# ENHANCING DECISION-MAKING OF LARGE LAN GUAGE MODELS VIA ACTOR-CRITIC

Anonymous authors

Paper under double-blind review

# ABSTRACT

Large Language Models (LLMs) have achieved significant advancements in natural language processing tasks, yet they encounter challenges in complex decisionmaking scenarios that require long-term reasoning and alignment with high-level objectives. This paper introduces a novel gradient-free LLM-based Actor-Critic framework, termed LAC, which addresses these limitations by integrating both action generation and action evaluation mechanisms. Our approach employs two distinct critics: a language-based critic that provides context-sensitive feedback and a value-based critic that offers quantitative assessments of expected long-term rewards. This dual-critic architecture enhances decision-making by leveraging the complementary strengths of both critics, enabling contextually appropriate and more robust action selection. Additionally, we propose a gradient-free policy improvement method that reduces computational overhead, facilitating efficient updates to the actor's policy without the complexities of gradient backpropagation. We validate the effectiveness of LAC across diverse environments that cover both high-level action space (ALFWorld) and low-level action space (BabyAI-Text), demonstrating its superior performance compared to existing state-of-the-art methods. Our method outperforms other state-of-the-art baselines using the same 7B/8B open-source LLMs and even exceeds a strong baseline ReAct using GPT-4 in most settings. Our findings highlight the efficacy and generality of the dual-critic Actor-Critic framework in enhancing LLM-based decision-making.

029 030 031

032

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

# 1 INTRODUCTION

Large Language Models (LLMs) (Touvron et al., 2023; Jiang et al., 2023; Team et al., 2024) have
demonstrated remarkable capabilities across a wide range of tasks in natural language processing,
from text generation to question answering and summarization. Despite their strengths, LLMs often
struggle in more complex decision-making tasks that require not only generating immediate action but
also reasoning over long horizons and aligning actions with high-level objectives (Ahn et al., 2022;
Yao et al., 2022b; Hao et al., 2023; Liu et al., 2023; Huang et al., 2022b). This raises a fundamental
question: how can we efficiently leverage the rich prior knowledge encoded in LLMs to enable more
reliable and effective sequential decision-making in diverse and complex environments?

Recent studies have explored various methods to improve LLM-based decision-making. Through the lens of reinforcement learning (RL) (Barto et al., 1989), these methods typically adopt either an *actor-only* or *critic-only* paradigm. In *actor-only* approaches, the LLM serves as an actor, generating actions based on its autoregressive next-token prediction capabilities (Ahn et al., 2022; Yao et al., 2022b; Huang et al., 2022a; Shinn et al., 2024). While such methods are simple and effective for short-term action generation, they often suffer from a lack of long-term planning. As a result, decisions may appear locally optimal but fail to achieve the overall task objective in more complex, multistep environments.

On the other hand, *critic-only* approaches use LLMs as critics to evaluate candidate actions based
on predicted future trajectories (Hao et al., 2023; Liu et al., 2023; Fu et al., 2024; Brooks et al.,
2024) and select the action with the best outcome. Although this allows for additional evaluation
of actions, it may lead to suboptimal action selection when the sampled candidate actions do not
include the optimal action and the predicted future trajectories deviate from reality. Furthermore,
such methods often prioritize numerical assessments of actions, ignoring important contextual and



Figure 1: An illustrative explanation of our method LAC in ALFWorld. The histogram on the right shows the action probabilities of different methods. While actor ( $\pi_{LLM}$ ) and critics (*lang-critic*  $C_{LLM}$ , *value-critic*  $Q_{LLM}$ ) make mistakes at different time steps, LAC (ours) can select the correct action by integrating actor and critics. The LAC Inference step is detailed in Figure 2.

qualitative language-based information embedded in the task instructions. As shown in Figure 1, these two paradigms fall short of delivering optimal performance, as they fail to balance immediate action generation with long-term action evaluation. While previous works have attempted simple combinations of these roles (Zhang et al., 2023a), they lack a systematic analysis of the interaction between actor and critic, which is crucial for effective decision-making.

092 To address these limitations, we propose a novel gradient-free LLM-based Actor-Critic (LAC) 093 framework. Unlike previous methods, our approach seamlessly integrates both action generation 094 (actor) and action evaluation (critic) to significantly enhance decision-making capabilities. In 095 LAC, the actor improves its policy by using two distinct critics: *lang-critic* and *value-critic*, which 096 provide complementary types of feedback. The *lang-critic* leverages the LLM's natural language 097 understanding capabilities to offer rich, interpretable, and context-sensitive feedback, ensuring that actions align with the task's high-level goals and that previous mistakes are avoided. The *value-critic*, 098 on the other hand, provides numerical evaluations, estimating the long-term reward and ensuring that action distributions are quantitatively optimized. To obtain accurate numerical assessments of actions, 100 we designed a novel *value-critic* estimation method that extracts the internal belief of LLMs on the 101 optimality of candidate actions. 102

For intuition of the dual critics, consider a task where an agent must navigate through a virtual room to pick up an apple:

105

082

083

084

085

The *lang-critic* might guide the agent with insights like "the apple is likely in the refrigerator in the kitchen", helping it prioritize actions that are contextually relevant (*e.g.*, move toward the refrigerator instead of wandering aimlessly).

quantitative judgment. 111 By integrating both critics, LAC selects actions that are both contextually appropriate (based on the 112 lang-critic) and highly likely to succeed (based on the value-critic). As illustrated in Figure 1, our 113 method outperforms previous actor-only and critic-only approaches, which often select suboptimal 114 actions at different time steps. The integration of actor and dual critics in LAC leads to more effective 115 decision-making and consistent task completion. 116

• The *value-critic* evaluates the actions by estimating the long-term reward: "picking up the

apple from the refrigerator has a 90% chance of completing the task", providing a precise

117 In addition, LAC introduces a novel gradient-free policy improvement mechanism that allows the 118 actor to update its policy in the direction suggested by the critics without the computational burden of backpropagation. This not only enhances scalability but also makes the system practical for 119 large-scale LLMs, enabling efficient decision-making in complex environments. 120

- 121 In summary, this work advances the state of LLM-based decision-making through the following key 122 contributions:
  - We introduce a dual-critic actor-critic framework, where a language-based critic provides contextual feedback and a value-based critic ensures quantitative optimization. This combination enables more robust decision-making by leveraging the complementary strengths of both critics.
    - We design an effective value-based critic estimation approach that extracts internal information from LLMs to provide robust numerical evaluation for candidate actions.
    - We propose a novel gradient-free policy improvement method that reduces computational overhead while effectively refining the actor's policy based on both qualitative and quantitative feedback from the dual critics.
    - We demonstrate the effectiveness and generality of LAC across diverse environments, including high-level decision-making tasks in (ALFWorld, (Shridhar et al., 2021)) and low-level action space (BabyAI-Text, (Carta et al., 2023b)). Empirical results show that our approach consistently outperforms state-of-the-art methods such as ReAct Yao et al. (2022b) with GPT-4 (Achiam et al., 2023), even when using much smaller 7B/8B LLMs.

# 137 138 139

140

108

109

110

123

124

125

126

127

128

129 130

131

132

133

134

135

136

#### 2 **RELATED WORK**

141 Large Language Models for Sequential Decision-Making Sequential decision-making is a 142 fundamental ability of intelligent agents, involving generating a series of actions to achieve the goal 143 (Barto et al., 1989; Littman, 1996; McCarthy et al., 1963; Bylander, 1994). Recently, LLM-based 144 agents have been widely used for decision-making in many areas, which only needs some instructions 145 or few-shot examples to generalize to completely new tasks (Huang et al., 2022b; Singh et al., 2023; 146 Ding et al., 2023), thanks to the pre-training on large-scale dataset. According to the functionality, 147 the LLMs that most previous work used mainly belong to two roles: actors, which take trajectories 148 as input and output actions, and critics, which take both trajectories and actions as input and output evaluations of actions. Based on this classification, most of the earlier work in this line of research is 149 actor-only (Ahn et al., 2022; Huang et al., 2022b; Yao et al., 2022b; Huang et al., 2022a; Shinn et al., 150 2024), *i.e.*, directly using the action generated by LLMs based on previous trajectory. Due to the 151 auto-regressive nature of LLM, it does not do reasoning and planning explicitly. Accordingly, LLM 152 with actor-only methods often struggles with complex tasks that require multiple steps of planning 153 and reasoning (Huang & Chang, 2022; Mialon et al., 2023). To overcome this hurdle, another line of 154 work, critic-only uses another LLM to evaluate each action by simulating the consequence of it and 155 then choose the action with the best-simulated outcome (Hao et al., 2023; Liu et al., 2023; Fu et al., 156 2024). However, both actor-only and critic-only methods ignore the interrelation between actor and 157 critic, prioritize one over the other, and insufficiently exploit the available knowledge from the actor 158 and critic. Previous work that tries to combine actor and critic (Zhang et al., 2023b) only uses the language-based outputs of critics. Some work in other fields, such as decoding (Xie et al., 2024), also 159 uses numerical outputs of critics, but it cannot be directly adapted to decision-making problems. To 160 address these limitations, our method LAC integrates prior actor-only and critic-only methods and 161 utilizes the merits of the actor-critic algorithm with the strengths of the LLMs.

162 Large Language Models with Reinforcement Learning Classical sequential decision-making 163 methods, such as Reinforcement Learning (RL), have been widely adopted in embodied environments 164 (Schulman et al., 2017; Fujimoto et al., 2018; Huang et al., 2020; Dong et al., 2022). However, these 165 RL-based methods are typically sample-inefficient and require lots of samples for training. On the 166 other hand, LLMs that contain rich prior knowledge about the world may alleviate this burden. To combine RL and LLM, one straightforward way is to use LLMs as base models and add policy/value 167 heads on top of LLMs (Carta et al., 2023a; Tan et al., 2024). Then use classical RL methods like 168 PPO (Schulman et al., 2017) for training (Szot et al., 2023; Zhou et al.). However, these methods still require lots of training samples of the same tasks, which reduces the benefits of using LLM to 170 some extent and contradicts our settings. There are also other paradigms for combining. RLEM 171 (Zhang et al., 2024) adopts Q-learning (Watkins & Dayan, 1992) and an experience memory to 172 update policies, but it may get stuck in the tasks with extremely sparse rewards like ALFWorld and 173 BabyAI-Text. Retroformer (Yao et al., 2023) trains a smaller LLM with PPO to generate suitable 174 prompts for a larger LLM for a specific task, while our method only needs a small model. ICPI 175 (Brooks et al., 2024) uses LLMs to implement policy iteration by predicting future trajectories and 176 accumulating future rewards, which may also struggle with sparse reward settings. We have compared 177 it empirically in Section 5.

178 179

# 3 PRELIMINARY

180 181 182

In this section, we describe the task setting, previous *actor-only* methods, and *critic-only* methods to
better understand prior work's limitations and the motivations for our LAC. For better understanding,
we compare the frameworks of these methods in Figure 2.

186 Task setup. Consider a general setup of an agent interacting with an environment for achieving 187 some given goals, e.g., goal g = "put a clean egg in microwave" (from benchmark ALFWorld) or 188 goal q = "pick up the green ball" (from benchmark BabyAI-Text). At time step t, the agent receives 189 an observation  $o_t \in \mathcal{O}$ , which is described in natural language in this work, from the environment. 190 The agent then takes an action  $a_t \in \mathcal{A}$  that is sampled from some policy  $\pi(a|q, h_t)$ , where  $h_t :=$  $(o_1, a_1, o_2, a_2, \cdots, o_t)$  is the history to the agent. During execution, there is no immediate reward. 191 The environment will give a signal about whether the task was completed successfully or not only at 192 the end of each episode. The agent has never seen the testing tasks before and can only try each task 193 once in this work. 194

195 Actor-only methods. To solve the above tasks with large language models, one simple method is to 196 directly use pre-trained LLMs as policy:  $a_t \leftarrow \arg \max_a \pi_{LLM}(a|g,h_t)$ , as shown in Figure 2 (a). We also provide detailed algorithm description of *actor-only* methods in Algorithm 2 of Appendix B.3, 197 which can be implemented by simply injecting instructions or few-show examples to the prompt as 198 shown in Yao et al. (2022b). Despite its simplicity, the actor-LLM  $\pi_{LLM}$  generates actions solely 199 relying on its auto-regression ability and it does not conduct long-term planning explicitly, which is 200 typically necessary for sequential decision-making tasks. Additionally, this issue will be exacerbated 201 when using smaller models like CodeLlama-7B (Roziere et al., 2023) and Mistral-7B (Jiang et al., 202 2023). This problem is verified in Section 5. 203

*Critic-only* methods. To handle the issue of lack of long-term planning in *actor-only* methods, 204 another line of research resorts to critic-only methods (Hao et al., 2023; Liu et al., 2023; Fu et al., 205 2024). The basic idea of *critic-only* methods is to first sample several candidate actions from actor 206  $\{a_t^1, a_t^2, \cdots, a_t^n\} \sim \pi_{LLM}(\cdot | g, h_t)$ , then self-evaluate each candidate action by other LLMs and 207 finally select the action with the highest evaluation value. We call it *critic-only* because only the 208 critic's output is considered when choosing the final action. The self-evaluation procedure is the key to 209 *critic-only* methods, which can adopt many approaches. For example, directly ask an LLM to evaluate 210 the action candidate (Fu et al., 2024), or predict the future trajectory  $u_t$  of each action candidate using 211 an LLM as a forward model  $f_{LLM}$  and use the future outcome as evaluation (shown in Figure 2 212 (b)), or use tree-search methods like Monte Carlo Tree Search (MCTS) (Kocsis & Szepesvári, 2006; 213 Coulom, 2006) to expand each action candidate Hao et al. (2023). We also provide detailed algorithm description of *critic-only* methods in Algorithm 3 of Appendix B.4. Despite this progress, *critic-only* 214 methods often neglect the knowledge of actor and the interaction between actor and critic, which may 215 lead to ineffective decision-making.



Figure 2: Framework of our LAC. At each time step, LAC selects an action via three steps: (1) given current goal and history, *lang-critic*  $C_{LLM}$  generates language-based judgments on previous actions; (2) the *actor*  $\pi_{LLM}$  samples candidate actions based on the judgments; (3) the *value-critic*  $Q_{LLM}$  provides numerical evaluations for candidate actions by predicting future trajectories. Finally, the action distribution that integrates the actor and critics can be calculated in a gradient-free way.

231

232

233

234

235

239

# 4 Method

240 In this section, we present our LLM-based Actor-Critic (LAC) algorithm that integrates actor and 241 critic to enhance the decision-making ability of large language models. The key idea behind LAC is to 242 improve the actor  $\pi_{LLM}$  based on the evaluations provided by two distinct but complementary critics: 243 the *lang-critic*  $C_{LLM}$  and the *value-critic*  $Q_{LLM}$ . The *lang-critic*  $C_{LLM}$  provides language-based 244 evaluations that contain richer and more interpretable information for assessing actions, but they are 245 difficult to quantify. The value-critic  $Q_{LLM}$  provides the numerical assessment in terms of long-term value or reward but may lack detailed explanations for assessment. By combining both natural 246 language-based evaluations and numerical assessments, LAC ensures that actions are contextually 247 relevant and optimized for long-term success. Figure 1 demonstrates how relying solely on one critic 248 (either  $C_{LLM}$  or  $Q_{LLM}$ ) may still lead to suboptimal actions, whereas combining both critics avoids 249 such mistakes. 250

The policy improvement process for  $\pi_{LLM}$  involves two main steps: (1) improving  $\pi_{LLM}$  with *lang-critic*'s language-based judgments over previous actions by injecting these judgments into  $\pi_{LLM}$ 's prompt so that  $\pi_{LLM}$  could avoid previous mistakes and sample better candidate actions (Section 4.1); (2) further refining  $\pi_{LLM}$  based on the *value-critic*'s numerical assessment through a gradient-free policy improvement procedure (Section 4.2). Our framework is shown in Figure 2. The overall algorithm is outlined in Algorithm 1. We also compare the frameworks of previous *actor-only* methods, *critic-only* methods, and our LAC, respectively, for better understanding in Figure 15 (a-c).

258 259

# 4.1 POLICY IMPROVEMENT WITH lang-critic

260

To enhance the actor's decision-making with contextually grounded feedback, we first prompt the LLM to generate judgments on previous action selections. Given the task goal g and history  $h_t$ , the *lang-critic* generates a short evaluation sentence  $c_t$  such as "I have found object-X. This step is GOOD" or "I should take object-X instead of object-Y first. This step is BAD." These judgments provide hints about whether and why the previous actions were appropriate, which helps guide the actor towards better candidate actions.

To improve the actor's policy  $\pi_{LLM}$ , we then condition  $\pi_{LLM}$  on language-based evaluation  $c_t$ before sampling candidate actions:  $\{a_t^1, a_t^2, \dots, a_t^n\} \sim \pi_{LLM}(\cdot|g, h_t, c_t)$ . By incorporating  $c_t$ , the actor can generate better candidate actions through in-context learning, aligning more closely with task objectives and avoiding previous mistakes. 270 Algorithm 1: LAC: LLM-based Actor-Critic algorithm. 271 **Input:** current task goal q, history  $h_t$ , actor  $\pi_{LLM}$ , forward model  $f_{LLM}$ , language-based critic 272  $C_{LLM}$ , value-based critic  $Q_{LLM}$ , hyperparameter  $\alpha$ , candidate action size n. 273 **Output:** selected action  $a_t^*$ 274  $\triangleright$  generate language-based evaluations (Section 4.1) 1  $c_t \leftarrow \mathcal{C}_{LLM}(g, h_t);$ 275 2  $\{a_t^1, a_t^2, \cdots, a_t^n\} \sim \pi_{LLM}(\cdot | g, h_t, c_t);$  $\triangleright$  generate candidate actions 276 3 for  $i \leftarrow 1, 2, \cdots, n$  do 277  $u_t^i \leftarrow f_{LLM}(g, h_t, a_t^i);$ 4  $\triangleright$  imagine future trajectory 278  $\mathcal{Q}_{LLM}(g, h_t, a_t^i, u_t^i) \leftarrow \log \frac{P(y=+1|g, h_t, a_t^i, u_t^i)}{P(y=-1|g, h_t, a_t^i, u_t^i)};$ ▷ calculate numerical evaluations 5 279 (Section 4.2, Equation (3)) 280 6 end 281  $\tau \pi(a_t^i|g, h_t, c_t) \leftarrow \pi_{LLM}(a_t^i|g, h_t, c_t) \exp(\alpha \mathcal{Q}_{LLM}(g, h_t, a_t^i, u_t^i));$  ▷ update action distribution 282 (Section 4.2, Equation (6)) 283 s  $a_t^* \leftarrow \arg \max_{a_t^i} \pi(a_t^i | g, h_t, c_t)$ 284 285

The key advantage of conditioning on  $c_t$  is that it acts as an intermediate variable (Prystawski et al., 2024), helping the actor perform more effective reasoning. Similar to the Chain-of-Thought mechanism (Wei et al., 2022; Kojima et al., 2022), this enables the actor to adjust its policy based on the feedback provided by the *lang-critic*. For more examples of language-based evaluations, please refer to Table 14 and Table 15 of Appendix B.

# 4.2 POLICY IMPROVEMENT WITH value-critic

Next, we refine the actor's policy using the *value-critic*  $Q_{LLM}$ , which provides numerical evaluations of each candidate action  $a_t^i$ . While the *lang-critic* focuses on providing contextually grounded feedback, the *value-critic* quantitatively estimates the probability of successfully completing the task after executing each action  $a_t^i$ . This value-based assessment is crucial for guiding the actor's decisions to align with long-term rewards, especially in tasks where immediate outcomes do not fully capture the consequences of an action.

In the following, we will first connect the value-critic to the agent's success probability of completing
 the task, then show how LLMs can be used to estimate this value-based evaluation, and finally derive
 a gradient-free policy improvement mechanism using the estimated value-based evaluation.

304 305

306

287

288

289

290

291 292 293

294

# 4.2.1 CONNECT value-critic TO AGENT'S SUCCESS PROBABILITY

Let  $Q_{LLM}(g, h_t, a_t^i)$  be the value-based evaluation of each candidate action  $a_t^i$  given the task goal g and history  $h_t$ . Ideally,  $Q_{LLM}(g, h_t, a_t^i)$  should represent the cumulative rewards an agent can acquire after executing  $a_t^i$ , analogous to the action-value function in conventional RL algorithms. However, in the benchmarks we consider, only binary success or failure signals are provided at the end of each episode, with no intermediate rewards.

To model  $Q_{LLM}(g, h_t, a_t^i)$  similarly to action-value in RL, and to make it easy to estimate using LLMs, we employ a logistic function (Jordan et al., 1995). Let  $P(y = +1|g, h_t, a_t^i) \in [0, 1]$  denote the probability of successfully completing the task goal g after executing action  $a_t^i$ , where y = +1represents a success signal at the end of the episode. Similarly, let  $P(y = -1|g, h_t, a_t^i)$  represent the failure probability. We use the following logistic function to relate  $P(y = +1|g, h_t, a_t^i)$  to  $Q_{LLM}(g, h_t, a_t^i)$ :

318 319

320

$$P(y = +1|g, h_t, a_t^i) = \frac{1}{1 + \exp\left(-\mathcal{Q}_{LLM}(g, h_t, a_t^i)\right)}.$$
(1)

This formulation indicates that the value-based evaluation  $Q_{LLM}(g, h_t, a_t^i)$  is positively correlated with the success probability  $P(y = +1|g, h_t, a_t^i)$ . Higher  $Q_{LLM}(g, h_t, a_t^i)$  values map to a greater likelihood of success, allowing the critic to guide the actor's policy toward actions that maximize long-term success. While other formulations could be used, we found that Equation (1) is both simple and effective for a wide range of tasks. For a comparison of alternative formulations, refer to Figure 8 in Appendix A.3.

# 4.2.2 ESTIMATE value-critic WITH LLMS

To estimate  $Q_{LLM}(g, h_t, a_t^i)$  using LLMs, we perform an equivalent transformation on Equation (1):

$$\mathcal{Q}_{LLM}(g, h_t, a_t^i) = \log \frac{P(y=+1|g, h_t, a_t^i)}{1 - P(y=+1|g, h_t, a_t^i)} = \log \frac{P(y=+1|g, h_t, a_t^i)}{P(y=-1|g, h_t, a_t^i)}.$$
(2)

331 332

350 351 352

353

354

355 356

357

360 361 362

364

369 370

375

327

328

329 330

333 With Equation (2), we can use the LLM to obtain value evaluation  $Q_{LLM}(g, h_t, a_t^i)$  via first esti-334 mating  $P(y = \pm 1 | g, h_t, a_t^i)$ . The basic idea is to prompting the LLM to predict the success/failure 335 probability given the current trajectory  $(g, h_t)$  and action  $a_t^i$ . Specifically, we use the generated 336 probabilities of special paired tokens that contain positive/negative meanings to indicate LLMs' belief in success/failure. For example, let the generated probability of "GOOD" or "SUCCESS" represent 337 positive results  $P(y = +1|g, h_t, a_t^i)$  and let the generated probability of "BAD" or "FAILURE" 338 represent negative results  $P(y = -1|q, h_t, a_t^i)$ . The insight is that if the agent selects actions cor-339 rectly, the LLM that are pre-trained via next-token prediction tends to increase the probability of 340 *positive* tokens internally. Otherwise, the probability of *negative* tokens will increase. Finally, using 341 Equation (2), we can calculate  $\mathcal{Q}_{LLM}(g, h_t, a_t^i)$  for action  $a_t^i$ . 342

To improve the accuracy of  $Q_{LLM}(g, h_t, a_t^i)$ , we introduce future trajectory rollouts using a forward model  $f_{LLM}$ , which can be implemented by prompting LLMs, *e.g.*, adding few-shot examples, or by fine-tuning on these examples. For each candidate action  $a_t^i$ , we roll out several future steps to predict the resulting trajectory  $u_t^i = f_{LLM}(g, h_t, a_t^i)$ . By considering the future trajectory  $u_t^i$ , we obtain more informed estimates of the success and failure probabilities,  $P(y = \pm 1|g, h_t, a_t^i, u_t^i)$ . This approach accounts for the delayed consequences of actions and ensures that  $Q_{LLM}(g, h_t, a_t^i)$ reflects the long-term value of each action:

$$\mathcal{Q}_{LLM}(g, h_t, a_t^i, u_t^i) = \log \frac{P(y = +1|g, h_t, a_t^i, u_t^i)}{P(y = -1|g, h_t, a_t^i, u_t^i)}.$$
(3)

Trajectory rollouts are especially important in tasks where the outcomes of actions may unfold over several steps. By simulating the future impact of actions, the value-critic provides a more accurate assessment, guiding the actor toward actions that maximize the probability of long-term success.

### 4.2.3 IMPROVE actor WITH value-critic

With the value-critic  $Q_{LLM}$ , we can improve the actor's policy using the following optimization problem:

$$\max_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i}) \right] - \frac{1}{\alpha} \mathbb{D}_{KL} \left[ \pi(a_{t}^{i}|g,h_{t},c_{t}) \| \pi_{LLM}(a_{t}^{i}|g,h_{t},c_{t}) \right],$$
(4)

where  $\alpha$  is a hyperparameter controlling the deviation from the original actor  $\pi_{LLM}$ . The KLdivergence term prevents the new actor  $\pi$  from deviating too far from the original policy, balancing the actor's prior knowledge and the value-critic's guidance.

Following prior work (Rafailov et al., 2024; Go et al., 2023; Peng et al., 2019; Jain et al., 2013; Peters & Schaal, 2007), we can show that the optimal solution to the KL-constrained maximization objective in Equation (4) takes the following form:

$$\pi(a_t^i|g, h_t, c_t) = \frac{1}{Z(g, h_t, c_t)} \pi_{LLM}(a_t^i|g, h_t, c_t) \exp\left(\alpha \mathcal{Q}_{LLM}(g, h_t, a_t^i, u_t^i)\right),$$
(5)

where  $Z(g, h_t, c_t) = \sum_{a_t^i} \pi_{LLM}(a_t^i | g, h_t, c_t) \exp(\alpha Q_{LLM}(g, h_t, a_t^i, u_t^i))$  is the partition function. Please refer to Appendix B.1 for a complete derivation. As the partition function does not depend on action  $a_t^i$ , we can ignore it in practice:

$$\pi(a_t^i|g, h_t, c_t) \propto \pi_{LLM}(a_t^i|g, h_t, c_t) \exp\left(\alpha \mathcal{Q}_{LLM}(g, h_t, a_t^i, u_t^i)\right).$$
(6)

We simply take the action with maximum proportion  $a_t \leftarrow \arg \max_{a_t^i} \pi(a_t^i | g, h_t, c_t)$ . It is worth mentioning that if we let  $\alpha = 0$  and remove  $c_t$  in Equation (6), we recover *actor-only* methods. And if we let  $\alpha \to +\infty$ , we recover *critic-only* methods that use our action-value estimation approach. There are two key advantages of using Equation (6). Firstly, it updates the action distribution of *actor*  $\pi_{LLM}$  in the direction suggested by *critic*  $Q_{LLM}$  in a gradient-free way, which achieves policy improvement with much lower computation burden compared to gradient-based methods, especially when the actor is realized by a large model. Secondly, the action distribution of the new actor  $\pi$  is a balanced integration of the actor's prior based on past information and the critic's posterior based on predicted future information.

384

386

389

390

391

392

393

394

# 5 EXPERIMENTS

387 388 In this s

In this section, we benchmark our method LAC on benchmarks that cover both high-level action space (ALFWorld (Shridhar et al., 2021)) and low-level action space (BabyAI-Text (Chevalier-Boisvert et al., 2019)). We evaluate the effectiveness of LAC by answering the following questions: (1) Can LAC outperform other decision-making with LLMs methods? (Section 5.2) (2) How does each component of LAC contributes to its performance? (Section 5.3) (3) How do different large language models influence performance? (Section 5.2 and Section 5.3) (4) Is our method computationally consuming? (Section 5.4).

395 5.1 EXPERIMENT SETUP

We compare our method with various decision-making with LLMs baselines, which can be largely classified into *actor-only* and *critic-only* methods. For more details of baselines, please refer to Appendix C.2.

Actor-only methods. ReAct (Yao et al., 2022b) combines reasoning and acting in the interaction
 with the environment and leverages the reasoning capabilities of LLMs to increase the probability of
 the LLM acting correctly as an actor.

Critic-only methods. RAP(Hao et al., 2023) utilizes LLMs as actor and world models and adopts tree-search planning methods to evaluate each possible action candidate. ICPI (Brooks et al., 2024) implements policy iteration using LLMs by predicting future trajectories and selecting the action with the highest predicted cumulative rewards. RAFA (Liu et al., 2023) evaluates each action candidate by tree-search and selects the action that may complete the most sub-goals.

We evaluate LAC on two decision-making benchmarks with high-level actions and low-level actions, respectively.

Benchmark with high-level actions: ALFWorld. ALFWorld (Shridhar et al., 2021) is a widely
used text-based household environment with 134 different tasks, which require the agent to achieve a
goal through a sequence of high-level actions, *e.g.* "go to place-X", "take object-Y from place-X", *etc.*. The main challenge of this benchmark is to locate the target object and fulfill household work
with commonsense knowledge of LLMs. Following ReAct, we evaluate all 134 unseen evaluation
games in a task-specific setup.

417 Benchmark with low-level actions: BabyAI-Text. BabyAI-Text (Carta et al., 2023b) is a Grid 418 World environment that extended from the BabyAI platform (Chevalier-Boisvert et al., 2019), in 419 which the agent and objects are placed in a room of  $8 \times 8$  tiles. The agent has 6 primitive actions: 420 turn left, turn right, go forward, pick up, drop, toggle, to solve a task described in natural language 421 (*e.g.* "Pick up the red box"). These tasks could be difficult because agents have to make a long-term 422 plan, avoid obstacles and find a short path to target objects based on partial observations that are 423 described in natural language

To show the stability of LAC, we adopt four open-source large language models from different organizations: CodeLlama-7B (Roziere et al., 2023), Mistral-7B (Jiang et al., 2023), Gemma-7B (Team et al., 2024), and Llama-3-8B (Meta, 2024a).

427 428

429

5.2 Performance

We report the results of our method LAC compared with other baselines in ALFWorld and BabyAIText in Figure 3 and Figure 6 of Appendix A.1. For all experiments, we set the temperature of LLMs to 0, hence the generation is deterministic. For this reason, there is no error bar in the figure.



Figure 3: Performance of our LAC compared with various baselines in benchmarks ALFWorld and BabyAI-Text.

LAC outperforms all other baselines on both ALFWorld and BabyAI-Text across different LLMs, and
 LAC is even better than GPT-4+ReAct in most settings, which validates our method's effectiveness
 and stability.

453 LAC's superior performance stems from its balanced integration of the strengths of both actor and 454 critic. While actor-only (e.g., ReAct) methods excel in short-term actions, they often struggle with 455 long-term reasoning. In contrast, critic-only (e.g., RAP) methods conduct explicit reasoning but 456 might mispredict future trajectories and lead to even worse action selection occasionally compared 457 with actor-only methods. LAC addresses these limitations by balancing the actor's action generation 458 and the critic's evaluation. We have provided illustrative examples for ALFWorld and BabyAI-Text in 459 Figure 1 and Figure 13 respectively. In summary, actor-only and critic-only methods make mistakes 460 at different time steps, our LAC can select the correct action.

Regarding the performance of LAC with different base models, we highlight two key findings: (1)
Our method is general and can be adapted to various base models, and (2) stronger base models, such as Gemma-7B, demonstrate higher performance when integrated with our approach. However, due to the incomplete public availability of training details for these base models, further in-depth analysis will require additional investigation.

# 5.3 ABLATION STUDIES





466 467

447

448 449

Figure 4: Ablation studies in benchmarks ALFWorld and BabyAI-Text.

To investigate the contributions of each component of LAC, we conduct elaborate ablation studies. There are two main components that characterize our method: (1) the integration of actor  $\pi_{LLM}$  and  $\begin{array}{ll} \mbox{486} & \mbox{lang-critic } \mathcal{C}_{LLM} \mbox{ before action generation and (2) the integration of actor } \pi_{LLM} \mbox{ and value-critic } \mathcal{Q}_{LLM} \mbox{ after action generation. Therefore, to show the contribution of each component, we design the following ablation studies: (1) LAC w/o lang-critic removes the lang-critic <math>\mathcal{C}_{LLM} \mbox{ from LAC as well as the integration before action generation; (2) LAC w/o value-critic removes the value-critic } \mathcal{Q}_{LLM} \mbox{ from LAC as well as the integration after action generation; (3) Value-critic-only only uses value-critic } \mathcal{Q}_{LLM} \mbox{ for decision-making.} \end{array}$ 

We report the result in Figure 4. LAC is better than all other variants in both ALFWorld and BabyAI-Text. Specifically, the performance decrease in LAC *w/o* lang-critic and LAC *w/o* value-critic compared with LAC verify the effectiveness of lang-critic  $C_{LLM}$  and value-critic  $Q_{LLM}$ , respectively. And the result that Value-critic-only performs worse than LAC also suggests the necessity for integrating actor and critic.

# 5.4 COMPUTATIONAL COST ANALYSIS

497 498

513

514

520 521

522 523

524

525

526 527 528

529



Figure 5: Computational cost analysis of LAC and baselines.

515 Our method integrates the actor and two critics, which may bring extra computational cost per step. 516 In Figure 5, we compare computational costs concerning the number of tokens spent and running 517 time between LAC and other baselines. Specifically, though LAC has a higher computational cost per 518 step due to the extra inference procedure of critics and the forward model, the total cost of LAC is 519 still lower than most LLM-based baselines because LAC requires fewer steps to finish each task.

# 6 DISCUSSION

In this work, we introduce a novel LLM-based Actor-Critic algorithm LAC that integrates the ability of actors and critics as well as exploits the strong prior knowledge in LLMs for sequential decision-making. Compared with previous *actor-only* and *critic-only* methods, LAC achieves high performance on ALFWorld and BabyAI-Text even using small open-source LLMs.

# References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea
  Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say:
  Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- Andrew Gehret Barto, Richard S Sutton, and CJCH Watkins. *Learning and sequential decision making*, volume 89. University of Massachusetts Amherst, MA, 1989.
- 539 Ethan Brooks, Logan Walls, Richard L Lewis, and Satinder Singh. Large language models can implement policy iteration. *Advances in Neural Information Processing Systems*, 36, 2024.

540 Tom Bylander. The computational complexity of propositional strips planning. Artificial Intelligence, 541 69(1-2):165-204, 1994. 542 Thomas Carta, Clément Romac, Thomas Wolf, Sylvain Lamprier, Olivier Sigaud, and Pierre-Yves 543 Oudeyer. Grounding large language models in interactive environments with online reinforcement 544 learning. In International Conference on Machine Learning, pp. 3676–3713. PMLR, 2023a. 546 Thomas Carta, Clément Romac, Thomas Wolf, Sylvain Lamprier, Olivier Sigaud, and Pierre-Yves 547 Oudeyer. Grounding large language models in interactive environments with online reinforcement 548 learning, 2023b. 549 Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, 550 Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence 551 modeling. Advances in neural information processing systems, 34:15084–15097, 2021. 552 Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Saharia, 553 Thien Huu Nguyen, and Yoshua Bengio. Babyai: A platform to study the sample efficiency of 554 grounded language learning, 2019. 555 556 Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In International conference on computers and games, pp. 72-83. Springer, 2006. 558 Yan Ding, Xiaohan Zhang, Chris Paxton, and Shiqi Zhang. Task and motion planning with large 559 language models for object rearrangement. In 2023 IEEE/RSJ International Conference on 560 Intelligent Robots and Systems (IROS), pp. 2086–2092. IEEE, 2023. 561 562 Heng Dong, Tonghan Wang, Jiayuan Liu, and Chongjie Zhang. Low-rank modular reinforcement 563 learning via muscle synergy. arXiv preprint arXiv:2210.15479, 2022. 564 Dayuan Fu, Jianzhao Huang, Siyuan Lu, Guanting Dong, Yejie Wang, Keqing He, and Weiran Xu. Pre-565 act: Predicting future in react enhances agent's planning ability. arXiv preprint arXiv:2402.11534, 566 2024. 567 Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-568 critic methods. In International conference on machine learning, pp. 1587–1596. PMLR, 2018. 569 570 Dongyoung Go, Tomasz Korbak, Germán Kruszewski, Jos Rozen, Nahyeon Ryu, and Marc Dymet-571 man. Aligning language models with preferences through f-divergence minimization. arXiv 572 preprint arXiv:2302.08215, 2023. 573 Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 574 Reasoning with language model is planning with world model. arXiv preprint arXiv:2305.14992, 575 2023. 576 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 577 and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. 578 579 Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey. 580 arXiv preprint arXiv:2212.10403, 2022. 581 Wenlong Huang, Igor Mordatch, and Deepak Pathak. One policy to control them all: Shared 582 modular policies for agent-agnostic control. In International Conference on Machine Learning, pp. 583 4455-4464. PMLR, 2020. 584 585 Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot 586 planners: Extracting actionable knowledge for embodied agents. In International Conference on Machine Learning, pp. 9118–9147. PMLR, 2022a. 588 Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through 590 planning with language models. arXiv preprint arXiv:2207.05608, 2022b. 591 Ashesh Jain, Brian Wojcik, Thorsten Joachims, and Ashutosh Saxena. Learning trajectory preferences 592 for manipulators via iterative improvement. Advances in neural information processing systems, 26, 2013.

594 595 596	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> , 2023.
597 598 599	Michael I Jordan et al. Why the logistic function? a tutorial discussion on probabilities and neural networks, 1995.
600 601	Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In <i>European conference</i> on machine learning, pp. 282–293. Springer, 2006.
603 604 605	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35: 22199–22213, 2022.
606 607	Michael Lederman Littman. Algorithms for sequential decision-making. Brown University, 1996.
608 609 610 611	Zhihan Liu, Hao Hu, Shenao Zhang, Hongyi Guo, Shuqi Ke, Boyi Liu, and Zhaoran Wang. Reason for future, act for now: A principled framework for autonomous llm agents with provable sample efficiency. <i>arXiv preprint arXiv:2309.17382</i> , 2023.
612	John McCarthy et al. Situations, actions, and causal laws. Comtex Scientific, 1963.
614	Meta. Meta llama 3. https://llama.meta.com/llama3/, 2024a.
615 616	Meta. Meta llama 3.1. https://ai.meta.com/blog/meta-llama-3-1/, 2024b.
617 618 619	Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ram Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, et al. Augmented language models: a survey. <i>arXiv preprint arXiv:2302.07842</i> , 2023.
620 621 622	Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. <i>arXiv preprint arXiv:1910.00177</i> , 2019.
623 624 625	Jan Peters and Stefan Schaal. Reinforcement learning by reward-weighted regression for operational space control. In <i>Proceedings of the 24th international conference on Machine learning</i> , pp. 745–750, 2007.
626 627 628	Ben Prystawski, Michael Li, and Noah Goodman. Why think step by step? reasoning emerges from the locality of experience. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
629 630 631	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
632 633 634 635	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. <i>arXiv preprint arXiv:2308.12950</i> , 2023.
636 637 638	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> , 2017.
639 640 641	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. <i>Advances in Neural Information Processing</i> <i>Systems</i> , 36, 2024.
642 643 644 645	Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 10740–10749, 2020.
646 647	Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning, 2021.

648 649 650 651	Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large language models. In 2023 IEEE International Conference on Robotics and Automation (ICRA), pp. 11523–11530. IEEE, 2023.
652 653 654 655 656	Andrew Szot, Max Schwarzer, Harsh Agrawal, Bogdan Mazoure, Rin Metcalf, Walter Talbott, Natalie Mackraz, R Devon Hjelm, and Alexander T Toshev. Large language models as generalizable policies for embodied tasks. In <i>The Twelfth International Conference on Learning Representations</i> , 2023.
657 658 659	Weihao Tan, Wentao Zhang, Shanqi Liu, Longtao Zheng, Xinrun Wang, and Bo An. True knowledge comes from practice: Aligning llms with embodied environments via reinforcement learning. <i>arXiv</i> preprint arXiv:2401.14151, 2024.
660 661 662 663	Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. <i>arXiv preprint arXiv:2403.08295</i> , 2024.
664 665 666	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023.
667	Christopher JCH Watkins and Peter Dayan. Q-learning. Machine learning, 8:279–292, 1992.
669 670 671	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.
672 673 674	Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, James Xu Zhao, Min-Yen Kan, Junxian He, and Michael Xie. Self-evaluation guided beam search for reasoning. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
675 676 677 678	Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. <i>Advances in Neural Information Processing Systems</i> , 35:20744–20757, 2022a.
679 680 681	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> , 2022b.
682 683 684 685	Weiran Yao, Shelby Heinecke, Juan Carlos Niebles, Zhiwei Liu, Yihao Feng, Le Xue, Rithesh Murthy, Zeyuan Chen, Jianguo Zhang, Devansh Arpit, et al. Retroformer: Retrospective large language agents with policy gradient optimization. <i>arXiv preprint arXiv:2308.02151</i> , 2023.
686 687 688	Bin Zhang, Hangyu Mao, Jingqing Ruan, Ying Wen, Yang Li, Shao Zhang, Zhiwei Xu, Dapeng Li, Ziyue Li, Rui Zhao, et al. Controlling large language model-based agents for large-scale decision-making: An actor-critic approach. <i>arXiv preprint arXiv:2311.13884</i> , 2023a.
689 690 691	Bin Zhang, Hangyu Mao, Jingqing Ruan, Ying Wen, Yang Li, Shao Zhang, Zhiwei Xu, Dapeng Li, Ziyue Li, Rui Zhao, et al. Controlling large language model-based agents for large-scale decision-making: An actor-critic approach. <i>arXiv preprint arXiv:2311.13884</i> , 2023b.
693 694	Danyang Zhang, Lu Chen, Situo Zhang, Hongshen Xu, Zihan Zhao, and Kai Yu. Large language models are semi-parametric reinforcement learning agents. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.

- Yifei Zhou, Andrea Zanette, Jiayi Pan, Sergey Levine, and Aviral Kumar. Archer: Training language model agents via hierarchical multi-turn rl. In Forty-first International Conference on Machine Learning.

# A EXTRA RESULTS



# A.1 RESULTS OF OTHER TASKS IN BABYAI-TEXT

Figure 6: Performance of our LAC compared with various baselines in all tasks from BabyAI-Text.





Figure 7: Performance of LAC in benchmark WebShop

#### 773 **RESULTS OF LAC IN WEBSHOP** A.2

771 772

774

786

787

803

804

We have conducted new experiments using the WebShop benchmark (Yao et al., 2022a), which 775 presents a scenario with a potentially infinite action space. This benchmark requires an agent to 776 purchase a product based on specific instructions (e.g. "I need a long clip-in hair extension which 777 is natural looking, and price lower than 20.00 dollars") through web interactions (e.g. search "long 778 clip-in hair extension", choose buttons such as "click [item ID]" or "back to search"). Within this 779 context, the "search" and "click" actions can indeed lead to an unbounded set of potential actions, as 780 the agent can continuously refine its queries and selections based on dynamic web results. 781

The results of our experiments are detailed in Figure 7. We found that our method, LAC, consistently 782 outperforms other baselines, in terms of both accumulated reward and success rate across various base 783 models. This demonstrates the robustness of our method in handling more complex and open-ended 784 action spaces. 785



#### Results of using other definition of $\mathcal{Q}_{LLM}$ A.3



Figure 8: Performance of LAC when using different definition of value-critic  $Q_{LLM}$ 

In LAC we define lang-critic  $Q_{LLM}$  as  $Q_{LLM}(g, h_t, a_t^i) = \log \frac{P(y=+1|g, h_t, a_t^i)}{P(y=-1|g, h_t, a_t^i)}$ . There are also other definitions, for example, the simplest variant is  $Q_{LLM}(g, h_t, a_t^i) = \log P(y=+1|g, h_t, a_t^i)$ . 805 807

In this subsection, we provide a performance comparison between them in Figure 8. LAC outperforms 808 the variant in most situations across tasks and models. We speculate that this is because LAC uses 809 more information, *i.e.*, both  $P(y = +1|g, h_t, a_t^i)$  and  $P(y = -1|g, h_t, a_t^i)$ , than the variant, and the evaluation might be more accurate and more stable. There might be other definitions of  $Q_{LLM}$  and among them, our  $Q_{LLM}$  is simple and effective.

812 813

814

827

828

829

# A.4 COMPARISON OF LAC WITH MORE BASELINES ON ALFWORLD



Figure 9: Performance of our LAC and LAC's variants compared with various baselines in benchmark ALFWorld.

In this subsection, we compare LAC with more baselines including some traditional RL methods implemented using LLMs on ALFWorld (Shridhar et al., 2021). The comparison is shown in Figure 9.

While in LAC we fine-tune the critic using a few trajectories, it is also possible to fine-tune the actor 833 to generate actions using those trajectories. Therefore, a potential baseline could be fine-tuning the 834 actor in actor-only method. To demonstrate the improvement brought by fine-tuning the actor, we 835 fine-tune the actor in ReAct (Yao et al., 2022b) and show the results in Figure 9. We also show 836 the results of LAC w/ fine-tuned actor in Figure 9. In brief, ReAct w/ fine-tuned actor is a strong 837 baseline compared with other baselines, but still inferior to our method LAC and LAC w/ fine-tuned 838 actor. Compared to LAC, the underperformance of LAC w/ fine-tuned actor arises from its tendency 839 to overfit the training trajectories. This overfitting causes the actor to favor actions that are more 840 frequent in the dataset, potentially leading to suboptimal action selection. 841

For example, in the ALFWorld training dataset, the action "take an apple from X" occurs frequently. After fine-tuning, the actor may disproportionately generate this action, even when it is irrelevant to the current goal. One case is that the current goal is to "heat some egg and put it in the garbage can". When the agent sees an "apple 2" in "fridge 1", it generates and selects an irrelevant action "take apple 2 from fridge 1", which does not align with the task.

This tendency towards overfitting arises because the complexity of the policy function, which maps states *s* to actions *a*, often exceeds that of the critic. The policy often has to capture a wide variety of potential actions for each state, particularly in complex environments. However, the quite limited training dataset in our setting restricts its ability to generalize effectively, resulting in memorization of specific actions rather than flexible decision-making. In contrast, our critic, which includes a world model for rollout and an evaluation function, focuses on capturing more predictable dynamics of the environment and simpler evaluation criteria. This typically requires simpler mappings than those needed for the policy, thus avoiding overfitting.

We also include some LLM-based RL variants as baselines to show the superiority of LAC over conventional RL algorithms. We design three LLM-based RL variants that are built upon pre-trained LLMs and directly extract actions/values information from LLMs without adding action/value heads, namely Conventional Policy Gradient, Conventional Policy Gradient w/ dense rewards and Conventional Actor-Critic in Figure 9.

For the implementation of the Conventional Policy Gradient, we need the probability of actions and the returns. To obtain the probability of actions, we directly use LLM to compute the conditional probability of each token in action  $a_i = [w_1, w_2, \dots, w_{|a_t|}]$  given the goal g, history  $h_t$  and then calculate their product:

$$\pi(a_t|g, h_t) = \prod_{j=1}^{|a_t|} P_{LLM}(w_j|g, h_t, w_{< j})$$

in which  $P_{LLM}(w_j|g, h_t, w_{< j})$  is the probability of token  $w_j$  given goal g, history  $h_t$  and previous tokens  $w_{< j}$  computed by LLM. Then we regard the cumulative future rewards as the return  $G_t$ , which is +1 for successful trajectories and -1 for failed trajectories in the tasks we considered. Finally, the gradient of policy is  $\mathbb{E}[\sum_t \nabla \log \pi(a_t|g, h_t)G_t]$ . Conventional Policy Gradient w/ dense rewards is similar to Conventional Policy Gradient except that we manually add intermediate rewards for each step, and then use the cumulative future rewards as the return  $G_t$ .

For the implementation of the Conventional Actor-Critic, we additionally need a critic to estimate action values. As it is possible to train a new value head using only 18 trajectories, we instead approximate the action value similar to  $Q_{LLM}$  in our method LAC, *i.e.*  $Q_{LLM}(g, h_t, a_t, u_t) =$  $\log \frac{P_{LLM}(y=\pm 1|g,h_t,a_t,u_t)}{P_{LLM}(y=\pm 1|g,h_t,a_t,u_t)}$ , in which  $P_{LLM}(y=\pm 1|g,h_t,a_t,u_t)$  is the output probability of special positive/negative tokens like GOOD or BAD that indicate positive/negative results as LLM's belief on success/failure. Finally, the gradient of policy is  $\mathbb{E}[\sum_t \nabla \log \pi(a_t|g,h_t)Q_{LLM}(g,h_t,a_t,u_t)]$ .

In summary, Conventional Policy Gradient exhibits almost all zero performance, which is due to the extremely sparse reward problems, compared with Conventional Policy Gradient w/ dense rewards. Conventional Actor-Critic demonstrates non-zero performance only on some stronger LLMs like Gemma-7B (Team et al., 2024), Llama-3-8B (Meta, 2024a) and Llama-3.1-8B (Meta, 2024b), which may be because the optimization method of conventional actor-critic is not suitable in insufficient data settings.

In addition to the aforementioned LLM-based RL variant, Decision Transformer (Chen et al., 2021) 883 is also a potential solution in combining RL and transformer-based LLMs. We fine-tune pretrained 884 LLMs in a similar way as conventional decision transformers. We construct a dataset using decision-885 transformers' trajectory representation:  $\tau = [R_1, s_1, a_1, R_2, s_2, a_2, \cdots]$ , in which  $R_t$  is return-to-go, 886 *i.e.*, +1 for successful trajectories and -1 for failed trajectories in our extremely sparse reward settings. 887 Then we fine-tune LLMs with next-token prediction loss on these trajectories. During execution, we insert +1 before state  $s_t$  to specify the desired outcome. The results are shown in Figure 9 as 889 Conventional Decision Transformer. In short, Conventional Decision Transformer exhibits a similar 890 performance to ReAct, which may be because the 18 trajectories are insufficient for fine-tuning 891 decision transformers.

Our method LAC is better than all considered baselines because of its ability to handle extremely
 sparse reward problems using LLM's prior knowledge and to fully utilize insufficient data.

895 896

897

# A.5 MORE ABLATIONS OF LAC IN ALFWORLD AND BABYAI-TEXT

The *value-critic* relies on future trajectory predictions to improve the accuracy of its evaluations. By predicting future trajectories, the critic considers long-term consequences and evaluates actions more effectively, which ultimately leads to better decision-making.

For a full comparison, here we conducted an extra experiment for LAC *w/o* rollout, in which the *value-critic* generate value-based evaluations without future trajectory predictions. The results, included in Figure 10 and Figure 11, show that LAC *w/o* rollout consistently underperforms compared to the full LAC across various base models. This finding emphasizes the importance of future trajectory predictions for accurate evaluations.

906 907

908

# A.6 COMPUTATIONAL COST ANALYSIS OF LAC WITH MORE BASELINES IN ALFWORLD

In this subsection, we demonstrate the computation cost of LAC and other baselines in Figure 12.
We show the success rate, steps per task, time used per task, and token cost per task respectively.
Specifically, though LAC has a higher computational cost per step due to the extra inference procedure of critics and the forward model, the total cost is still lower than most LLM-based baselines because LAC has a higher success rate and requires fewer steps to finish each task.

914 915

# A.7 ILLUSTRATION OF BABYAI-TEXT

916 917

We should the illustrative example of BabyAI-Text in Figure 13.



# Figure 11: More ablation studies of LAC in BabyAI-Text

Gemma-7B

Different Large Language Models

Llama-3-8B

Mistral-7B

# A.8 RESULTS OF DIFFERENT CRITIC IMPROVEMENT METHODS

CodeLlama-7B

946

947 948

949

950

951 952

953 954 955

956 957

958

959

960

961

962

963

964 965 0.2

0.1

0.0

To improve the critics given only several examples, we can fine-tune the open-source models via nexttoken prediction. Please refer to Appendix B.2 for more fine-tuning details. To show the effectiveness of fine-tuning, we present the performance of LAC and other variants on task "go to" and "pick up" from BabyAI-Text when we just add these examples into the prompt, *i.e.*, in-context learning, in Table 1 and Table 2. We also show the the performance improvement if we do fine-tuning in the parentheses of Table 1. This result indicates that (1) fine-tuning can incorporate extra knowledge into LLMs better than in-context learning in our case (2) both of our two critics can benefit from fine-tuning. It is worth mentioning that our LAC still outperforms baselines without fine-tuning.

Table 1: Performance of two critic improvement methods: in-context learning or fine-tuning.

066	10010 11 1 011011100100				ng of this tani
900		Codel lama 7B	Gamma 7B	I lama 3 8P	Mistral 7B
967		Couellania-/D	Gennia-/B	Liama-3-0D	Wilsu al- / D
968	LAC	0.30 († 0.16)	0.62 († 0.14)	0.32 († 0.34)	0.24 († 0.46)
969	LAC w/o lang-critic	0.30 († 0.02)	0.58 († 0.02)	0.38 († 0.10)	0.26 († 0.28)
970	LAC w/o value-critic	0.28 († 0.06)	0.48 († 0.04)	0.42 († 0.10)	0.10 († 0.28)
971	Value-critic-only	0.42 († 0.04)	0.40 († 0.24)	0.34 († 0.22)	0.38 († 0.24)



Figure 12: Computational cost analysis of our LAC compared with various baselines in benchmarks ALFWorld. Though LAC may have a higher computational cost per step due to the extra inference procedure of critics and the forward model, the total cost of LAC is still lower than most LLM-based baselines because LAC requires fewer steps to finish each task.

Table 2: Performance of two critic improvement methods: in-context learning or fine-tuning.

	*			0
	CodeLlama-7B	Gemma-7B	Llama-3-8B	Mistral-7B
LAC	0.20 († 0.06)	0.22 († 0.20)	0.34 († 0.08)	0.20 († 0.16)
LAC w/o lang-critic	0.16 († 0.08)	0.32 († 0.04)	0.32 († 0.04)	0.22 († 0.06)
LAC w/o value-critic	0.12 († 0.14)	0.36 († 0.20)	0.28 († 0.06)	0.26 († 0.04)
Value-critic-only	0.22 († 0.04)	0.24 († 0.26)	0.16 († 0.16)	0.16 († 0.26)

A.9 ANALYSIS OF THE FINE-TUNING PROCESS IN LAC

In order to improve the quality of the language-based critic generated by LLM, we finetune the LLM that generates the critic. In this section, we analyze the finetuning process, showing the impact of finetuning on critic prediction, as well as the impact of different data amounts and positive and negative sample ratios on task success rates. The comparison can be seen in Figure 14.

In Figure 14 (a), we show the influence of fine-tuning data size. We use 9, 18, 27 and 36 trajectories to fine-tune LLMs, and show the final success rate on 134 evaluation tasks. In summary, larger data sizes (27 or 36 trajectories) generally bring higher success rate, while small data sizes (18 and even 9 trajectories in some cases) are already enough for LAC to achieve outperformance.

Figure 14 (b) shows the influence of different positive/negative sample ratio (positive:negative = 0:1, 1:3, 1:1, 3:1 and 1:0) on final performance. We keep the total number of samples the same and just change positive/negative ratio. In short, our LAC is robust to reasonable positive/negative ratios (*e.g.* 1:3, 1:1, 3:1), while LAC based on CodeLlama-7B (Roziere et al., 2023) and Gemma-7B (Team et al., 2024) even perform better when given all positive samples (1:0).

Figure 14 (c) shows the learning curves of the fine-tuning process. We plot the next prediction loss and positive/negative tokens prediction accuracy for CodeLlama-7B (Roziere et al., 2023). In short, as the next token prediction loss decreases during fine-tuning, the accuracy of predicting the special tokens (GOOD or BAD) increases, which exhibits the effect of the fine-tuning process.

1018

985

986

987

988

989 990 991

- 1019 1020 A.10 Results of different hyper-parameter  $\alpha$
- 1021 1022 The hyper-parameter  $\alpha$  in Equation (4) controls the deviation from the original actor  $\pi_{LLM}$ . In this 1023 subsection, we grid-search this hyper-parameter over  $\{1/2, 1, 2, 5, 10\}$  in task "go to" of BabyAI-1024 Text, then we fix  $\alpha$  for other tasks:  $\alpha = 1$  for model CodeLlama-7B,  $\alpha = 2$  for model Gemma-7B,  $\alpha = 2$  for model Llama-3-8B and  $\alpha = 10$  for model Mistral-7B.
  - As for benchmark ALFWorld, we fixed  $\alpha = 1$  in all experiments.



Figure 13: An illustrative explanation of our method LAC in BabyAI-Text. The histogram on the right shows the action probabilities of different methods. While actor ( $\pi_{LLM}$ ) and critics (*lang-critic*  $C_{LLM}$ , value-critic  $Q_{LLM}$ ) make mistakes at different time steps, LAC (ours) can select the correct action by integrating actor and critics. Please refer to Table 16 for the full trajectory.

Table 3: Results of different hyper-parameter  $\alpha$ 

	CodeLlama-7B	Gemma-7B	Llama-3-8B	Mistral-7B
LAC ( $\alpha = 1/2$ )	0.46	0.54	0.62	0.68
LAC ( $\alpha = 1$ )	0.46	0.62	0.64	0.58
LAC ( $\alpha = 2$ )	0.44	0.76	0.66	0.64
LAC ( $\alpha = 5$ )	0.46	0.72	0.62	0.64
LAC ( $\alpha = 10$ )	0.40	0.58	0.60	0.70

1072 B METHOD DETAILS

B.1 DERIVING THE SOLUTION OF THE KL-CONSTRAINED MAXIMIZATION OBJECTIVE

In this subsection, we will derive Equation (5). We optimize the following objective:

$$\max_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i}) \right] - \frac{1}{\alpha} \mathbb{D}_{KL} \left[ \pi(a_{t}^{i}|g,h_{t},c_{t}) \| \pi_{LLM}(a_{t}^{i}|g,h_{t},c_{t}) \right].$$



Figure 14: Analysis regarding the fine-tuning process of our LAC. (a) Influence of the fine-tuning data size. Larger data sizes (27, 36 trajectories) generally bring higher performance, but small data sizes (18 and even 9 trajectories) are already enough for our method to achieve outperformance. (b) Influence of the positive/negative data ratio. LAC is robust to reasonable positive/negative ratios (1:3, 1:1, 3:1) while CodeLlama-7B and Gemma-7B-based LAC even perform better given all positive data (1:0). (c) Learning curves of next-token prediction loss and positive/negative tokens prediction accuracy for CodeLlama-7B and ALFWorld.

1099 We now have:

$$\begin{aligned} & \max_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i}) \right] - \frac{1}{\alpha} \mathbb{D}_{KL} \left[ \pi(a_{t}^{i}|g,h_{t},c_{t}) \| \pi_{LLM}(a_{t}^{i}|g,h_{t},c_{t}) \right] \\ & = \max_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i}) - \frac{1}{\alpha} \log \frac{\pi(a_{t}^{i}|g,h_{t},c_{t})}{\pi_{LLM}(a_{t}^{i}|g,h_{t},c_{t})} \right] \\ & = \min_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \log \frac{\pi(a_{t}^{i}|g,h_{t},c_{t})}{\pi_{LLM}(a_{t}^{i}|g,h_{t},c_{t})} - \alpha \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i}) \right] \\ & = \min_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \log \frac{\pi(a_{t}^{i}|g,h_{t},c_{t})}{\pi_{LLM}(a_{t}^{i}|g,h_{t},c_{t})} - \alpha \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i}) \right] \\ & = \min_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \log \frac{\pi(a_{t}^{i}|g,h_{t},c_{t})}{\frac{1}{Z(g,h_{t},c_{t})}\pi_{LLM}(a_{t}^{i}|g,h_{t},c_{t})} \exp(\alpha \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i}))} - \log Z(g,h_{t},c_{t}) \right] \\ & = \min_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \log \frac{\pi(a_{t}^{i}|g,h_{t},c_{t})}{\frac{1}{Z(g,h_{t},c_{t})}\pi_{LLM}(a_{t}^{i}|g,h_{t},c_{t})} \exp(\alpha \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i}))} - \log Z(g,h_{t},c_{t}) \right] \\ & = \max_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \log \frac{\pi(a_{t}^{i}|g,h_{t},c_{t})}{\frac{1}{Z(g,h_{t},c_{t})}\pi_{LLM}(a_{t}^{i}|g,h_{t},c_{t})} \exp(\alpha \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i}))} \right] \\ & = \max_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \log \frac{\pi(a_{t}^{i}|g,h_{t},c_{t})}{\frac{1}{Z(g,h_{t},c_{t})}\pi_{LLM}(a_{t}^{i}|g,h_{t},c_{t})} \exp(\alpha \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i}))} \right] \\ & = \max_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \log \frac{\pi(a_{t}^{i}|g,h_{t},c_{t})}{\frac{1}{Z(g,h_{t},c_{t})}\pi_{LLM}(a_{t}^{i}|g,h_{t},c_{t})} \exp(\alpha \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i})} \right] \\ & = \max_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \exp(\alpha \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i})} \exp(\alpha \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i})} \right] \\ & = \max_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \exp(\alpha \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i})} \exp(\alpha \mathcal{Q}_{LM}(g,h_{t},a_{t}^{i},u_$$

1109 1110 1111

1112 1113

1117

1098

where we have the partition function:

$$Z(g, h_t, c_t) = \sum_{a_t^i} \pi_{LLM}(a_t^i | g, h_t, c_t) \exp\left(\alpha \mathcal{Q}_{LLM}(g, h_t, a_t^i, u_t^i)\right)$$

Since the partition function is a function of only  $g, h_t$  and the original actor  $\pi_{LLM}$ , but does not depend on the optimized actor  $\pi$ , we define

$$\pi^*(a_t^i|g, h_t, c_t) = \frac{1}{Z(g, h_t, c_t)} \pi_{LLM}(a_t^i|g, h_t, c_t) \exp\left(\alpha \mathcal{Q}_{LLM}(g, h_t, a_t^i, u_t^i)\right).$$

This definition of policy if a valid probability distribution as  $\pi^*(a_t^i|g, h_t, c_t)$  for all  $a_t^i$  and  $\sum_{a_t^i} \pi^*(a_t^i|g, h_t, c_t) = 1$ . As  $Z(g, h_t, c_t)$  is not a function of  $a_t^i$ , we can then re-organize the objective as:

$$\begin{split} \min_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \log \frac{\pi(a_{t}^{i}|g,h_{t},c_{t})}{\frac{1}{Z(g,h_{t},c_{t})}\pi_{LLM}(a_{t}^{i}|g,h_{t},c_{t})\exp\left(\alpha \mathcal{Q}_{LLM}(g,h_{t},a_{t}^{i},u_{t}^{i})\right)} - \log Z(g,h_{t},c_{t}) \right] \\ = \min_{\pi} \mathbb{E}_{a_{t}^{i} \sim \pi(a_{t}^{i}|g,h_{t},c_{t})} \left[ \log \frac{\pi(a_{t}^{i}|g,h_{t},c_{t})}{\pi^{*}(a_{t}^{i}|g,h_{t},c_{t})} - \log Z(g,h_{t},c_{t}) \right] \\ = \min_{\pi} \mathbb{D}_{KL} \left[ \pi(a_{t}^{i}|g,h_{t},c_{t}) \| \pi^{*}(a_{t}^{i}|g,h_{t},c_{t}) \right] - \log Z(g,h_{t},c_{t}). \end{split}$$

1129 Then since  $Z(g, h_t, c_t)$  does not depend on  $\pi$ , we can only care about the KL-divergence, which is 1130 minimized at 0 if and only if the two distributions are identical. Therefore, the optimal solution is

$$\pi(a_t^i|g, h_t, c_t) = \pi^*(a_t^i|g, h_t, c_t) = \frac{1}{Z(g, h_t, c_t)} \pi_{LLM}(a_t^i|g, h_t, c_t) \exp\left(\alpha \mathcal{Q}_{LLM}(g, h_t, a_t^i, u_t^i)\right),$$

$$(7)$$

which completes the derivation.



1154

Figure 15: Frameworks comparison. (a) *Actor-only* methods directly select the action with the highest probability generated by actor-LLM  $\pi_{LLM}$ , which may result in a lack of long-term planning and non-optimal action selection; (b) *Critic-only* methods self-evaluate each candidate action with another critic-LLM  $Q_{LLM}$  by first predicting candidate's future trajectory  $u_t^i$  and then directly select the action with the best-predicted outcome, which may ignore the prior knowledge in actor; (c) Our LLM-based Actor-Critic (LAC) algorithm integrate actor and dual critics: *lang-critic*  $C_{LLM}$  and *value-critic*  $Q_{LLM}$  to enhance the decision-making ability of LLMs.

# 1164 B.2 CRITIC IMPROVEMENT OF LAC

1165 The *lang-critic*  $C_{LLM}$ , *value-critic*  $Q_{LLM}$ , and forward model  $f_{LLM}$  we used can be easily imple-1166 mented by prompting LLMs via providing instructions or few-shot examples from similar tasks like 1167 prior work (Yao et al., 2022b; Liu et al., 2023). However, empirically, we found that they can be 1168 further improved via fine-tuning LLMs with simple next-token prediction loss on several samples 1169 collected from training tasks. In this work, we consider 18 trajectories for each benchmark for 1170 fine-tuning. Though 18 trajectories are significantly fewer than what is required for conventional 1171 reinforcement learning algorithms, they are generally enough for our method. Each trajectory has the following format:  $(g, o_0, a_1, o_1, c_1, \cdots, a_H, o_H, c_H)$ , where H is the episode length and  $c_i$  is a 1172 language-based evaluation of action  $a_t$ . Each  $c_t$  includes an explanation about the action  $a_t$  (e.g., "I 1173 have found object-X. This step is " or "I should take object-X instead of object-Y first. This step is ") 1174 and a special token that indicates positive/negative judgment (e.g., "GOOD" or "BAD"). 1175

1176 Practically, we just fine-tune the LLM once and use it to construct all the *lang-critic*  $C_{LLM}$ , value*critic*  $Q_{LLM}$ , and forward model  $f_{LLM}$ , thanks to the fine-tuning with the above data format. 1177 Specifically, when minimizing the loss of predicting future trajectories, the forward model  $f_{LLM}$ 1178 is improved. When minimizing the loss of generating language-based evaluations, the lang-critic 1179  $C_{LLM}$ , value-critic  $Q_{LLM}$  are both improved. The latter is because language-based evaluations 1180 also contain special tokens that indicate positive/negative judgments, whose generated probabilities 1181 are used to calculate  $Q_{LLM}$  in Equation (3). We analyze this fine-tuning process in Appendix A.9. 1182 Some examples of the labeled trajectories in ALFWorld and BabyAI-Text are shown in Table 14 and 1183 Table 15 respectively. 1184

1185

# 1186 B.3 Actor-only METHODS

1187

We compare our LAC with actor-only and critic-only methods in Figure 15.

1188 We detail the general *actor-only* methods in Algorithm 2. 1189 1190 Algorithm 2: Actor-only methods. 1191 **Input:** current task goal g, history  $h_t$ , actor  $\pi_{LLM}$ , candidate action size n. 1192 **Output:** selected action  $a_t^*$ 1193 1  $\{a_t^1, a_t^2, \cdots, a_t^n\} \sim \pi_{LLM}(\cdot | g, h_t)$ 1194  $a_t^* \leftarrow \arg \max_{a_t^i} \pi(a_t^i | g, h_t)$ 1195 1196 1197 **B**.4 Critic-only METHODS 1198 1199 We detail the general *critic-only* methods in Algorithm 3. Note that critic-only may use different  $Q_{Other}$  to estimate numerical assessment of actions. 1201 1202 Algorithm 3: Critic-only methods. 1203 **Input:** current task goal g, history  $h_t$ , actor  $\pi_{LLM}$ , forward model  $f_{LLM}$ , value-based critic 1204  $Q_{Other}$ , candidate action size n. 1205 **Output:** selected action  $a_t^*$ 1206 1  $\{a_t^1, a_t^2, \cdots, a_t^n\} \sim \pi_{LLM}(\cdot | g, h_t);$  $\triangleright$  generate candidate actions 1207 <sup>2</sup> for  $i \leftarrow 1, 2, \cdots, n$  do 1208  $u_t^i \leftarrow f_{LLM}(g, h_t, a_t^i);$ ▷ imagine future trajectory 3 1209 calculate  $\mathcal{Q}_{Other}(g, h_t, a_t^i, u_t^i)$ 4 1210 5 end 1211 6  $a_t^* \leftarrow \arg \max_{a_t^i} \mathcal{Q}_{Other}(g, h_t, a_t^i, u_t^i)$ 1212 1213 1214 С EXPERIMENT DETAILS 1215 1216 C.1 **BENCHMARK DETAILS** 1217 1218 C.1.1 ALFWORLD: BENCHMARK WITH HIGH-LEVEL ACTIONS 1219 We choose ALFWorld (Shridhar et al., 2021), a text-based household environment, to demonstrate 1220 the effectiveness of LAC on high-level planning. ALFWorld is a synthetic text-based game aligned with ALFRED (Shridhar et al., 2020) benchmark. There are 6 types of tasks in this environment, 1222 which require the agent to achieve a high-level goal through a sequence of high-level actions, e.g. 1223 "go to place-X", "take object-Y from place-X", etc. The details about the 6 task types in ALFWorld 1224 are shown in Table 4. 1225 1226 Table 4: All the task types and the corresponding goals for ALFWorld 1227 Туре Description 1228 Pick & Place The agent needs to put a target object to a target place, e.g. put some spraybottle on toilets, find some apple 1229 and put it in sidetable, etc. Clean & Place The agent needs to find a target object, clean it and put it to a target place, e.g. clean some apple and put it in 1230 sidetable, put a clean lettuce in diningtable, etc 1231 Heat & Place The agent needs to find a target object, heat it and put it to a target place, e.g. heat some egg and put it in 1232 diningtable, put a hot apple in fridge, etc. Cool & Place The agent needs to find a target object, cool it and put it to a target place, *e.g.* cool some pan and put it in 1233 stoveburner, put a cool mug in shelf, etc. Examine & Place The agent needs to find a target object, and examine it with desklamp, e.g. look at bowl under the desklamp, examine the pen with the desklamp, etc. Pick Two & Place The agent needs to put two target objects to a target place, e.g. put two saltshaker in drawer, find two pen and put them in dresser, etc.

1230

A challenge built into ALFWorld is that the agent needs to explore the environment to find a target object. The commonsense knowledge in LLMs about the likely locations for common household items makes this environment suitable for LLMs to solve. The reward is 1 only when the agent reaches the goal. Following ReAct, we evaluate 134 unseen evaluation games in a task-specific setup.

#### 1242 C.1.2 BABYAI-TEXT: BENCHMARK WITH LOW-LEVEL ACTIONS 1243

1244 For decision-making tasks with low-level planning, we adopt BabyAI-Text (Carta et al., 2023b) as our test-bed. BabyAI-Text is a text-only version environment extended from the BabyAI platform 1245 (Chevalier-Boisvert et al., 2019). BabyAI-Text is a Grid World environment, in which the agent and 1246 objects are placed in a room of  $8 \times 8$  tiles. The agent has 6 primitive actions: turn left, turn right, go 1247 forward, pick up, drop, toggle, to solve a task described in natural language (*e.g.* Pick up the red box). 1248 The agent has access to a  $7 \times 7$  partial view, which means it can only observe the objects belonging 1249 to the  $7 \times 7$  grid in front of it. In addition to objects relevant to completing a given task, there are 1250 also other distractors in the room. All the task types in BabyAI-Text are shown in Table 5. 1251

1252 1253

1263

Table 5: All the task types and the corres	sponding goals for BabyAI-Text

254	Туре	Description
255	go to	The agent needs to find target object and go to it, e.g. go to the green key, go to the red ball, etc.
056	pick up	The agent needs to find target object, go to it and pick up it, e.g. pick up the blue key, pick up the purple ball,
200		etc.
257	go to after pick up	The agent needs to find and pick up one object, then go to another object, <i>e.g.</i> go to the blue key after you
050		pick up the green key <i>etc</i> .
200	pick up then go to	The agent needs to find and pick up one object, then go to another object, e.g. pick up the green box, then go
259		to the purple box <i>etc</i> .
260	put next to	The agent needs to find and pick up one object, then go to another object and put the first object next to it,
200		<i>e.g.</i> put the grey key next to the yellow ball <i>etc.</i>
261	open door	The agent needs to know which key to pick up, then find and pick up it to open the door, <i>e.g.</i> open the door,
262		open the blue door, <i>etc</i> .
202		

Unlike ALFWorld, an agent interacting with BabyAI-Text needs to find out the suitable low-level 1264 action to execute at each step. We evaluate on the test environment in BabyAI-Text. The objects 1265 in a room are randomly chosen, and their position, as well as the agent's position, are also random. 1266 Considering the time and computational resource constraints, we evaluate on 50 tasks for each task 1267 type, yielding 300 tasks total.

1268 1269

#### C.2 IMPLEMENTATION DETAILS OF BASELINES 1270

1271 C.2.1 REACT 1272

1273 ReAct (Yao et al., 2022b) is an actor-only method that combines reasoning and acting during the 1274 interaction with the environment. It leverages the reasoning capabilities of LLMs to increase the probability of the LLM acting correctly as an actor. 1275

1276 Table 6 shows an example of a trajectory interleaving between actions and observations in ALFWorld. 1277 To leverage the reasoning ability of LLMs, ReAct adds an extra thinking action to the original action 1278 set in ALFWorld, allowing LLM agent to engage in reasoning about task decomposition and goal 1279 achievement during task execution. An annotated trajectory for ALFWorld given by ReAct is shown 1280 in Table 7. Text labeled in light green is the annotated thought to help LLM complete a task better.

Table 6. An even	nle of succes	oful trajectory	for ALEW/orld
-1 and $0$ . All $0$ all	IDIC OF SUCCES	SIULUAICCIOLV	

1283	Table 6: An example of successful trajectory for ALFWorld
1284	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
1285	handtowelholder 2, a handtowelholder 1, a sinkbasin 2, a sinkbasin 1, a toilet 1,
1286	a toiletpaperhanger 1, and a towelholder 1. Your task is to:  put some spraybottle on toilet.
1287	> go to cabinet 1
1288	On the cabinet 1, you see a cloth 1, a soapbar 1, a soapbottle 1. > go to cabinet 2
1289	The cabinet 2 is closed.
1290	> open cabinet 2 You open the cabinet 2. The cabinet 2 is open. In it, you see a candle 1, and a
1291	spraybottle 2.
1292	> take spraybottle 2 from cabinet 2 You pick up the spraybottle 2 from the cabinet 2.
1293	> go to toilet 1
1294	On the toilet 1, you see a soapbottle 2. > put spraybottle 2 in/on toilet 1
1295	You put the spraybottle 2 in/on the toilet 1.

1297	Table 7: An example of annotated successful trajectory for ALFWorld in ReAct
1298	You are in the middle of a room. Looking quickly around you, you see a cabinet
1299	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.
1300	a toiletpaperhanger 1, and a towelholder 1.
1301	Your task is to: put some spraybottle on toilet.
1302	toilet.
1303	OK. > think: First I need to find a spraybottle. A spraybottle is more likely to appear
1304	in cabinet (1-4), countertop (1), toilet (1), sinkbasin (1-2), garbagecan (1). I can
1305	check one by one, starting with cabinet 1. OK.
1306	> go to cabinet 1
1307	On the cabinet 1, you see a cloth 1, a soapbar 1, a soapbottle 1.
1000	> go to cabinet 2 The cabinet 2 is closed.
1308	> open cabinet 2
1309	You open the cabinet 2. The cabinet 2 is open. In it, you see a candle 1, and a
1310	spraybottle 2.
1311	<pre>&gt; think: Now I lind a spraybottle (2). Next, I need to take it. &gt; take spraybottle 2 from cabinet 2</pre>
1312	You pick up the spraybottle 2 from the cabinet 2.
1212	> think: Now I take a spraybottle (2). Next, I need to put it in/on toilet 1.
1010	UK.
1314	On the toilet 1, you see a soapbottle 2.
1315	> put spraybottle 2 in/on toilet 1
1316	You put the spraybottle 2 in/on the toilet 1.

1323 1324

1296

During the interaction with the environment, two successful trajectories like in Table 7 are provided to LLMs to leverage the in-context learning ability of LLMs. The LLMs will generate an action to be executed in the environment or conduct some reasoning about how to achieve the final goal.
The chosen action and new observation are appended to the prompt for LLMs to form a sequential decision-making process.

Table 8: An example of successful trajectory for BabyAI-Text

1225	Table 6. An example of successful trajectory for DabyAI-Text		
1323	Goal of the agent: go to the green key		
1326	Observation:You see a wall 2 steps right, You see a wall 3 steps forward, You see a		
1327	grey box 3 steps left, You see a grey key 3 steps left and 1 step forward, You see a blue box 3 steps left and 2 steps forward		
1328	Action:turn left		
1329	Observation:You see a wall 3 steps right, You see a blue key 3 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a green ball		
1330	1 step left and 3 steps forward, You see a grey box 3 steps forward, You see a grey		
1331	key 1 step right and 3 steps forward, You see a blue box 2 steps right and 3 steps forward		
1332	Action:go forward		
1333	Observation:You see a wall 3 steps right, You see a blue key 3 steps left and 1 step forward, You see a green key 2 steps left, You see a green ball 1 step left and 2		
1334	steps forward, You see a grey box 2 steps forward, You see a grey key 1 step right		
1335	and 2 steps forward, You see a blue box 2 steps right and 2 steps forward Action:turn left		
1336	Observation:You see a wall 3 steps left, You see a green box 2 steps left and 3 steps		
1337	forward, You see a green box 2 steps left and 1 step forward, You see a green key 2 steps forward, You see a blue key 1 step right and 3 steps forward, You see a green		
1338	ball 2 steps right and 1 step forward, You see a grey box 2 steps right		
1339	Action:go forward Observation:You see a green box 2 steps left and 2 steps forward, You see a green box		
1340	2 steps left, You see a green key 1 step forward, You see a blue key 1 step right and		
1341	2 steps forward, You see a green ball 2 steps right		

Similar to ALFWorld, we annotate some successful trajectories in BabyAI-Text to teach a ReAct agent how to complete a task in BabyAI-Text. Examples of successful trajectories without thoughts and with thoughts are shown in Table 8 and Table 9 respectively.

# 1346 C.2.2 RAP 1347

RAP(Hao et al., 2023) can be treated as a *critic-only* method, in which it utilizes LLMs as actor and
 world model, and adopts tree-search planning methods to evaluate each possible action candidate.
 Trajectories interleaving between actions and observations like in Table 6 and Table 8 are provided

1000	
1351	Table 9: An example of annotated successful trajectory for BabyAI-Text in ReAct
1352	Goal of the agent: go to the green key
1353	grey box 3 steps left, You see a grey key 3 steps left and 1 step forward, You see a
1354	blue box 3 steps left and 2 steps forward
1355	Action:think: To complete the task, I need to find a green key and go to it. Observation:OK
1356	Action:think: I have not seen a green key. Firstly, I need to find a green key.
1357	Observation:UK Action:turn left
1358	Observation:You see a wall 3 steps right, You see a blue key 3 steps left and 2 steps
1359	forward, You see a green key 2 steps left and 1 step forward, You see a green ball 1 step left and 3 steps forward, You see a grey box 3 steps forward, You see a grey
1360	key 1 step right and 3 steps forward, You see a blue box 2 steps right and 3 steps
1361	forward Action:think: I have seen a green key 2 steps left and 1 step forward. Next, I need
1362	to go to the green key.
1363	Observation:OK Action:go forward
1364	Observation:You see a wall 3 steps right, You see a blue key 3 steps left and 1 step
1365	forward, You see a green key 2 steps left, You see a green ball 1 step left and 2 steps forward, You see a grey box 2 steps forward, You see a grey key 1 step right
1366	and 2 steps forward, You see a blue box 2 steps right and 2 steps forward
1367	Action:turn left Observation:You see a wall 3 steps left. You see a green box 2 steps left and 3 steps
1368	forward, You see a green box 2 steps left and 1 step forward, You see a green key 2
1369	steps forward, You see a blue key 1 step right and 3 steps forward, You see a green ball 2 steps right and 1 step forward. You see a grev box 2 steps right
1370	Action:go forward
1371	Observation:You see a green box 2 steps left and 2 steps forward, You see a green box 2 steps left. You see a green key 1 step forward. You see a blue key 1 step right and
1372	2 steps forward, You see a green ball 2 steps right

1250

1374

1375

to LLMs to learn how to choose action as an actor and the dynamics of the environment as a world model.

1378 The assessment of each step is performed by a reward function, which can be the log probability 1379 of the action or self-evaluation given by LLMs, or based on some task-specific heuristics. In our 1380 implementation, we adopt the log probability of actions given by LLMs as the reward. For simplicity 1381 of implementation, we adopted a greedy approach to expand the tree, generating only one action at a time. More specifically, at each step, LLMs will sample some action candidates. For each action 1382 candidate, LLMs will generate a rollout trajectory until a maximum step or terminal state. The 1383 summation of log probabilities of all the actions on the rollout accessed by LLMs are used as Q value 1384 for each action candidate. The candidate with the highest Q value is chosen to be executed in the 1385 environment. 1386

1387

1388 1389

# C.2.3 ICPI

1390 1391

ICPI (Brooks et al., 2024) proposes to implement policy iteration using LLMs through in-context learning. At each step, the actor in ICPI will sample some action candidates and the critic will compute the Q values for each action candidates. The action candidates with the highest Q values is chosen to be executed.

The actor is implemented using LLMs, and successful trajectories like in Table 6 and Table 8 are provided to it.

As to the critic, ICPI prompts LLMs to give the numerical reward for each step directly. Given the current history and an action candidate, the critic in ICPI will rollout a trajectory starting from the action candidate. Apart from predicting the observations, the critic will give the numerical reward for each step on the rollout trajectory. The (discounted) return on the rollout is treated as the Q value for the action candidate. For both ALFWorld and BabyAI-Text, we define the reward as 1 when the agent reaches the goal. All other steps will have a reward 0. The examples provided to the critic are like in Table 10 for ALFWorld and Table 11 for BabyAI-Text.

_	Table 10: An example provided to critic in ICPI for ALFWorld
	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 tolletpaperhanger 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. Reward:0
	> 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.
_	Reward:1
	Table 11: An example provided to critic in ICPI for BabyAI-Text
	Goal of the agent: go to the green key
	Observation: You see a wall 2 steps right, You see a wall 3 steps forward, You see a grey box 3 steps left. You see a grey key 3 steps left and 1 step forward. You see
	blue box 3 steps left and 2 steps forward
	Action:turn left Observation:You see a wall 3 steps right. You see a blue key 3 steps left and 2 st
	forward, You see a green key 2 steps left and 1 step forward, You see a green ball
	1 step left and 3 steps forward, You see a grey box 3 steps forward, You see a gree key 1 step right and 3 steps forward. You see a blue box 2 steps right and 3 steps
	forward
	Reward:0
	Action:go lorward Observation:You see a wall 3 steps right, You see a blue key 3 steps left and 1 st
	forward, You see a green key 2 steps left, You see a green ball 1 step left and 2
	steps forward, You see a grey box 2 steps forward, You see a grey key 1 step right and 2 steps forward. You see a blue box 2 steps right and 2 steps forward
	Reward:0
	Action:turn left
	forward, You see a green box 2 steps left and 1 step forward, You see a green key
	steps forward, You see a blue key 1 step right and 3 steps forward, You see a gree
	<pre>pail 2 steps right and 1 step forward, You see a grey box 2 steps right Reward:0</pre>
	Action:go forward
	Observation: You see a green box 2 steps left and 2 steps forward, You see a green 2 steps left. You see a green key 1 step forward. You see a blue key 1 step right
	2 steps forward, You see a green ball 2 steps right
	Reward:1

The framework of RAFA (Liu et al., 2023) is also like RAP or ICPI. The main difference is the critic used.

RAFA implements tree-search using LLM to evaluate each action candidate. Different from ICPI,
RAFA uses the task completion progress as the value for each step. They have the LLMs decompose
a goal into sub-goals, and use the completion status of the sub-goals after each step as the value for
the step. RAFA evaluates the completion status of sub-goals based on the predicted observations.
Examples provided to critic in RAFA are like in Table 12 for ALFWorld and Table 13 for BabyAI-Text.

1458	Table 12: An example provided to critic in PAEA for ALEWorld
1439	You are in the middle of a room. Looking guickly around you, you see a cabinet
1460	4, a cabinet 3, a cabinet 2, a cabinet 1, a countertop 1, a garbagecan 1, a
1461	handtowelholder 2, a handtowelholder 1, a sinkbasin 2, a sinkbasin 1, a toilet 1,
1462	a toiletpaperhanger 1, and a towelholder 1. Your task is to: put some spraybottle on toilet
1463	> critic: My task requires two sub-goals in order: take a spraybottle and put the
1464	
1465	> 0K.
1466	On the cabinet 1, you see a cloth 1, a soapbar 1, a soapbottle 1. > OK.
1467	The cabinet 2 is closed.
1468	> OK.
1469	Sou open the cabinet 2. The cabinet 2 is open. In it, you see a candle 1, and a spraybottle 2.
1470	> UK. You pick up the spravhottle 2 from the cabinet 2.
1471	> critic: Now I take a spraybottle. My current state satisfies the first of the two sub-goals: take a spraybottle. The value is $1/2-0.5$
1472	> OK.
1473	On the toilet 1, you see a soapbottle 2.
1474	> OK.
1475	<pre>&gt; critic: Now I put the spraybottle on the toilet. My current state satisfies all the two sub-goals. The value is 2/2-1</pre>
1476	the two sub-goals. The value is 2/2-1.
1477	
1478	
1479	Table 13: An example provided to critic in RAFA for BabyAI-Text
1/00	Goal of the agent: go to the green key You see a wall 2 steps right. You see a wall 3 steps forward. You see a grey box 3
1481	steps left, You see a grey key 3 steps left and 1 step forward, You see a blue box 3 steps left and 2 steps forward
1482	>critic: My task requires two sub-goals in order: find the green key, and go to
1483	
1484	>0K.
1485	You see a wall 3 steps right, You see a blue key 3 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a green ball 1 step left
1486	and 3 steps forward, You see a grey box 3 steps forward, You see a grey key 1 step right and 3 steps forward, You see a blue box 2 steps right and 3 steps forward
1487	
1489	>OK.
1490	see a green key 2 steps left, You see a green ball 1 step left and 2 steps forward,
1491	You see a grey box 2 steps forward, You see a grey key 1 step right and 2 steps
1492	-orward, fou see a brue box 2 steps right and 2 steps forward >OK.
1493	You see a wall 3 steps left, You see a green box 2 steps left and 3 steps forward,
1494	forward, You see a blue key 1 step right and 3 steps forward, You see a green key 2 steps
1495	2 steps right and 1 step forward, You see a grey box 2 steps right
1496	You see a green box 2 steps left and 2 steps forward, You see a green box 2 steps
1497	left, You see a green key 1 step forward, You see a blue key 1 step right and 2 steps forward. You see a green ball 2 steps right
1498	>critic: Now I go to the green key. My current state satisfies all the two
1499	sub-goals. The value is 2/2=1.
1500	

# C.3 IMPLEMENTATION DETAILS OF LAC, ABLATION STUDIES, AND VARIANTS

We build our method based on REACT (Yao et al., 2022b). Specifically, we adopt the trajectories 1503 with thoughts from ReAct and annotate language-based evaluations for those trajectories. For each 1504 task type, we annotate two trajectories from ReAct as successful examples to teach LLMs how 1505 to complete this type of tasks. Some examples of the annotated trajectories for ALFWorld and 1506 BabyAI-Text are shown in Table 14 and Table 15 respectively. The language-based evaluations we 1507 annotated are labeled in light orange. We also use ReAct algorithm to collect some failed trajectories 1508 from the training set for each task type. The collected failed trajectories are also annotated with 1509 language-based evaluations, aiming to teach LLMs to distinguish between steps that are beneficial, 1510 harmful, or indeterminate for achieving a goal. 1511

We annotate the language-based evaluations using the following criteria:

1512	
1513	Table 14: An example of annotated successful trajectory for ALFWorld in LAC
1514	You are in the middle of a room. Looking quickly around you, you see a cabinet
1515	handtowelholder 2, a handtowelholder 1, a sinkbasin 2, a sinkbasin 1, a toilet 1,
1516	a toiletpaperhanger 1, and a towelholder 1.
1517	> think: To solve the task, I need to find and take a sparybottle, then put it on
1518	toilet.
1519	OK. ==> critic: The thought is reasonable. This step is GOOD.
1520	> think: First I need to find a spraybottle. A spraybottle is more likely to appear
1521	check one by one, starting with cabinet 1.
1522	OK.
1523	
1524	> go to cabinet 1
1525	==> critic: I have not seen spraybottle in this place. This step is UNKNOWN.
1526	> go to cabinet 2
1527	==> critic: The cabinet 2 is closed. This step is UNKNOWN.
1528	> open cabinet 2
1529	You open the capinet 2. The capinet 2 is open. In it, you see a candle 1, and a spraybottle 2.
1530	==> critic: I have seen spraybottle 2 in this place. This step is GOOD.
1531	> think: Now I find a spraybottle (2). Next, I need to take it. OK.
1532	
1533	This step is GOOD. > take spraybottle 2 from cabinet 2
1534	You pick up the spraybottle 2 from the cabinet 2.
1535	==> critic: I have taken spraybottle 2 in this place. This step is GOOD. > think: Now I take a spraybottle (2). Next, I need to put it in/on toilet 1.
1536	OK.
1537	
1538	> go to toilet 1
1539	On the toilet 1, you see a soapbottle 2. ==> critic: I should go to toilet if I want to put the spraybottle on it. This step
1540	is GOOD.
1541	> put spraybottle 2 in/on toilet 1 You put the spraybottle 2 in/on the toilet 1.
1542	==> critic: I have put the spraybottle 2 in/on the toilet 1. This step is GOOD.
1543	
1544	(1) If an action on a trajectory is necessary for achieving the final goal, the following language-
1545	hased evaluations will end with "This step is GOOD." to indicate the completion of an
1546	essential step, such as action "take spraybottle 2 from cabinet 2" for the goal "put some
1547	spraybottle on toilet".
1548	(2) If an action is crucial for obtaining essential information to complete a goal, the subsequent
1549	language-based evaluations will end with "This step is GOOD.", such as action "open
1550	cabinet 2" in Table 14 which finds the target object.
1551	(3) If an action is definitely useless or incorrect for achieving a goal, the following language-
1552	hased evaluations will end with "This step is BAD" to indicate this action is undesirable
1553	such as action "take cloth 1 from cabinet 1" for the task in Table 14.
1554	(4) If an action can not be evaluated as good or had from the history, the following language
1555	(4) If all action can not be evaluated as good of bad from the instory, the following language- hased evaluations will end with "This step is UNKNOWN" such as action "go to cabinet 1"
1556	or "go to cabinet 2" in Table 14
1557	
1558	All the annotated successful and failed trajectories are used to construct fine-tuning examples to
1559	fine-tune LLMs to generate better language-based evaluations.
1560	Considering the computational cost of fully finetuning of LLMs, we use $L_0RA$ (Hu et al. 2021) to
1561	finetune our models. In ALFWorld, with two successful trajectories and one failed trajectory for
1562	each task type, we have 485 (input,output) pairs in total of six task types to finetune models. In
1563	BabyAI-Text, the number of finetuning (input,output) pairs is 418. We finetune models for 1.000
1564	steps with learning rate 2.5e-5 and batch size 2. We use A100 GPU with 80GB memory to fine-tune
1565	our model. With just about 400-500 (input,output) pairs and 1,000 fine-tuning steps, we can complete
	the fine-tuning process within one and a half hours.

1566	
1567	Table 15: An example of annotated successful trajectory for BabyAI-Text in LAC
1568	Goal of the agent: go to the green key
1569	critic box 3 steps left, You see a critic key 3 steps left and 1 step forward, You
1570	see a blue box 3 steps left and 2 steps forward
1571	Action:think: To complete the task, I need to find a green key and go to it. Observation:OK
1572	Critic:The thought is reasonable. This step is GOOD.
1573	Action:think: I have not seen a green key. Firstly, I need to find a green key. Observation:OK
1574	
1575	GOOD. Action:turn left
1576	Observation:You see a wall 3 steps right, You see a blue key 3 steps left and 2 steps
1577	forward, You see a green key 2 steps left and 1 step forward, You see a green ball 1 step left and 3 steps forward. You see a critic box 3 steps forward. You see a critic
1578	key 1 step right and 3 steps forward, You see a blue box 2 steps right and 3 steps
1579	forward Critic:I have turned left and seen a green key 2 steps left and 1 step forward. This
1580	
1581	Action:think: I have seen a green key 2 steps left and 1 step forward. Next, I need
1582	Observation:OK
1583	
1584	Action:go forward
1585	Observation:You see a wall 3 steps right, You see a blue key 3 steps left and 1 step
1596	forward, You see a green key 2 steps left, You see a green ball I step left and 2 steps forward, You see a critic box 2 steps forward, You see a critic key 1 step
1500	right and 2 steps forward, You see a blue box 2 steps right and 2 steps forward
1507	Critic:1 have gone forward. The green key is 2 steps left. This step is UNKNOWN. Action:turn left
1500	Observation:You see a wall 3 steps left, You see a green box 2 steps left and 3 steps
1509	forward, You see a green box 2 steps left and 1 step forward, You see a green key 2 steps forward, You see a blue key 1 step right and 3 steps forward, You see a green
1590	ball 2 steps right and 1 step forward, You see a critic box 2 steps right
1591	Critic:I have turned left. The green key is 2 steps forward. This step is UNKNOWN. Action:go forward
1592	Observation:You see a green box 2 steps left and 2 steps forward, You see a green box
1593	2 steps left, You see a green key 1 step forward, You see a blue key 1 step right and 2 steps forward. You see a green ball 2 steps right
1594	Critic:I have gone forward. The green key is 1 step forward. This step is GOOD.
1595	
1590	
1597	
1590	
1599	During testing, the fine-tuned models are used to generate language-based evaluations after executing
1600	an action in the environment, as well as to forecast the potential outcomes of each action candidate.
1601	We sample $n = 5$ action candidates from $\pi_{TTM}$ at each time step. The $\pi_{TTM}$ is augmented by
1602	leveraging the language-based evaluations generated by $C_{IIM}$ . After sampling action candidates.
1603	we use the fine-tuned model to predict future outcomes for each action candidate. The model needs
1604	to predict the possible observation and generate language-based evaluations for each predicted step.
1605	We set the maximum prediction step as 4, the model will continue the prediction until it generates
1606	a language-based evaluation ending with "This step is GOOD." or "This step is BAD", or when it
1607	reaches the maximum prediction step.
1608	For the optimization of $\pi_{xxxx}$ , we solve an optimization problem in Equation (4) with a hyper
1609	parameter $\alpha$ which balances the generating probabilities of $\pi_{LLM}$ and the values given by $Q_{LLM}$
1610	For ALFWorld, we set $\alpha$ as 1, which yields superb performance over baselines. For BabyAI-Text,
1611	we conduct a grid-search over $\{1/2, 1, 2, 5, 10\}$ for $\alpha$ , finding that different LLMs will have best
1612	performance with different $\alpha$ . The results can be seen in Table 3.
1613	We get the movimum having a langth to 40 few AI EW-ald and 20 (a) Data AI That I (the second se
1614	we set the maximum norizon length to 40 for ALF world and 50 to BabyAl-lext. If the agent has not
1615	reached the final goal after 40 or 50 steps, this episode will be marked as failure.
1616	We use A100 GPU with 80GB memory to evaluate our method. For LAC, the execution time for
1617	ALFWorld is about 10 hours for 134 tasks using single A100 GPU. And for BabyAI-Text, the
1618	execution time can be varied for different task types, ranging from 4 to 10 hours for 50 tasks using
1619	one A100 GPU. The GPU memory usage may range from 15GB to over 70GB during the interaction
	according to the length of inputs to LLMs.

1620 We compare our method with all the aforementioned baselines, demonstrating the effectiveness of our 1621 method on decision-making tasks with both high-level actions and low-level actions. To demonstrate 1622 the effectiveness of each component in our method, we conduct ablation studies on each component. 1623 We removes the lang-critic  $C_{LLM}$  from LAC as well as the integration during pre-action-generation 1624 phase. This variant is called LAC w/o lang-critic. We also evaluate the role of  $Q_{LLM}$  by removing it from LAC as well as the integration during post-action-generation phase. This variant is called 1625 LAC w/o value-critic. We also demonstrate the role of the action prior given by LLM policy by using 1626 only value-critic  $Q_{LLM}$  for decision-making. We call this variant as Value-critic-only. The execution 1627 time of those variants during evaluation can be varied according to its performance because a method 1628 having poor performance typically will cost more time to execute. On ALFWorld, it may be 10-20 1629 hours. The comparisons between those variants are shown in Figure 4. 1630

We found that each component in LAC is crucial for the superb performance. Removing some components may lead to wrong choice of action candidates. Such an example is shown in Table 16.
LAC can complete this task successfully, while eliminating some components in LAC will lead to failure. The comparison is shown in Figure 1.

1635

1637

# 6 D LIMITATIONS

1638 Our work has two limitations. Firstly, *lang-critic* of LAC is only used before action generation and it 1639 can also be applied after action generation. For example, it can provide language-based evaluations 1640 for predicted future trajectories to re-sample candidate actions if the previous candidate actions all 1641 fail to complete the target task. Secondly, though the *value-critic* of LAC can also adopt tree-search 1642 to provide a more accurate assessment of candidate actions, in this paper, we only expand one node 1643 for each candidate action for simplicity.

1644

# 1645 E BROADER IMPACTS

1646

Our method is built upon open-source large language models (LLMs). Like other methods that use LLMs, our method also inherits some benefits and challenges from LLMs. For the benefits, our method directly exploits the prior knowledge from LLMs, which may reduce potential carbon costs compared with training policies from scratch. For the challenges, our method might be susceptible to producing unintended output when confronted with harmful input, such as unethical text or input intended for adversarial attacks. To solve this problem, we suggest a thoughtful deployment of our method, such as adding a filtering component.

- 1654
- 655
- 1656

1657

- 1658 1659
- 1660
- 1661
- 1662
- 1663 1664
- 166
- 166
- 1667
- 1669
- 1670
- 1671
- 1672
- 1673

<pre>Table 16: An example of successful trajectory in test environment in BabyAl-Text given by L2 Goal of the agentigo to the green ball Observation:You see a wall 5 steps forward, You see a blue ball is top left and 4 steps forward, You see a green key 2 steps right and 2 steps forward, You see a green key 2 steps right and 2 steps forward, You see a green key 2 steps right and 2 steps forward, You see a green key 2 steps right and 2 steps forward, You see a green key 2 steps right and 2 steps forward, You see a green key 2 steps right and 2 steps forward, You see a green key 2 steps right and 2 steps forward, You see a green key 2 steps right and 2 steps forward, You see a green key 2 steps right and 2 steps forward, You see a green key 1 steps right and 2 steps forward, You see a green key 2 steps right and 2 steps forward, You see a green key 1 step left and 3 steps forward, You see a green key 2 steps left and 2 steps forward, You see a green key 1 step left and 2 steps forward, You see a green key 1 step left and 3 steps forward, You see a green key 1 step left and 2 steps forward, You see a green key 1 step left and 2 steps forward, You see a green key 1 step left and 2 steps forward, You see a green key 1 step left and 3 steps forward, You see a green key 1 step left and 2 steps forward, You see a green key 1 step left and 3 steps forward, You see a green key 1 step left and 3 steps forward, You see a green key 1 step left and 3 steps forward, You see a green key 1 step left and 3 steps left and 3 steps forward, You see a green key 1 step left and 2 steps forward, You see a green key 1 step left and 3 steps forward, You see a green key 1 step left and 3 steps forward, You see a green key 1 step left and 3 steps forward, You see a green key 1 step left and 3 steps forward, You see a green key 1 step left. The green ball is 3 steps left and 3 steps forward, You see a green key 1 step left. Thus g</pre>
<pre>Table 16: An example of successful trajectory in test environment in BabyAl-Text given by L2 Goal of the agentigo to the green ball Observation: You see a wall 5 steps forward, You see a wall 2 steps left, You see a yellow ball 1 step left and 4 steps forward, You see a wall 2 steps right and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a blue key 1 steps if the task, i need to find a green ball and go to it. Observation: You see a read to task for a steps forward. Not see a yellow ball is reasonable. This step is COOL. Action: think: I have seen a green ball is top left and 1 step forward. Next, I need to go to the green ball. Observation: You see a wall 5 steps forward, You see a wall 2 steps right, You see a green ball 3 step ielf and 4 steps forward, You see a gree key 2 steps i forward, You see a wall 5 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 4 steps forward, You see a gree key 2 steps i forward, You see a blue key 1 step left and 3 steps forward, You see a green ball 3 steps iforward, You see a gree key 2 steps i forward, You see a blue key 1 step left and 3 steps forward, You see a gree key 2 steps i forward, You see a blue key 1 step left and 3 steps forward, You see a gree key 2 steps i forward. This step is forward, You see a steps key 1 step left and 3 steps forward, You see a wall 4 steps forward, You see a gree key 1 step left and 3 step forward. This step is forward. This step is forward. This step is forward. This steps is forward. This step is</pre>
<pre>Table 16: An example of successful trajectory in test environment in BabyAI-Text given by 12 Cost of the agentigo to the green ball a steps forward, You see a red ball is teps left and 1 step forward, You see a green key 2 steps right and 2 steps forward, You see a green key 2 steps right and 2 steps forward, You see a green key 2 steps right and 1 step forward, You see a red ball is tep left and 1 step forward. Mation:thinkin awa been a green table step is COOO. Adjoint the word is reasonable. This step is COOO. Adjoint the word is reasonable. I can go to the green ball after seeing it. This step is COOD. Adjoint the agent is reasonable. I can go to the green ball after seeing it. This step is COOD. Adjoint the word is reasonable. I can go to the green ball after seeing it. This step is COOD. Adjoint the agent agent and 1 steps forward, You see a green key 2 steps left and 2 steps forward, You see a wall 5 steps forward, You see a green ball a steps forward. Mouse a green ball 3 steps left and 4 steps forward, You see a green key 1 step left and 2 steps forward, You see a green ball 3 steps forward. Mouse a green ball 3 steps forward, You see a green ball 3 steps left and 3 steps forward. Mouse a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward. Mouse a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward. Mouse a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward. Mouse a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward. Mouse a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward. Mouse a wall 4 steps forward, You see a green ball 3 steps left and 3 step forward. Mouse a wall 4 steps forward, You see a green ball 3 steps left. Mouse a step is forward. Mouse a bale ball all step right. The green ball is 3 steps left and 1 step forward. Mouse a superise is MONON. Acti</pre>
<pre>Table 16: An example of successful trajectory in test environment in BabyAl-Text given by Li Goal of the agent:go to the green ball Observation:You see a wall 5 steps forward, You see a wall 2 steps left, You see a yellow ball 1 step left and 4 steps forward, You see a blue ball step left and 3 steps forward, You see a red ball 1 step left and 1 step forward, You see a grey key 2 steps right and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 3 steps right and 1 step forward, You see a green ball step left and 4 steps forward. Next, I need to go to the green ball. Observation:YO CriticiThis thought is reasonable. I can go to the green ball after seeing it. This step is GOOL Action:think:To complete the tasts, I need to 1 find a green ball after seeing it. This step is GOOL Action:think:To complete the tasts is the step left and 3 steps forward, You see a green ball after seeing it. This step is GOOL Action:think:To complete the tasts is the step left and 3 steps forward, You see a green ball after seeing it. This step is GOOL Action:think:To was a wall 5 steps forward, You see a green ball after seeing it. This step is GOOL Action:town right Observation:You see a blue key 1 step left and 3 steps forward, You see a green ball 3 steps forward, You see a green ball 1 step right and 1 step forward CriticiT have turned right. The green ball is 3 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a wall 4 steps forward, You see a green ball 3 steps left, You see a green ball 1 step right Critic:T have gone forward. The green ball is 3 steps left and 1 step forward. You see a blue ball Steps left and 2 steps forward, You see a green key 1 step left and 3 steps forward Observation:You see a wall 4 steps forward, You see a green key 1 step left Critic:T have gone forward. The green ball is 3 steps left and 1 step forward. This step is UWKNOM Action; yo forward Ob</pre>
Table 16: An example of successful trajectory in test environment in BabyAI-Text given by L. Goal of the agent:go to the green ball observation:You see a wall 5 steps forward, You see a wall 2 steps left, You see a yellow ball 1 step left and 4 steps forward, You see a blue ball 1 step left and 1 step forward, You see a red ball 1 step left and 1 step forward, You see a red ball 1 step left and 1 step forward, You see a red ball 1 step left and 1 step forward, You see a red ball 1 step left and 1 step forward, You see a red ball 1 step left and 1 step forward, You see a red ball 1 step left and 1 step forward, You see a red ball 1 step left and 1 step forward, You see a preen ball and go to it. Observation:0% Critic:Inis thought is reasonable. This step is GOOD. Action:think:I have seen a green ball is the left and 6 steps forward. Next, I need to go to the green ball. Observation:0% Critic:This thought is reasonable. I can go to the green ball after seeing it. This step is GOOD. Action:turn right Observation:You see a wall 5 steps forward, You see a wall 2 steps right, You see a green hey 1 step left and 2 steps forward, You see a blue ball 1 step right and 1 step forward. Critic: have turned right. The green ball is 3 steps left and 4 steps forward. This step is OKXONN. Action:go forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a grey key 2 steps left and 3 steps forward. Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward. You see a blue ball 1 step right Critic: have gone forward. The green ball is 3 steps left and 3 steps forward. You see a blue ball 1 step right Critic: have gone forward. The green ball is 3 steps left and 3 steps forward. You see a wall 3 steps left and 2 steps forward, You see a green key 1 step left. You see a green ball 3 steps left and 2 steps forward. Motion:Y
Table 16: An example of successful trajectory in test environment in BabyAI-Text given by Li Goal of the agent:go to the green ball Observation:You see a weall 5 steps forward, You see a wall 2 steps left, You see a yellow ball 1 step left and 4 steps forward, You see a blue ball 1 step forward, You see a gree key 2 steps right and 2 steps forward, You see a green key 2 steps right and 1 step forward, You see a to blue key 3 steps right and 1 step forward. You see a blue key 3 steps right and 1 step forward. You see a blue key 3 steps right and 1 step forward. Not the green ball. Observation:00 Critic:This thought is reasonable. This step is GOOD. Action:think:I have seen a green ball 1 step left and 4 steps forward. Next, I need to go to the green ball. Observation:00 Critic:This thought is reasonable. I can go to the green ball after seeing it. This step is GOOD. Action:turn right Observation:You see a blue key 1 step left and 3 steps forward, You see a green key 1 step left and 2 steps forward, You see a green key 2 steps left and 2 steps forward, You see a blue key 1 step left and 3 steps forward, You see a green key 1 step left and 2 steps forward, You see a green ball 3 steps forward. This step is forward, You see a green key 1 step left and 4 steps forward. This step is forward, You see a green key 1 step left and 4 steps forward. This step is forward, You see a green key 1 step left and 3 steps forward. This step is UKNONN. Action:go forward Observation:You see a wall 3 steps forward, You see a green key 1 step left and 3 steps forward. This step is UKNONN. Action:go forward Observation:You see a wall 3 steps forward, You see a green key 1 step left. This step is UKNONN. Action:go forward Observation:You see a wall 3 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 2 steps forward, You see a green key 1 step left. This step is UKNONN. Action:go forward Observation:You see a wall 3 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step
<pre>Goal of the agentige to the green ball Observation:You see a wall 5 steps forward, You see a wall 2 steps left, You see a yellow ball 1 step left and 4 steps forward, You see a blue ball 1 step left and 3 steps forward, You see a red ball 1 step left and 1 step forward Action:think:In complete the task, I need to find a green ball and go to ft. Observation:OK CriticiThis thought is reasonable. This step is GOOD. Action:think:In complete the task, I need to find a green ball and go to ft. Observation:OK CriticiThis thought is reasonable. This step is GOOD. Action:think:In complete the task, I need to find a green ball after seeing it. This step is GOOD. Action:think:In complete the task of need to find a green ball after seeing it. This step is GOOD. Action:think:I have seen a green ball step left and 4 steps forward. Next, I need to go to the green ball. Observation:You see a wall 5 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 4 steps forward, You see a greek y 2 steps left and 2 steps forward, You see a blue key 1 step left and 3 steps forward, You see a green key 1 step left and 2 steps forward, You see a blue ball 1 step right and 1 step forward OrificiThave turned right. The green ball is 3 steps left and 4 steps forward. This atep is DUKNOM. Action:go forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward. Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward Observation:You see a wall 3 steps forward, You see a green ball 3 steps left. You see a green ball 3 steps left and 2 steps forward, You see a green key 1 step left and 3 steps forward. This step is UUKNOM. Action:go forward Observation:You see a wall 3 steps forward, You see a awall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps forward. This step is UUKNOM. Action:go forward Observation</pre>
<pre>Deservation:You see a wall 5 steps forward, You see a wall 2 steps left, You see a yellow hall 1 step left and 4 steps forward, You see a fue ball 1 step left and 3 steps forward, You see a red ball 1 step left and 1 step forward, You see a blue key 3 steps right and 1 step forward Action:think:Ho complete the task, I need to find a green ball and go to it. Observation:OK Critic:This thought is reasonable. This step is GOOD. Action:think:Hought is reasonable. I can go to the green ball after seeing it. This step is GOOD. Action:think:Hought is reasonable. I can go to the green ball after seeing it. This step is GOOD. Action:turn right Observation:OX usee a wall 5 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 4 steps forward, You see a wall 2 steps left and 2 steps forward, You see a vall 5 steps forward, You see a green key left and 1 step forward Observation:You see a wall 5 steps forward, You see a green key left and 1 step forward Critic:Thave turned right. The green ball is 3 steps left and 4 steps forward. This step is OURMOWM. Action:go forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a yeek by 2 steps left and 1 step forward. This step is OURMOWM. Action:go forward Observation:You see a grey key 2 steps left and 1 step forward. You see a blue ball 1 step right Critic: I have gone forward. The green ball is 3 steps left and 3 steps forward. This step is OURMOWM. Action:go forward Observation:You see a wall 3 steps forward, You see a green key 1 step left Critic: I have gone forward. The green ball is 3 steps left and 3 steps forward. This step is OURMOWM. Action:go forward Observation:You see a wall 2 steps forward, You see a green key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is OURMOWM. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You</pre>
<pre>Notward, Hou see a blue key 3 steps Fight and 1 step forward Action:think:16 complete the task, I need to find a green ball and go to it. Observation:0K Critic:This thought is reasonable. This step is GOOD. Action:think:I have seen a green ball 1 step left and 4 steps forward. Next, I need to go to the green ball. Observation:0K Critic:This thought is reasonable. I can go to the green ball after seeing it. This step is GOOD. Action:turn right Observation:You see a wall 5 steps forward, You see a greey key 2 steps left and 2 steps forward, You see a blue key 1 step left and 3 steps forward, You see a green key 1 step left and 2 steps forward, You see a green ball 1 step right and 1 step forward Critic:I have turned right. The green ball is 3 steps left and 4 steps forward. This step is UNKNOWN. Action:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a targe key 2 steps left and 1 step forward. This step is UNKNOWN. Action:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward. You see a green ball is 3 steps left and 1 step forward. You see a blue ball 1 step right Critic:I have gone forward. You see a green ball 3 steps right, You see a blue key 1 step left and 2 steps forward, You see a green ball 3 steps right. You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNONN. Action:go forward Observation:You see a wall 2 steps forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNONN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You</pre>
<pre>Observation:OK CriticiThis thought is reasonable. This step is GOOD. Action:think:I have seen a green ball 1 step left and 4 steps forward. Next, I need to go to the green ball. Observation:OK Critic:This thought is reasonable. I can go to the green ball after seeing it. This step is GOOD. Action:turn right Observation:You see a wall 5 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 4 steps forward, You see a grey key 2 steps left and 2 steps forward, You see a blue key 1 step left and 3 steps left and 1 step forward Critic: Have turned right. The green ball is 3 steps left and 4 steps forward. This step is UNKNONN. Action:go forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward. Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward Observation:You see a wall 3 steps forward, You see a green key 1 step left and 2 steps forward. This step is UNKNONN. Action:go forward Observation:You see a wall 3 steps forward, You see a green key 1 step left and 2 steps forward. This step is UNKNONN. Action:go forward Observation:You see a wall 3 steps forward, You see a green key 1 step left Critic: I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNONN. Action:go forward Observation:You see a wall 2 steps forward, You see a green key 1 step left Critic: I have gone forward. The green ball is 3 steps left and 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left. This step is UNKNONN. Action::go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps</pre>
<ul> <li>Action:thinkil have seen ball: Step left and 4 steps forward. Next, I need to go to the green ball.</li> <li>Observation:OK</li> <li>Critic:This thought is reasonable. I can go to the green ball after seeing it. This step is GOOD.</li> <li>Action:turn right</li> <li>Observation:You see a wall 5 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 4 steps forward, You see a green key 1 step left and 3 steps forward, You see a green key 1 step left and 2 steps forward.</li> <li>Critic:Thiswithwown.</li> <li>Action:go forward</li> <li>Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 2 steps forward.</li> <li>Action:go forward</li> <li>Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward.</li> <li>Action:go forward</li> <li>Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a grey key 2 steps left and 1 step forward.</li> <li>You see a blue ball 1 step right</li> <li>Critic:Thave gone forward. The green ball is 3 steps left and 3 steps forward.</li> <li>Action:go forward</li> <li>Observation:You see a wall 3 steps forward, You see a green key 1 step left and 2 steps forward.</li> <li>Action:go forward</li> <li>Observation:You see a wall 3 steps forward, You see a green key 1 step left.</li> <li>Critic:Thave gone forward. The green ball is 3 steps left and 3 steps forward.</li> <li>Action:go forward</li> <li>Observation:You see a wall 2 steps forward, You see a green key 1 step left</li> <li>Critic:Thave gone forward. The green ball is 3 steps left and 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a sull 2 steps right. You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left and 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left and 2 steps forward, You see a blue key 1 step left and 1 step forward, You see a</li></ul>
<pre>to go to the green ball. Observation:OK Critic:This thought is reasonable. I can go to the green ball after seeing it. This step is GOOD. Action:Turn right Observation:You see a wall 5 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 4 steps forward, You see a greey key 2 steps left and 2 steps forward, You see a blue key 1 step left and 3 steps forward, You see a green key 1 step left and 2 steps forward, You see a blue ball 1 step right and 1 step forward Critic:I have turned right. The green ball is 3 steps left and 4 steps forward. Mation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a grey key 2 steps left and 1 step forward, You see a blue key 1 step left and 2 steps forward, You see a green key 1 step left and 1 step forward, You see a blue ball 1 step right Critic:I have gone forward. You see a green key 1 step left and 1 step forward, You see a blue ball 1 step right Critic:I have gone forward. The green ball is 3 steps left and 3 steps right, You see a green ball 3 steps left and 2 steps forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps right, You see a blue key 1 step left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNNNOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left. This step is UNNNOW. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left. This step is UNNNOW. Action:go forward Observation:You see a wall 1 step forward, You see a green key 1 step left. This step is UNNNOWN. Action:Unn left Observation:You see a wall 1 step right, You</pre>
<pre>Doservation:0k CriticiThis thought is reasonable. I can go to the green ball after seeing it. This step is GOOD. Action:turn right Observation:You see a wall 5 steps forward, You see a greey key 2 steps left and 2 steps forward, You see a blue key 1 step left and 3 steps forward, You see a green key 1 step left and 2 steps forward, You see a blue ball 1 step right and 1 step forward Critic:I have turned right. The green ball is 3 steps left and 4 steps forward. This step is UNKNONN. Action:go forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a grey key 2 steps left and 1 step forward, You see a blue key 1 step left and 2 steps forward, You see a green key 1 step left and 1 step forward, You see a blue ball 1 step right Critic:I have gone forward. The green ball is 3 steps left and 3 steps forward. This step is UNKNONN. Action:go forward Observation:You see a wall 3 steps forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps right, You see a blue key 1 step left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNONN. Action:Go forward Observation:You see a wall 1 step right, You see a green key</pre>
<pre>step is GOOD. Action:turn right Observation:You see a wall 5 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 4 steps forward, You see a grey key 2 steps left and 2 steps forward, You see a blue key 1 step left and 3 steps forward. You see a green key 1 step left and 2 steps forward, You see a blue ball 1 step right and 1 step forward Critic:I have turned right. The green ball is 3 steps left and 4 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a grey key 2 steps left and 1 step forward, You see a blue key 1 step left and 2 steps forward, You see a green key 1 step left and 1 step forward, You see a blue ball 1 step right Critic:I have gone forward. The green ball is 3 steps left and 3 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 3 steps forward, You see a green key 2 steps left, You see a green ball 3 steps left and 2 steps forward, You see a green key 2 steps right, You see a green ball 3 steps left and 2 steps forward, You see a green key 2 steps right, You see a green ball 3 steps left and 2 steps forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is steps left and 1 step right, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward,</pre>
Action:turn right Observation:You see a wall 5 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 4 steps forward, You see a greey key 2 steps left and 2 steps forward, You see a blue key 1 step left and 3 steps forward, You see a green key 1 step left and 2 steps forward, You see a blue ball 1 step right and 1 step forward Critic:I have turned right. The green ball is 3 steps left and 4 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a grey key 2 steps left and 1 step forward, You see a blue key 1 step left and 2 steps forward, You see a green key 1 step left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 3 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 2 steps forward. The green ball is 3 steps left and 3 steps left. You see a blue key 1 step left and 2 steps forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 3 steps left. You see a green ball 3 steps left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left. This green ball is 3 steps left and 1 step forward. This step is (UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left. Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN.
<pre>green ball 3 steps left and 4 steps forward, You see a grey key 2 steps left and 2 steps forward, You see a blue key 1 step left and 3 steps forward, You see a green key 1 step left and 2 steps forward, You see a blue ball 1 step right and 1 step forward Critic:I have turned right. The green ball is 3 steps left and 4 steps forward. This step is UNNNOW. Action:go forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a grey key 2 steps left and 1 step forward, You see a blue key 1 step left and 2 steps forward, You see a green key 1 step left and 1 step forward, You see a blue ball 1 step right Critic:I have gone forward. The green ball is 3 steps left and 3 steps forward. This step is UNENOWN. Action:go forward Observation:You see a wall 3 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 2 steps forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNENOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNENOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNENOWN. Action:Turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a green ball 3 steps forward</pre>
<pre>steps forward, You see a blue key 1 step left and 3 steps forward, You see a green key 1 step left and 2 steps forward, You see a blue ball 1 step right and 1 step forward Critic:I have turned right. The green ball is 3 steps left and 4 steps forward. This step is UNNNOWN. Action:go forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a grey key 2 steps left and 1 step forward, You see a blue key 1 step left and 2 steps forward, You see a green key 1 step left and 1 step forward, You see a blue ball 1 step right Critic:I have gone forward. The green ball is 3 steps left and 3 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 3 steps forward, You see a green key 1 step left and 2 steps left and 2 steps forward, You see a grey key 2 steps left. You see a blue key 1 step left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNONN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNONN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward. This step is ball 3 steps forward</pre>
<pre>forward Critic:I have turned right. The green ball is 3 steps left and 4 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a grey key 2 steps left and 1 step forward, You see a blue key 1 step left and 2 steps forward, You see a green key 1 step left and 1 step forward, You see a blue ball 1 step right Critic:I have gone forward. The green ball is 3 steps left and 3 steps forward. Action:go forward Observation:You see a wall 3 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 2 steps forward, You see a green key 1 step left Critic:I have gone forward. Observation:You see a wall 3 steps forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is</pre>
<pre>Critic:T have turned right. The green ball is 3 steps left and 4 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a grey key 2 steps left and 1 step forward, You see a blue key 1 step left and 2 steps forward, You see a green key 1 step left and 1 step forward, You see a blue ball 1 step right Critic:I have gone forward. The green ball is 3 steps left and 3 steps forward. This step is UNKNOMN. Action:go forward Observation:You see a wall 3 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 2 steps forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps right, You see a blue key 1 step left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOMN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOMN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward</pre>
Action:go forward Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a grey key 2 steps left and 1 step forward, You see a blue key 1 step left and 2 steps forward, You see a green key 1 step left and 1 step forward, You see a blue ball 1 step right Critic:I have gone forward. The green ball is 3 steps left and 3 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 3 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 2 steps forward, You see a greey key 2 steps left, You see a blue key 1 step left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward
<ul> <li>Observation:You see a wall 4 steps forward, You see a green ball 3 steps left and 3 steps forward, You see a grey key 2 steps left and 1 step forward, You see a blue key 1 step left and 2 steps forward, You see a green key 1 step left and 1 step forward, You see a blue ball 1 step right</li> <li>Critic:I have gone forward. The green ball is 3 steps left and 3 steps forward. This step is UNKNOWN.</li> <li>Action:go forward</li> <li>Observation:You see a wall 3 steps forward, You see a green key 1 step left</li> <li>Critic:I have gone forward. The green ball is 3 steps left and 2 steps right, You see a green ball 3 steps left and 2 steps forward, You see a green key 2 steps left, You see a blue key 1 step left and 1 step forward, You see a green key 1 step left</li> <li>Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNOWN.</li> <li>Action:go forward</li> <li>Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left</li> <li>Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN.</li> <li>Action:go forward</li> <li>Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left</li> <li>Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN.</li> <li>Action:go forward</li> <li>Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left</li> <li>Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN.</li> <li>Action:turn left</li> <li>Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 1 step left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a green ball 3 steps forward</li></ul>
<pre>steps left and 2 steps forward, You see a green key 1 step left and 1 step forward, You see a blue ball 1 step right Critic:I have gone forward. The green ball is 3 steps left and 3 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 3 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 2 steps forward, You see a grey key 2 steps left, You see a blue key 1 step left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is</pre>
<pre>You see a blue ball 1 step right Critic:1 have gone forward. The green ball is 3 steps left and 3 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 3 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 2 steps forward, You see a grey key 2 steps left, You see a blue key 1 step left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is</pre>
<pre>Children have gone forward. The green ball is 3 steps left and 3 steps forward. Action:go forward Observation:You see a wall 3 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 2 steps forward, You see a grey key 2 steps left, You see a blue key 1 step left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is</pre>
<pre>Action:go forward Observation:You see a wall 3 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 2 steps forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is</pre>
<pre>doservation:You see a wall 3 steps lotward, You see a wall 2 steps light, You see a green ball 3 steps left and 2 steps forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is</pre>
<pre>a blue key 1 step left and 1 step forward, You see a green key 1 step left Critic:I have gone forward. The green ball is 3 steps left and 2 steps forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is</pre>
<pre>Children have gone forward. The green ball is 5 steps fert and 2 steps forward. Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is</pre>
Action:go forward Observation:You see a wall 2 steps forward, You see a wall 2 steps right, You see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is
<pre>conservation.iou see a wall 2 steps lotward, You see a wall 2 steps light, Hou see a green ball 3 steps left and 1 step forward, You see a blue key 1 step left Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is</pre>
Critic:But I have gone forward. The green ball is 3 steps left and 1 step forward. This step is UNKNOWN. Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is
Action:go forward Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is
Observation:You see a wall 1 step forward, You see a wall 2 steps right, You see a green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is
<pre>green ball 3 steps left Critic:But I have gone forward. The green ball is 3 steps left. This step is UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is</pre>
UNKNOWN. Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is
Action:turn left Observation:You see a wall 1 step right, You see a grey key 2 steps left and 2 steps forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is
forward, You see a green key 2 steps left and 1 step forward, You see a blue key 1 step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is
step left and 1 step forward, You see a red ball 4 steps forward, You see a green ball 3 steps forward Critic:But I have turned left. The green ball is 3 steps forward. This step is
Critic:But I have turned left. The green ball is 3 steps forward. This step is
UNKNOWN. Action:go forward
Observation:You see a wall 1 step right, You see a grey key 2 steps left and 1 step
forward, You see a green key 2 steps left, You see a blue key 1 step left, You see a red ball 3 steps forward.
Critic:But I have gone forward. The green ball is 2 steps forward. This step is
UNKNOWN.
Action:go forward Observation:You see a wall 1 step right, You see a grey key 2 steps left, You see a
red ball 2 steps forward, You see a green ball 1 step forward