

INTERACTIVECOT: ALIGNING DYNAMIC CHAIN-OF-THOUGHT PLANNING FOR EMBODIED DECISION-MAKING

Anonymous authors

Paper under double-blind review

ABSTRACT

Vision-Language Models (VLMs) are increasingly being employed as the decision-making “brains” of embodied agents. Effectively harnessing their powerful generalization capabilities in dynamic, context-specific tasks remains a significant challenge. Chain-of-Thought (COT) prompting is often utilized for complex task execution, but existing methods either rely on static strategies that fail to adapt to changing environments or fine-tune on offline datasets, which are insufficient for optimizing agent decision-making through interaction. In this paper, we propose a novel approach that focuses on optimizing the COT reasoning process rather than just the final action tokens. By aligning the COT process through preference-based reinforcement learning, specifically Direct Preference Optimization (DPO), we enhance the agent’s ability to make accurate decisions in dynamic environments while mitigating model degradation during fine-tuning. Our method models the environment as a Markov decision process, requiring the agent to reflect on the current state in real time to generate adaptive plans and actions. By prioritizing the optimization of the COT process over the final actions, we enhance the agent’s reasoning adaptability while effectively mitigating model degradation during fine-tuning. Experiments in the ALFWorld environment demonstrate an average success rate of **26.67%**, which is a **6%** improvement over RL4VLM, and show that our method effectively mitigates model degradation post fine-tuning. These results highlight the potential of integrating preference-based reinforcement learning techniques with COT processes to enhance the decision-making capabilities of vision-language models in embodied agents.

1 INTRODUCTION

Large Language Models (LLMs) and Large Multimodal Models (LMMs) have achieved remarkable success in natural language understanding and generation tasks (Brown, 2020; Achiam et al., 2023). Recent studies have explored how LLMs can be leveraged to manage other AI models and tools for complex language or multimodal tasks (Shen et al., 2024; Lu et al., 2024), assist in playing sophisticated games such as TextWorld (Yao et al., 2022), Handi (Hu & Sadigh, 2023), and Minecraft (Wang et al., 2023a), or be deployed on robots for real-world interactions (Ahn et al., 2022; Driess et al., 2023). Recently, large multimodal models have garnered significant attention due to their ability to process various input modalities (text, images, videos, etc.). This has spurred increased research in embodied AI, where language-vision models are employed for decision-making and task planning in both simulated environments and the real physical world. While LLMs and Vision-Language Models (VLMs) can provide insightful suggestions for complex generation tasks, they often fail in solving simple decision-making tasks due to misalignment issues (Ahn et al., 2022).

To enhance decision-making capabilities, utilizing Chain-of-Thought (COT) reasoning has become a common approach. COT has been demonstrated to improve model performance in logical reasoning by facilitating the output of correct results through step-by-step reasoning. Mu et al. (2023) enhanced static planning capabilities by fine-tuning models on the EgoCOT dataset, integrating high-level task planning with low-level task control in a closed-loop manner, achieving promising performance in multiple specific tasks. In contrast, dynamic re-planning for decision-making has been shown to be more adaptive than static generation. Song et al. (2023b) introduced a few-shot planning method

054 leveraging in-context learning and a grounded re-planning mechanism to dynamically adjust high-
055 level plans based on environmental observations.

056 Nevertheless, planning based solely on a model’s generative capabilities is insufficient, especially
057 in complex tasks, partially observable scenarios, and multi-task environments. Agent models must
058 possess the ability for continual learning, continuously deriving insights from failures and aligning
059 online within specific task environments to make more accurate decisions. Aligning through rein-
060 forcement learning (RL) is a common approach. RL learns agents’ policies from scratch through trial
061 and error in environments (Sutton, 2018), ensuring that LMM-based agents are well-aligned with
062 their environments. A notable example is RL4VLM (Zhai et al., 2024), which combines Proximal
063 Policy Optimization (PPO) with COT reasoning to fine-tune vision-language models for decision-
064 making tasks. This integration allows the model to learn more effectively from task rewards through
065 interaction, improving exploration, adaptability, and reasoning. RL4VLM proposes to mitigate the
066 effect of the COT reasoning tokens by focusing the primary optimization target on the final action
067 tokens. Most RL methods start with random policies, which are updated based on returns from the
068 environment. This leads to poor sample efficiency, as initial policies perform poorly during the early
069 stages of learning. One way to improve sample efficiency is to incorporate prior knowledge into the
070 policy initialization and exploration during training (Kumar et al., 2022). LLMs are ideal sources
071 of prior knowledge for RL agents, as they are trained on vast amounts of data from diverse corpora.
072 Therefore, leveraging RL to align LLMs with embodied environments for decision-making tasks
073 can simultaneously address the misalignment issues in LLMs and the sample efficiency challenges
074 in RL.

075 Unlike RL4VLM, we believe that the COT process holds the key to optimization. Since the action
076 is the outcome of the COT process and is closely related to it, we focus more on the consistency
077 between the action and the COT. Our framework is based on Direct Preference Optimization (DPO).
078 DPO has recently emerged as a prominent method due to its efficient alignment without the need for
079 reward design, and it is widely used in the post-SFT stage of large models. To our knowledge, there
080 is no precedent for its use in embodied agent tasks. Therefore, we consider introducing the DPO
081 algorithm to efficiently learn strategies from sparse or no-reward interactions. Furthermore, we have
082 made improvements to DPO with a focus on optimizing the COT process and ensuring consistency
083 in model responses, thereby adapting it to our algorithmic framework for interactive alignment of
084 VLMs.

085 In summary, our contributions can be summarized in the following four points:

- 086 1. We propose an algorithmic framework, **InteractiveCOT**, for online alignment of the COT pro-
087 cess in embodied agents through interaction with the environment, supporting both PPO and DPO
088 alignment schemes.
- 089 2. We have made adaptive adjustments to DPO, designing a data sampling and sample pair construc-
090 tion framework tailored to the interaction characteristics of embodied agents, thereby improving the
091 sample utilization efficiency of the alignment algorithm.
- 092 3. We emphasize that aligning the COT is more important than aligning the final action in align-
093 ment tasks. Based on this, we have improved the DPO algorithm to enhance output consistency,
094 alleviating the issue of output degradation during model training.
- 095 4. We validate our approach through experiments in the ALFWorld environment, demonstrating
096 a **6%** increase in average success rates compared to baseline methods. Our results highlight the
097 potential of integrating preference-based reinforcement learning techniques with COT processes
098 to enhance the decision-making capabilities of vision-language models in embodied agents. This
099 advancement highlights the importance of optimizing the thought process itself to achieve better
100 performance and adaptability in complex, dynamic tasks.

101 102 103 2 RELATED WORK

104 **Embodied Agent with LLMs** The successful integration of language as a semantically rich input
105 for interactive decision-making highlights the crucial role of LLMs in facilitating interaction and
106 decision-making processes (Abramson et al., 2020; Karamcheti et al., 2022; Li et al., 2022). LLMs
107 are also applied in various environments to aid robot navigation (Parisi et al., 2022; Hong et al.,

2021; Majumdar et al., 2020) and manipulation (Jiang et al., 2022; Ren et al., 2023; Karamcheti et al., 2022). Recently, there have been a large number of methods that utilize LLMs to enhance agents’ planning and reasoning capabilities in embodied agents. SayCan (Ahn et al., 2022) assesses the affordance of candidate actions by multiplying their probabilities under LLMs with a value function. (Zeng et al., 2022) combine the LLM with a visual-language model and a pre-trained language-conditioned policy (Shridhar et al., 2022) to enable open vocabulary robotic tasks. (Huang et al., 2022a) demonstrate that LLMs can be employed for planning and executing simple household tasks. They ground LLM-generated actions by comparing their embeddings with a pre-defined list of acceptable actions. To incorporate environment feedback, Inner Monologue (Huang et al., 2022b) extends SayCan using a closed-loop principle. This principle is also applied in related works such as (Yao et al., 2023; Huang et al., 2022b; Kim et al., 2024; Singh et al., 2023; Liang et al., 2023; Shinn et al., 2023; Wang et al., 2023b) to continuously monitor agent behaviors and refine and adjust plans accordingly for tasks such as computer automation, Minecraft, etc. Furthermore, there are approaches that prompt LLMs to generate temporal-abstracted actions (Zheng et al., 2023). (Dasgupta et al., 2023) employ the LLM as a planner and success detector for an agent with their actor module necessitates pre-training with RL to enable the agent to follow natural language instructions. While these works demonstrate impressive results, they rely too heavily on the inherent capabilities of powerful LLMs, like GPT4 and PaLM (Chowdhery et al., 2023), which are difficult to apply to smaller LLMs with weaker reasoning abilities, like LLaMA-7B.

Concurrent to our work, GLAM (Carta et al., 2023) utilizes RL finetuning to achieve functional grounding of LLMs. However, they focus on simple primitive actions (turn left, turn right, go forward, etc.) evaluated in toy environments, BabyAI (Chevalier-Boisvert et al., 2018) with a much smaller encoder-decoder LLM, Flan-T5-780M. These primitive actions have a similar number of tokens and less meaningful semantics, resulting in underutilizing the capabilities of LLMs, and failing to observe the impact of prompt design and address the unbalance over action space, resulting in additional instability and poor robustness.

Preference Learning Preference learning has become a pivotal area in machine learning, aiming to develop predictive models that capture human preferences from observational data. Recent advances in deep learning and optimization algorithms have driven significant progress in this field, particularly in applications such as recommender systems, information retrieval, and personalized user interfaces.

Current preference learning methods can be categorized into pointwise, pairwise, and listwise approaches. Among these, Direct Preference Optimization (DPO) (Rafailov et al., 2024) has emerged as a novel and efficient paradigm, directly optimizing user preferences without intermediary ranking steps. DPO achieves more precise alignment with user preferences by constructing loss functions that directly reflect these preferences. Chen et al. (2024) introduces OPTune, an efficient method for online preference tuning in RLHF. By selectively regenerating low-reward responses and using a weighted DPO loss to focus on response pairs with larger reward gaps, OPTune improves training speed and model alignment while reducing computational costs. Recent pioneering studies have further expanded DPO’s applications and effectiveness. Step-DPO Lai et al. (2024) stands out as a significant advancement over Direct Preference Optimization (DPO) for tasks requiring long-chain reasoning, such as mathematical problem-solving. Unlike DPO, Step-DPO optimizes individual reasoning steps. By focusing on pinpointing the first erroneous step in a sequence and optimizing for more fine-grained accuracy, Step-DPO improves both factuality and reasoning in large language models. Pal et al. (2024), in their DPO-Positive study, advanced practical applications of DPO by focusing on positive direct preference optimization in sentiment-aware recommendations. The DPO-Positive method not only enhances user satisfaction but also incorporates sentiment information into the recommendation process, resulting in more accurate and user-aligned outcomes.

3 METHODS

3.1 ONLINE TRAINING OF REPLANNING FRAMEWORK

In previous work (Ahn et al., 2022; Song et al., 2023a), long-term planning using large language models (LLM) or large multimodal models (LMM) has typically been approached as static planning, the transition and planning between the initial and final states of a task is accomplished

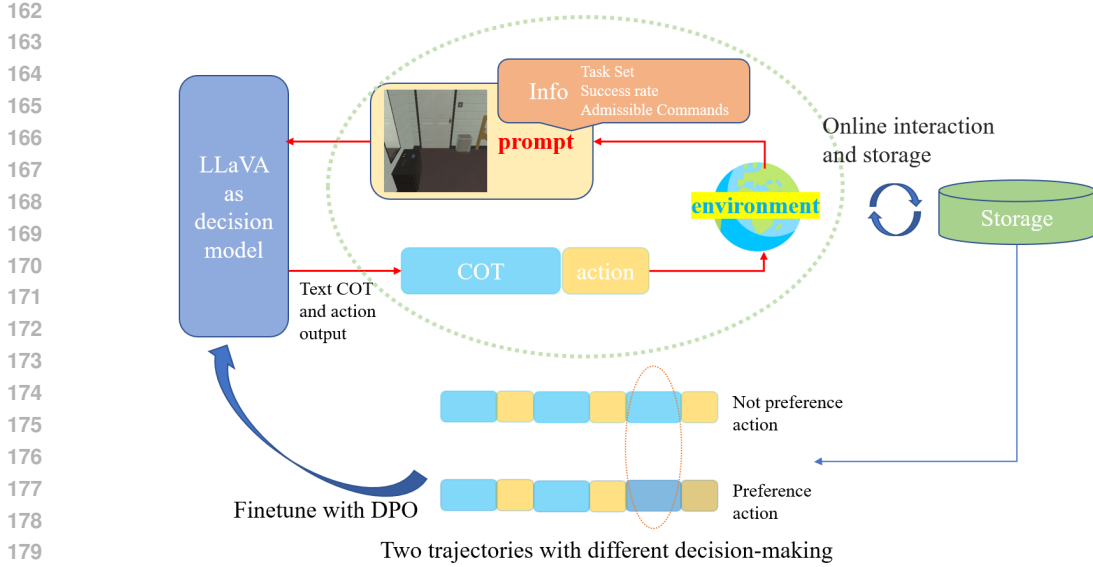


Figure 1: Main framework of our method. In our approach, we sample two different trajectories for the same step x and assess their preference based on task completion rates. This allows us to determine the preference level of different actions in the current state. We then fine-tune the vision-language model (VLM) using preference methods such as Direct Preference Optimization (DPO).

through a single planning instance. For once planning instance, planner output the planning results $\{planner : s_0, a_1, a_2, \dots, a_n, s_{goal}\}$, where s denote states and a for actions. This static planning often considers the completeness of adjacent decisions and planning costs. However, it frequently lacks the capability to timely correct erroneous plans. Consequently, the feedback provided by the environment following each planning action may not be utilized in a timely manner to adjust subsequent actions. The primary distinction between re-planning and static planning processes lies in the ability to make timely adjustments based on environmental feedback. Re-planning captures factors that change dynamically within the environment, providing different responses based on various states on each decision-making steps and generating new execution plans. Compared to static planning, this approach offers greater adaptability and robustness. Additionally, dynamic planning involves deeper interaction with the environment.

Our algorithm is designed based on the re-planning framework, which can be seen in Figure 1. Specifically, in each natural step of interaction between the agent and the environment, the planning result of next several time-steps is regenerated according to the current observation. We incorporate camera images and environmental feedback into the design of prompts, providing feedback to the agent at each step, requiring it to give subsequent plans step-by-step based on observations. Under the re-planning framework, a well-fine-tuned model base can already perform quite well. However, there are still some complex situations that the agent cannot handle effectively, such as navigation tasks where target objects are not observable and complex tasks with numerous steps. We ponder *whether the agent can learn and improve its planning abilities through interaction with the environment*. Reinforcement learning algorithms are a good choice, but they require a precisely designed reward function and must also consider potential reward hacking phenomena. Based on this, our framework aligns through the construction of preference sample pairs. During the interaction process, we sample and use the success rate of trajectories as a preference for alignment. When applying DPO to finetune VLM, the loss function is

$$L(\theta) = -\mathbb{E}_{\zeta} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(a_{win}^t | \tau_1^{t-1})}{\pi_{ref}(a_{win}^t | \tau_1^{t-1})} - \beta \log \frac{\pi_{\theta}(a_{lose}^t | \tau_2^{t-1})}{\pi_{ref}(a_{lose}^t | \tau_2^{t-1})} \right) \right] \quad (1)$$

where ζ is a tuple of $(o_t, s_{t-1}, a_{win}, a_{lose})$, $s_{1 \sim t-1}$ is the state-action decision trajectory from time 0 to time $t-1$, a_{win} and a_{lose} represent the preference and not preference decision-making actions at current timestep t . π_{θ} and π_{ref} denote the policy generated by VLM, π_{ref} refers to the policy generated by the unrefined output of VLM. This approach enables a more granular alignment of

216 preferences, focusing on fine-tuning for each individual decision rather than aligning preferences
 217 for the entire trajectory. Consequently, while introducing step-wise preference information, we
 218 improve upon previous methods that used the PPO algorithm for fine-tuning VLMs by adopting
 219 DPO for preference alignment. This allows us to align preferences for different decisions at each
 220 segment of the trajectory.

222 3.2 TRAINING COT WITH DPO

224 In completing each task, the input prompt for the VLM includes observations and action trajectories.
 225 Through the design of input prompts and instruction-following mechanisms, VLM can generate de-
 226 cision actions based on the current state, produce feasible actions, and provide a textual description
 227 of the current observation along with the formatted output. This structured approach enables the
 228 model to maintain context and make informed decisions effectively. The response output by the
 229 VLM includes a chain-of-thought reasoning process for the current action, immediately followed by
 230 the keyword “action”, which indicates the model’s current decision action. This structured format
 231 allows for clarity in the decision-making process, ensuring that the reasoning is explicitly linked to
 232 the chosen action.

233 3.2.1 THE METHOD OF CONSTRUCTING SAMPLE PAIRS

235 Compared to classic MLP-based policy networks, a advantage of VLM policies is that they can
 236 output neutral language, thus leverage COT reasoning for efficient exploration by performing inter-
 237 mediate reasoning steps that lead to the final decision. However, training a VLM policy π_θ with
 238 RL presents additional challenges. First, due to the sparse rewards obtained from the online interac-
 239 tions between VLM and the environment, many state transition processes receive a reward feedback
 240 scalar value of 0. In the case of state transition samples with a reward value of 0, employing the PPO
 241 (Proximal Policy Optimization) architecture for fine-tuning the VLM makes it challenging for the
 242 model to learn effective strategies for interacting with the environment. Consequently, the sample
 243 efficiency of fine-tuning the VLM using these state transition samples is relatively low. In some
 244 studies, researchers often design reward functions manually to mitigate the issue of sparse rewards.
 245 On the other hand, preference-based methods can construct preference pairs using different reward
 246 values $\{\tau_{win}^t = \{a_t^1, r_t^1, \tau_1^{t-1}\}, \tau_{lose}^t = \{a_t^2, r_t^2, \tau_1^{t-1}\}\}$, whereby trajectories with higher reward
 247 values can be treated as preferred trajectories, for example, τ_{win} has a higher reward r_t^1 and τ_{lose}
 248 has a lower reward r_t^2 . This approach allows for a more nuanced representation of preferences, fa-
 249 cilitating the learning process in environments characterized by sparse feedback. By employing this
 250 method, preference-based approaches can effectively leverage state transition samples with lower
 251 reward values, thereby enhancing sample efficiency. This strategy allows the model to learn from
 252 a broader range of experiences, improving its ability to identify and optimize preferred trajectories
 253 within the environment.

254 3.2.2 THINKING IS MORE IMPORTANT THAN DECISION-MAKING

255 It is worth noting that by outputting the text of the chain of thought, we enable the VLM to produce
 256 reasonable actions through autonomous reasoning. However, the reasoning ability of the VLM stems
 257 from its training on massive datasets. Jointly fine-tuning the COT text and the text actions output by
 258 the VLM is proven to be a better method adapting the VLM to embodied scenarios. This method
 259 often focuses on optimizing action output, and when designing loss functions, it tends to minimize
 260 the impact of the COT process or only consider the final action decision. The fine-tuning process
 261 breaks the coherence of the language output formed during pre-training, leading to model collapse.
 262 To solve this problem, we used two methods to constrain the fine-tuning process. First, we constrain
 263 the distributional distance between the fine-tuned output text and the unfine-tuned reference model
 264 output text, thus ensuring that the model does not deviate too much from the logicity of the original
 265 model language output due to fine-tuning. Similar to the derivation process of the DPO model, we
 266 set

$$267 Q(s, a) = \beta \log \frac{\pi_\theta(a|s)}{\pi_{ref}(a|s)} \quad (2)$$

268 Then, we can fine-tune the output strategy of VLM by optimizing the Q value, while limiting the
 269 output distance between the fine-tuned model and the reference model without fine-tuning by adding

a regularization term of KL divergence to the optimization objective, which is as follows

$$\max_{\pi_{\theta}} \mathbb{E}_{s \sim D, a \sim \pi_{\theta}(a|s)} [Q(s, a)] - \beta D_{\text{KL}}[\pi_{\theta} \parallel \pi_{\text{ref}}] \quad (3)$$

Based on this optimization objective, combined with some mathematical derivations of Yang et al. (2024), we can derive the following step-wise optimization formula:

$$L = -\mathbb{E}_{\zeta} \log \sigma \left(\beta \log \frac{p(a_1^t | \mathcal{T}_1^t) \pi_{\theta}(\mathcal{T}_1^t | \tau_1^{t-1})}{\pi_{\text{ref}}(a_1^t, \mathcal{T}_1^t | \tau_1^{t-1})} - \beta \log \frac{p(a_2^t | \mathcal{T}_2^t) \pi_{\theta}(\mathcal{T}_2^t | \tau_2^{t-1})}{\pi_{\text{ref}}(a_2^t, \mathcal{T}_2^t | \tau_2^{t-1})} \right) \quad (4)$$

$$\nabla_{\theta} L = -\beta \mathbb{E}_{\zeta} [\Lambda [\nabla_{\theta} \log \pi_{\theta}(\mathcal{T}_1^t | \tau_1^{t-1}) - \nabla_{\theta} \log \pi_{\theta}(\mathcal{T}_2^t | \tau_2^{t-1})]] \quad (5)$$

The detailed derivation can be found in Appendix A. In the formula, \mathcal{T}_i^t is the output text of the thinking chain at step t , and a is the output action after the thinking chain. However, by calculating the gradient of 4, we can see that the gradient term in Equation 5 directly eliminates the influence of action probabilities, where $\Lambda = \sigma(\hat{Q}_{\theta}(a_1^t, \mathcal{T}_1^t, \tau_1^{t-1}) - \hat{Q}_{\theta}(a_2^t, \mathcal{T}_2^t, \tau_2^{t-1}))$. Therefore, in practice, we adopt the following action probability weighting (APW) form:

$$\tilde{L} = -\mathbb{E}_{\zeta} \log \sigma \left(\beta p(a_1^t | \mathcal{T}_1^t) \log \frac{\pi_{\theta}(\mathcal{T}_1^t | \tau_1^{t-1})}{\pi_{\text{ref}}(a_1^t, \mathcal{T}_1^t | \tau_1^{t-1})} - \beta p(a_2^t | \mathcal{T}_2^t) \log \frac{\pi_{\theta}(\mathcal{T}_2^t | \tau_2^{t-1})}{\pi_{\text{ref}}(a_2^t, \mathcal{T}_2^t | \tau_2^{t-1})} \right) \quad (6)$$

We will later analyze the errors of both and demonstrate that our approach is feasible. Intuitively, the gradient term of the ‘‘action’’ probability adds weight to the ‘‘thoughts’’ probability. Actions with higher output probabilities after COT indicate a better alignment with the thoughts process, while lower probabilities suggest greater randomness in action generation. The weighting term can reduce the generation of highly random positive samples and encourage the generation of deterministic positive samples.

3.2.3 MAKE DECISIONS THAT ARE MORE CONSISTENT WITH REASONING

In the second method, we consider adding a regularization term to the final action text output. Fine-tuning the reasoning chain may alter the model’s language output conventions, potentially leading to model collapse. Aligning the action text output with the reference standard model ensures that the model adheres to the prompt’s formatting requirements, thus generating valid actions. This approach helps maintain the integrity of the output while allowing for effective task execution, the effects of this regularization can be observed in Figure 2.

To further ensure the consistency between the COT process of generating text and the final action, we consider adding a stronger constraint to the above formula. We believe that the pre-trained model has already been well-optimized for modeling the process from thoughts to actions. Therefore, in the subsequent interaction phase, we will maintain alignment with the pre-training results in the dimension of generating actions based on thoughts. We will add a mean square error (MSE) regularization term of action policy consistency constraint (APC), specifically:

$$L_{\text{InteractiveCOT}} = \tilde{L} + \kappa \text{MSE}(\pi_{\theta}(a_1^t | \mathcal{T}_1^t), \pi_{\text{ref}}(a_1^t | \mathcal{T}_1^t)) \quad (7)$$

Where κ is a hyper-parameter that regulates the strength of constraints. Note that this is different from the KL divergence used in the DPO derivation process with the reference model; Here, the focus is more on aligning the process of deriving actions from thoughts. In the comparison experiment of action consistency constraint, we extracted an output sample at the training step of 2000, both with and without the action probability consistency constraint, as shown in Figure 2. The sample clearly demonstrates the difference between the COT process and the final decision action output. In the sample with the action consistency constraint, the model’s COT process provides a clear analysis and identifies the valid action to be taken next. In contrast, in the sample without the consistency constraint, although the agent provides reasonable thoughts and analysis, it ultimately outputs an irrelevant action in the final decision, which does not fall within the scope of valid actions.

Finally, we provide a simple illustration of Equation 6 to demonstrate that our approximation is reasonable. We assume that the pre-trained model has achieved good alignment, so $p(a|\mathcal{T})$ will be close to 1. We have:

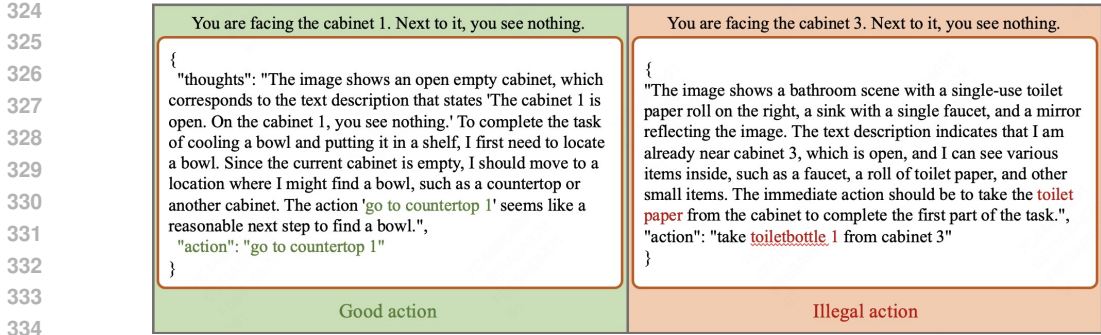


Figure 2: Use or not use action policy consistency constraint cases. When employing the use action policy consistency constraint strategy, the output actions are more likely to be valid actions. This approach helps to ensure that the generated actions align with the established policy, thereby enhancing the reliability and appropriateness of the actions in the context of the task being performed.

$$\Delta(p(a_i^t|\mathcal{T}_i^t)) = \log \frac{p(a_i^t|\mathcal{T}_i^t)\pi_\theta(\mathcal{T}_i^t|\tau_i^{t-1})}{\pi_{ref}(a_i^t, \mathcal{T}_i^t|\tau_i^{t-1})} - p(a_i^t|\mathcal{T}_i^t) \log \frac{\pi_\theta(\mathcal{T}_i^t|\tau_i^{t-1})}{\pi_{ref}(a_i^t, \mathcal{T}_i^t|\tau_i^{t-1})} \quad (8)$$

$$= \log p(a_i^t|\mathcal{T}_i^t) + (1 - p(a_i^t|\mathcal{T}_i^t)) \log \frac{\pi_\theta(\mathcal{T}_i^t|\tau_i^{t-1})}{\pi_{ref}(a_i^t, \mathcal{T}_i^t|\tau_i^{t-1})} \quad (9)$$

This variable will approach zero as $p(a|\mathcal{T})$ approaches 1. In practice, we have calculated the approximate distribution of action probabilities and demonstrated that our assumption is well-founded, which can be shown in Figure 6b.

4 EXPERIMENTS

In this part we perform experiments to validate three key questions:

- How does our framework enhance the decision-making capabilities of VLM?
- Can the regularization term effectively constrain the action distribution to prevent deviation from the original policy?
- Does action-weighting mitigate the issue of degradation?

we conducted experiments in the ALFWorld environment and recorded improvements in the visual semantic reasoning capabilities of the Vision-Language Model. ALFWorld encompasses six types of household planning tasks: *Pick & Place*, *Pick Two & Place*, *Clean & Place*, *Cool & Place*, *Heat & Place*, and *Examine in Light*. For convenience, we will refer to them as **Pick**, **Pick2**, **Clean**, **Cool**, **Heat** and **Look** hereafter. During the experiments, the agent captures a visual observation through egocentric view in the current state and a textual instruction describing the task to be completed. The agent must plan and navigate based on the visual information to accomplish the specified tasks. We instantiate our method on top of the *llava-v1.6-mistral-7b* model, and build the agent based on this model. During interactive, we package the observation picture into a special prompt to get LLaVA’s answer.

Prompt Our COT prompt consists of the following parts: First, we clarify the task requirements. The tasks in ALFWorld are semantically rearranged. For example, both “examine the pillow with the deskclam” and “look at the pillow under the deskclam” indicate that the agent needs to find the pillow, pick it up, then locate and navigate to the deskclamp. Secondly, we specify the range of valid actions. Each state in ALFWorld environment is accompanied by different valid action transitions. For instance, if the action is to pick up pillow 1, the prerequisite is that the agent must be close enough to reach the pillow. If the action is to put down an object, it must have previously executed the pick-up action. Therefore, one of the criteria for evaluating the agent’s capability in the experiment

is whether the actions it outputs are valid. Finally, we specify the output format of the agent’s LLaVA model, which must strictly follow the JSON format containing “thoughts” and “action”. The action must be derived from the thoughts and should not produce irrelevant actions. Our prompt design is shown in Figure 3:

```
Your are an expert in the ALFRED Embodied Environment.
Your task is to * task name *. You are also given the following text description of the current scene: * obs *.
Your admissible actions of the current situation are: [* reformatted admissible actions *].
Your response should be a valid Json file in the following format:
"thoughts": "{first describe what do you see in the image using the text description, then carefully think about
which action to complete the task. }",
"reflections": "{reflect on your historical trajectory and carefully think about which action to complete the task.}",
"action": "{an admissible action}"
your actions should be based solely on the analysis provided by your thoughts!
your output need to be in 60 words!
```

Figure 3: Prompt used in ALFWorld tasks. The prompt provides the embodied agent with several key components: the task to be completed, the current egocentric observations, the feasible actions available, the output format for the thought process in the reasoning chain, the format for action text output, and constraints on the length of the output text. This structured approach helps ensure that the VLM can generate coherent and contextually relevant responses, facilitating effective decision-making and task execution.

Implementation Our implementation consists of two parts. First, we conduct one epoch of model SFT (supervised fine-tuning) on the open-source dataset LEVI-Project/sft-data (Zhai et al., 2024) to ensure the model’s ability of formatted output . The LEVI-Project/sft-data dataset is an expert trajectory dataset sampled by a GPT-4-based agent, containing 45k different state samples, each adhering to the JSON format of COT outputs. After SFT, we employ the model to interact with the environment, optimizing its COT capabilities during these interactions and monitoring performance in real-time during training.

4.1 HOW MUCH BETTER WE ARE AT MAKING DECISIONS

The aim of experiments in this section is to validate the performance of the InteractiveCOT method. To evaluate whether the algorithm can consistently generate decisions through the COT process, we use the success rate of task execution as a reference and select PPO from the RL4VLM (Zhai et al., 2024) framework as the baseline. ALFWorld does not provide a reward function during interactions; it only indicates whether the current task is successfully executed and returns the task’s progress. For instance, if a task requires checking an object under a table lamp, finding and picking up the object results in a 50% progress update. Given that such progress updates are sparse in a larger action space, we construct preference criteria for preference learning. The preference score for each trajectory is calculated using Equation 10:

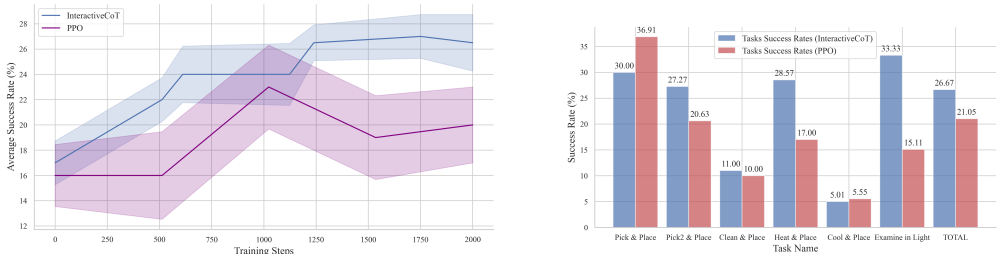
$$P = 50 * success\ rate - \mathbb{K}_{\{invalid\}} \tag{10}$$

$$\mathbb{K}_{\{invalid\}} = \begin{cases} 1 & \text{if } action \text{ not in } admissible\ action \\ 0 & \text{otherwise} \end{cases}$$

where $\mathbb{K}_{\{invalid\}}$ represents our stronger rejection of illegal actions given the same success rate. During the exploration phase, the agent collects trajectory data and constructs sample pairs based on the six task types mentioned above. Higher preference scores indicate greater sample preference. In practice, considering the achievement of long-term goals, we calculate preference scores using a method similar to discount factor weighting in reinforcement learning returns. Due to the high randomness of ALFWorld, we set up experimental environments with different seeds and calculated mean and variance of each results.

We use Equation 7 for the model weight update with $\kappa = 0.1$, measure the agent’s performance by the average success rate of each task. The final results are shown in Figure 4.

432
433
434
435
436
437
438
439
440



(a) Average success rate during training (all 6 tasks) (b) Average success rate of each task at 2000 steps

441
442
443
444
445
446
447
448
449
450

Figure 4: We demonstrate that fine-tuning the vision-language model (VLM) using Interactive Chain of Thought (InteractiveCOT) and Proximal Policy Optimization (PPO) results in varying task completion rates and average task completion rates in ALFworld. Our findings indicate that, for the majority of tasks, fine-tuning the model using preference methods yields better results than using reinforcement learning approaches. Additionally, we observe that through online interaction with the environment, the preference method achieves the same average task completion rate with fewer interaction steps, indicating higher sample efficiency and more minimal model degradation.

451
452
453
454
455
456

ALFWorld gives task randomly so we calculate the overall success rate as the weighted average of success rates under all tasks. InteractiveCOT shows an improvement in the overall success rate, indicating that our algorithm can learn more efficiently from interactions. In our experiments, we used approximations such as $\log(\pi(a|\mathcal{T}, \tau)\pi(\mathcal{T}|\tau)) \approx \pi(a|\mathcal{T}, \tau) \log(\pi(\mathcal{T}|\tau))$ when $\pi(a|\mathcal{T}, \tau) \rightarrow 1$. We calculated the occurrence probability distribution of action tokens in the experiments to demonstrate that our approximations are reasonable.

457
458

4.2 WHAT ROLE DOES ACTION POLICY CONSISTENCY CONSTRAINT PLAY?

459
460
461
462
463
464
465
466

We pointed out that during training, to enhance stability, we introduced the regularization of action token probabilities between finetune model and reference model. This section will explore the impact of regularization on the results and investigate its role. We designed ablation experiments, where we conducted trials with different regularization weight values κ under the same parameter settings, and recorded the average success rate of the agent during training. In this experiment, we use Equation 1 with the regular term as the loss function, with other conditions the same as in Section 4.1.

467
468
469
470
471
472
473
474
475
476
477
478
479
480

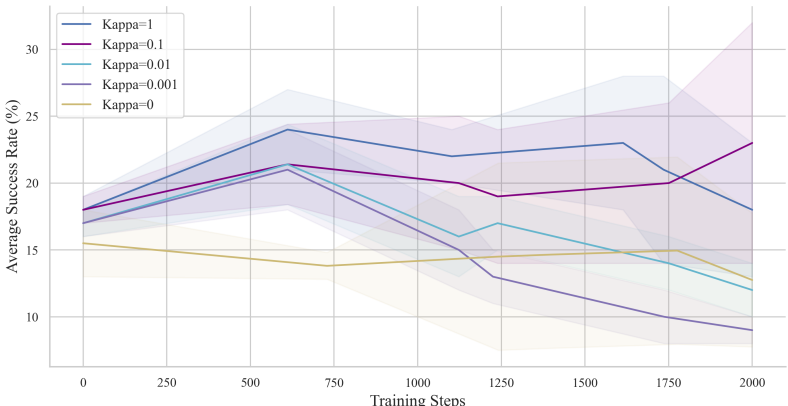


Figure 5: Parameter study of APC

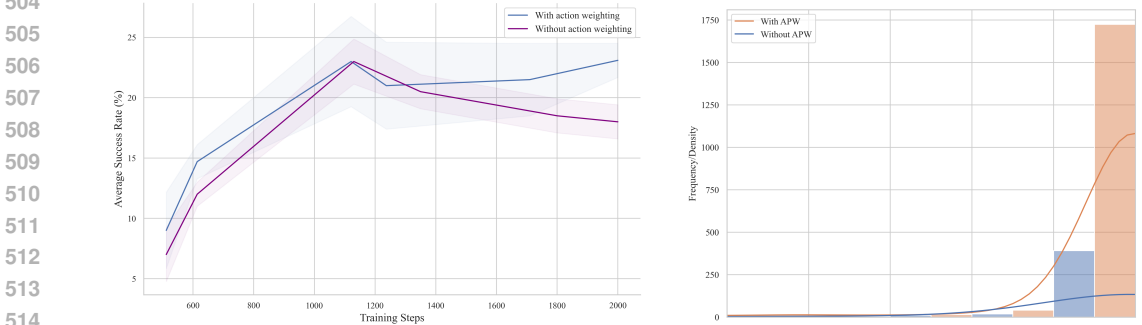
481
482
483
484
485

The results in Figure 5 show that different values of κ significantly impact the success rate. As the parameter increases, the action policy consistency constraint strengthens, leading to improved model performance. This validates the importance of regularization. However, when κ is set to 1, the algorithm’s performance declines, indicating that κ should neither be too large nor too small, with

486 a value around 0.1 yielding near-optimal performance. Given the importance of the κ parameter, its
 487 optimal value may vary across different environments or tasks. Due to space constraints, we do not
 488 explore this further in this paper.
 489

490 4.3 MORE CERTAIN, MORE STABLE
 491

492 In Section 3.2.2, we mentioned that the gradient weighted by action probability would prefer more
 493 certain successful strategies, which indirectly achieves the unification of thoughts and actions—the
 494 larger the conditional probability of an action, the more closely it is linked to the content of the COT.
 495 To verify this idea, we conducted following ablation experiments. We used loss functions with and
 496 without action probability weighting, Equation 6 and Equation 4, keeping all other settings identical
 497 to the main experiment. Figure 6a shows the comparison between the two sets of experiments,
 498 and Figure 6b presents the probability distribution of all actions in the first 2000 training steps.
 499 It is evident that under the APW condition, the probability distribution of action tokens is mostly
 500 concentrated around 1, indicating more certain and robust decision-making, which also leads to a
 501 higher success rate. In contrast, the results without weighting show a more dispersed distribution
 502 of action probabilities, with some probabilities falling below 0.8, which is not conducive to the
 503 convergence of the algorithm.
 504



515 (a) Comparison of the average success rates between APW
 516 and non-APW methods

(b) Action tokens probability distribution

517 Figure 6: Validate the impact of APW in interactions.
 518

519
 520
 521 5 CONCLUSIONS, LIMITATIONS AND FUTURE DIRECTIONS
 522

523 This work introduces an algorithmic framework, InteractiveCOT, for online interactive fine-tuning
 524 of multimodal models during the COT process, supporting both PPO and DPO algorithms. Based
 525 on LLaVA-7B, we execute household tasks in embodied scenarios through dynamic replan, achiev-
 526 ing better decision-making by aligning COT capabilities. We emphasize the core importance of
 527 COT, moving away from previous approaches that primarily focused on training actions. Instead,
 528 we maintain the consistency between COT and actions through APW and APC. Empirical results
 529 demonstrate that InteractiveCOT outperforms reinforcement learning algorithms in average perfor-
 530 mance within ALFWorld. Ablation studies further confirm the critical role of APW and APC in the
 531 algorithm’s convergence effectiveness.

532 One limitation of this study is the lack of validation across a broader range of environments and tasks,
 533 which will be addressed in future work. We aim to further optimize the generalization performance.
 534 Another limitation is the consideration of non-Markovian processes. Since the pre-training datasets
 535 in the SFT phase are all Markovian, our interaction experiments were conducted under the same
 536 conditions. Non-Markovian processes are more common in complex decision-making tasks, and
 537 effectively handling historical information is a crucial capability for agents. In future work, we
 538 will first deploy our algorithm framework in more simulated environments and datasets to enrich the
 539 experimental results. Additionally, we will consider modeling non-Markovian processes, focusing
 on the agent’s performance with long historical information.

REFERENCES

- 540
541
542 Josh Abramson, Arun Ahuja, Iain Barr, Arthur Brussee, Federico Carnevale, Mary Cassin, Rachita
543 Chhaparia, Stephen Clark, Bogdan Damoc, Andrew Dudzik, et al. Imitating interactive intelli-
544 gence. *arXiv preprint arXiv:2012.05672*, 2020.
- 545
546 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
547 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
548 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 549
550 Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea
551 Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say:
552 Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- 553
554 Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- 555
556 Thomas Carta, Clément Romac, Thomas Wolf, Sylvain Lamprier, Olivier Sigaud, and Pierre-Yves
557 Oudeyer. Grounding large language models in interactive environments with online reinforcement
558 learning. In *International Conference on Machine Learning*, pp. 3676–3713. PMLR, 2023.
- 559
560 Lichang Chen, Jiuhai Chen, Chenxi Liu, John Kirchenbauer, Davit Soteliya, Chen Zhu, Tom Gold-
561 stein, Tianyi Zhou, and Heng Huang. Optune: Efficient online preference tuning, 2024. URL
562 <https://arxiv.org/abs/2406.07657>.
- 563
564 Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Saharia,
565 Thien Huu Nguyen, and Yoshua Bengio. Babyai: A platform to study the sample efficiency of
566 grounded language learning. *arXiv preprint arXiv:1810.08272*, 2018.
- 567
568 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
569 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
570 Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):
571 1–113, 2023.
- 572
573 Ishita Dasgupta, Christine Kaeser-Chen, Kenneth Marino, Arun Ahuja, Sheila Babayan, Felix Hill,
574 and Rob Fergus. Collaborating with language models for embodied reasoning. *arXiv preprint*
575 *arXiv:2302.00763*, 2023.
- 576
577 Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter,
578 Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multi-
579 modal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- 580
581 Yicong Hong, Qi Wu, Yuankai Qi, Cristian Rodriguez-Opazo, and Stephen Gould. Vln bert: A
582 recurrent vision-and-language bert for navigation. In *Proceedings of the IEEE/CVF conference*
583 *on Computer Vision and Pattern Recognition*, pp. 1643–1653, 2021.
- 584
585 Hengyuan Hu and Dorsa Sadigh. Language instructed reinforcement learning for human-ai coordi-
586 nation. In *International Conference on Machine Learning*, pp. 13584–13598. PMLR, 2023.
- 587
588 Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot
589 planners: Extracting actionable knowledge for embodied agents. In *International conference on*
590 *machine learning*, pp. 9118–9147. PMLR, 2022a.
- 591
592 Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan
593 Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through
594 planning with language models. *arXiv preprint arXiv:2207.05608*, 2022b.
- 595
596 Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen, Li Fei-
597 Fei, Anima Anandkumar, Yuke Zhu, and Linxi Fan. Vima: General robot manipulation with
598 multimodal prompts. *arXiv preprint arXiv:2210.03094*, 2(3):6, 2022.
- 599
600 Siddharth Karamcheti, Megha Srivastava, Percy Liang, and Dorsa Sadigh. Lila: Language-informed
601 latent actions. In *Conference on Robot Learning*, pp. 1379–1390. PMLR, 2022.

- 594 Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks.
595 *Advances in Neural Information Processing Systems*, 36, 2024.
596
- 597 Aviral Kumar, Rishabh Agarwal, Xinyang Geng, George Tucker, and Sergey Levine. Offline q-
598 learning on diverse multi-task data both scales and generalizes. *arXiv preprint arXiv:2211.15144*,
599 2022.
- 600 Xin Lai, Zhuotao Tian, Yukang Chen, Senqiao Yang, Xiangru Peng, and Jiaya Jia. Step-dpo: Step-
601 wise preference optimization for long-chain reasoning of llms, 2024. URL <https://arxiv.org/abs/2406.18629>.
602
603
- 604 Shuang Li, Xavier Puig, Chris Paxton, Yilun Du, Clinton Wang, Linxi Fan, Tao Chen, De-An Huang,
605 Ekin Akyürek, Anima Anandkumar, et al. Pre-trained language models for interactive decision-
606 making. *Advances in Neural Information Processing Systems*, 35:31199–31212, 2022.
- 607 Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and
608 Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE*
609 *International Conference on Robotics and Automation (ICRA)*, pp. 9493–9500. IEEE, 2023.
610
- 611 Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu,
612 and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language mod-
613 els. *Advances in Neural Information Processing Systems*, 36, 2024.
- 614 Arjun Majumdar, Ayush Shrivastava, Stefan Lee, Peter Anderson, Devi Parikh, and Dhruv Batra.
615 Improving vision-and-language navigation with image-text pairs from the web. In *Computer*
616 *Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings,*
617 *Part VI 16*, pp. 259–274. Springer, 2020.
- 618 Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang, Mingyu Ding, Jun Jin, Bin Wang, Jifeng
619 Dai, Yu Qiao, and Ping Luo. Embodiedgpt: Vision-language pre-training via embodied chain of
620 thought, 2023. URL <https://arxiv.org/abs/2305.15021>.
621
- 622 Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White.
623 Smaug: Fixing failure modes of preference optimisation with dpo-positive, 2024. URL <https://arxiv.org/abs/2402.13228>.
624
- 625 Simone Parisi, Aravind Rajeswaran, Senthil Purushwalkam, and Abhinav Gupta. The unsurprising
626 effectiveness of pre-trained vision models for control. In *international conference on machine*
627 *learning*, pp. 17359–17371. PMLR, 2022.
628
- 629 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and
630 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model,
631 2024. URL <https://arxiv.org/abs/2305.18290>.
- 632 Allen Z Ren, Bharat Govil, Tsung-Yen Yang, Karthik R Narasimhan, and Anirudha Majumdar.
633 Leveraging language for accelerated learning of tool manipulation. In *Conference on Robot*
634 *Learning*, pp. 1531–1541. PMLR, 2023.
635
- 636 Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugging-
637 gpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information*
638 *Processing Systems*, 36, 2024.
- 639 Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and
640 Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023. URL
641 <https://arxiv.org/abs/2303.11366>.
- 642 Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic
643 manipulation. In *Conference on robot learning*, pp. 894–906. PMLR, 2022.
644
- 645 Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter
646 Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using
647 large language models. In *2023 IEEE International Conference on Robotics and Automation*
(ICRA), pp. 11523–11530. IEEE, 2023.

648 Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su.
649 Llm-planner: Few-shot grounded planning for embodied agents with large language models. In
650 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2998–3009,
651 2023a.

652 Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M. Sadler, Wei-Lun Chao, and Yu Su. Llm-
653 planner: Few-shot grounded planning for embodied agents with large language models, 2023b.
654 URL <https://arxiv.org/abs/2212.04088>.

655 Richard S Sutton. Reinforcement learning: An introduction. *A Bradford Book*, 2018.

656
657 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan,
658 and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models.
659 *arXiv preprint arXiv:2305.16291*, 2023a.

660
661 Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe,
662 explain, plan and select: Interactive planning with large language models enables open-world
663 multi-task agents. *arXiv preprint arXiv:2302.01560*, 2023b.

664
665 Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiabin Chen, Qimai Li, Weihang Shen, Xiaolong Zhu,
666 and Xiu Li. Using human feedback to fine-tune diffusion models without any reward model, 2024.
667 URL <https://arxiv.org/abs/2311.13231>.

668
669 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
670 React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*,
671 2022.

672
673 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
674 React: Synergizing reasoning and acting in language models, 2023.

675
676 Andy Zeng, Maria Attarian, Brian Ichter, Krzysztof Choromanski, Adrian Wong, Stefan Welker,
677 Federico Tombari, Aveek Purohit, Michael Ryoo, Vikas Sindhwani, et al. Socratic models: Com-
678 posing zero-shot multimodal reasoning with language. *arXiv preprint arXiv:2204.00598*, 2022.

679
680 Yuexiang Zhai, Hao Bai, Zipeng Lin, Jiayi Pan, Shengbang Tong, Yifei Zhou, Alane Suhr, Saining
681 Xie, Yann LeCun, Yi Ma, et al. Fine-tuning large vision-language models as decision-making
682 agents via reinforcement learning. *arXiv preprint arXiv:2405.10292*, 2024.

683
684 Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-exemplar
685 prompting with memory for computer control. In *The Twelfth International Conference on Learn-*
686 *ing Representations*, 2023.

687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702 A DERIVATION OF FORMULAS

703
704 We provide a simple derivation of Equation 4. During the RL phase with reward model, the object
705 of training is to maximize returns. Following prior works the optimization is formulated as:

$$706 \max_{\pi_\theta} \mathbb{E}_{s \sim D, a \sim \pi_\theta(a|s)} [Q(s, a)] - \beta \text{D}_{\text{KL}}[\pi_\theta \parallel \pi_{ref}] \quad (11)$$

707
708 which can be rewritten as:

$$\begin{aligned} 709 & \max_{\pi_\theta} \mathbb{E}_{s \sim D, a \sim \pi_\theta(a|s)} [Q(s, a)] - \beta \text{D}_{\text{KL}}[\pi_\theta \parallel \pi_{ref}] \\ 710 & = \max_{\pi_\theta} \mathbb{E}_{s \sim D, a \sim \pi_\theta(a|s)} [Q(s, a) - \beta \log \frac{\pi(a|s)}{\pi_{ref}(a|s)}] \\ 711 & = \min_{\pi_\theta} \mathbb{E}_{s \sim D, a \sim \pi_\theta(a|s)} [\log \frac{\pi(a|s)}{\pi_{ref}(a|s)} - \frac{1}{\beta} Q(s, a)] \\ 712 & = \min_{\pi_\theta} \mathbb{E}_{s \sim D, a \sim \pi_\theta(a|s)} [\log \frac{\pi(a|s)}{\pi_{ref}(a|s) \exp(\frac{1}{\beta} Q(s, a))}] \\ 713 & = \min_{\pi_\theta} \mathbb{E}_{s \sim D} [\text{D}_{\text{KL}}[\pi(a|s) \parallel \tilde{\pi}(a|s)]] \end{aligned}$$

714
715 where $\tilde{\pi}(a|s) = \pi_{ref}(a|s) \exp(\frac{1}{\beta} Q(s, a))$. KL-divergence is minimized at zero if and only if the
716 two distributions are identical. Therefore, in the case of the optimal solution we get:

$$717 \pi(a|s) = \tilde{\pi}(a|s) = \pi_{ref}(a|s) \exp(\frac{1}{\beta} Q(s, a))$$

718
719 A simple transformation yields:

$$720 Q(s, a) = \beta \log \frac{\pi(a|s)}{\pi_{ref}(a|s)} \quad (12)$$

721
722 We can know from Yang et al. (2024) that the Q-value form of Bradley-Terry preference distribution
723 can be expressed as:

$$724 p(\tau_1 > \tau_2 | a_i^t, s_i^t, a_i^{t-1}, \dots, s_i^0)_{i \in \{1, 2\}} = \frac{\exp(Q(s_1^t, a_1^t))}{\sum_{i \in \{1, 2\}} \exp(Q(s_i^t, a_i^t))} \quad (13)$$

725
726 Combining Eq. 12 and Eq. 13, replacing s_i^t with τ_i^{t-1} and a_i^t with (a_i^t, \mathcal{T}_i^t) , we derive the following
727 loss function:

$$728 L = -\mathbb{E}_\zeta \log \sigma \left(\beta \log \frac{\pi_\theta(a_1^t, \mathcal{T}_1^t | \tau_1^{t-1})}{\pi_{ref}(a_1^t, \mathcal{T}_1^t | \tau_1^{t-1})} - \beta \log \frac{\pi_\theta(a_2^t, \mathcal{T}_2^t | \tau_2^{t-1})}{\pi_{ref}(a_2^t, \mathcal{T}_2^t | \tau_2^{t-1})} \right) \quad (14)$$

729
730 which is similar to Eq. 4