - Ensure B is accessible
        def check_subgoals(state, a, b):
1230
            subgoals_satisfied = 0
            if f"You pick up the {a}" in state:
                subgoals_satisfied += 1
            if f"You arrive at" in state and b in state:
                subgoals_satisfied += 1
            if f"{b} is open" in state:
1235
                subgoals_satisfied += 1
            return subgoals_satisfied
```

1236 1237 1238

1223

1225

1227

1228

1229

1231

1232

1233

1234

1239 Since the LLM is only prompted with one example, it makes a mistake at the last step, which can be 1240 fixed by prompting the LLM with "No need to check if b is open since we only need to put a on b.". 1241 Then the LLM generates a correct rule-based value estimator as follows:

1242	Correct rule-based value estimator
1243	<pre>def check_subgoals(state, a, b):</pre>
1244	<pre>subgoals_satisfied = 0</pre>
1245	if f"You pick up the {a}" in state:
1240	subgoals_satisfied += 1 if f"You arrive at" in state and b in state.
1247	subgoals satisfied += 1
1240	if f"{a} is in/on {b}" in state:
1249	<pre>subgoals_satisfied += 1</pre>
1250	return subgoals_satisfied
1252	
1253	
1254	
1255	
1256	
1257	
1258	
1259	
1260	
1261	
1262	
1263	
1264	
1265	
1266	
1267	
1268	
1209	
1270	
1271	
1273	
1274	
1275	
1276	
1277	
1278	
1279	
1280	
1281	
1282	
1283	
1284	
1285	
1286	
1207	
1200	
1205	
1291	
1292	
1293	
1294	
1295	