

HINDSIGHT PLANNER: A CLOSED-LOOP FEW-SHOT PLANNER FOR EMBODIED INSTRUCTION FOLLOWING

Anonymous authors

Paper under double-blind review

ABSTRACT

This work focuses on building a task planner for Embodied Instruction Following (EIF) using Large Language Models (LLMs). Previous works typically train a planner to imitate expert trajectories, treating this as a supervised task. While these methods achieve competitive performance, they often lack sufficient robustness. When encountering a suboptimal action, the planner may encounter an out-of-distribution state, which can lead to task failure. In contrast, we frame the task as a Partially Observable Markov Decision Process (POMDP) and aim to develop a robust planner under a few-shot assumption. Thus, we propose a closed-loop planner with an adaptation module and a novel hindsight method, aiming to use as much information as possible to assist the planner. Our experiments on the ALFRED dataset indicate that our planner achieves competitive performance under a few-shot assumption. For the first time, our few-shot agent’s performance approaches and even surpasses that of the full-shot supervised agent.

1 INTRODUCTION

With the development of AI and robotics, many previous works have combined them to handle Embodied Instruction Following (EIF). Among them, the *Action Learning From Realistic Environments and Directives* (ALFRED) benchmark (Shridhar et al., 2020) is particularly challenging because it requires an agent to learn a long-horizon policy that maps egocentric images and language instructions into a sequence of actions. In each task, the agent will be given a natural instruction (e.g. “Put a heated mug down on a table”) and an egocentric visual observation at each step. The agent is required to output low-level actions (e.g. MoveAhead, RotateRight, etc.) based on the observation to complete the task. These tasks are usually challenging due to the sparse reward settings. For such a reason, many works have adopted a hierarchical structure to deal with it (Song et al., 2023; Min et al., 2021; Blukis et al., 2021; Kim et al., 2024). The high-level module decomposes the whole task into several sub-tasks, the low-level module outputs actions to finish each sub-task. Previously, sub-goal planners are trained on human-annotated dataset through supervised learning. However, they require large amounts of data and often lack robustness (Min et al., 2021; Blukis et al., 2021; Kim et al., 2024).

With recent advancements in Large Language Models (LLMs), many studies have explored using LLMs as sub-goal planners, utilizing their in-context learning abilities (Song et al., 2023; Shin et al., 2024; Ahn et al., 2022). Although these methods have achieved competitive performance under the few-shot assumption, a critical limitation is that these approaches all study the problem from a supervised learning perspective. They merely attempt to imitate the ground truth trajectories, which results in a lack of robustness within their agents. EIF benchmarks, on the other hand, require long-horizon planning ability. For example, the task “Put a warmed apple in the fridge” requires 12-step planning. Assuming that after applying in-context learning, the distribution of the agent’s output actions becomes closer to that of the Oracle, with an accuracy of 0.9, the overall accuracy of the entire planning task decreases to $0.9^{12} = 0.28$. Traditionally, a large amount of data is required to mitigate such an issue (Blukis et al., 2021; Kim et al., 2024). However, under the few-shot assumption, in-context learning methods rely heavily on the reasoning ability of pretrained LLMs (Brown et al., 2020; Dong et al., 2024). The hallucination problem of LLMs (Zhang et al., 2023) suggests that supervised methods through in-context learning are limited.

To address this issue, we approach the ALFRED task (Shridhar et al., 2020) as a Partially Observable Markov Decision Process (POMDP), where the planner makes decisions based on its current state.

Each task begins with a natural language description. At each step, the planner receives an egocentric RGB image and returns a high-level sub-goal. The planner can only receive reward signals (Success or Fail) at the end of the task. There are three major challenges in building a robust planner: (1) The sparse reward settings make it difficult for the planner to learn and make accurate decisions. (2) The planner can only receive an egocentric picture and cannot detect the whole state. (3) Under the few-shot assumption, the planner cannot obtain enough information from trajectories.

For the first problem, we adopt an actor-critic framework (Liu et al., 2024) which consists of two actors, one critic, and one generator. At each step, the planner receives a new state and performs a tree search with the actors and generator to plan future trajectories, rather than directly outputting a sub-goal. The critic is then used to select the best rollout and return its initial action. Thus, the planner can optimize the output over the long horizon to address the issue of sparse reward. For the second difficulty, we design an adaptation module instantiated by LLMs. Upon receiving an egocentric image, the adaptation module aims to predict the invisible latent PDDL variables of the task, which could help the planner better understand the environment. For the third challenge, we propose a novel hindsight method. It collects suboptimal trajectories from the agent in the training environment and relabels them to complete the task. This approach provides the planner with additional information. During the deployment phase, the relabeled trajectories can guide the planner in adjusting its policy when incorrect actions are proposed and executed.

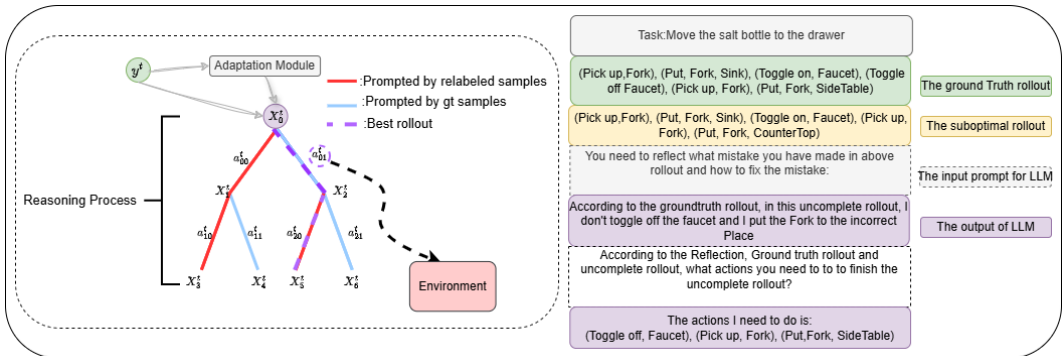


Figure 1: **Left:** The illustration of the Hindsight Planner: at each time step t , the planner receives a partial observation y^t from the environment. The adaptation module estimates the latent variable and concatenates it with y^t to produce the complete state. $\text{ACTOR}_{\text{hind}}$ and ACTOR_{gt} are prompted with different samples and make decisions. The Critic is utilized to evaluate the actions. The best rollout $(x_t, a_t^*, x_{t+1}^*, a_{t+1}^* \dots)$ is selected, and a_t^* is returned. **Right:** An example of the relabeling process for the $\text{ACTOR}_{\text{hind}}$: after collecting a suboptimal rollout, the LLM is prompted to generate a reflection on the previously taken actions. Following this reflection, the LLM is then prompted to complete the suboptimal rollout.

In summary, our contributions are threefold:

- (1) We study ALFRED (Shridhar et al., 2020) from a POMDP perspective for the first time and propose a closed-loop actor-critic planner to solve it.
- (2) We propose a novel hindsight prompting method and demonstrate that our method is theoretically superior to previous approaches.
- (3) Experiments on ALFRED (Shridhar et al., 2020) show that our method achieves state-of-the-art performance under few-shot assumptions. Specifically, the success rates for the “Test Seen” and “Test Unseen” splits are 25.51 and 18.77, respectively, representing a 60% and 39% improvement over the previous baseline.

2 RELATED WORK

2.1 LARGE LANGUAGE MODEL (LLM) AND IN-CONTEXT LEARNING (ICL)

Large language models (LLMs) have shown incredible reasoning ability (Vaswani et al., 2023; Wei et al., 2022; Touvron et al., 2023; OpenAI et al., 2024) across a wide range of tasks. A crucial way to enhance this reasoning ability is through in-context learning (ICL) (Brown et al., 2020; Dong et al., 2024), which allows LLMs to solve complex tasks with only a few samples. Furthermore,

108 this approach removes the need for fine-tuning, which can be time-consuming and computationally
 109 expensive. To utilize the ICL ability better, many studies propose certain frameworks aimed at
 110 enhancing the reasoning capabilities of LLMs (Yao et al., 2023; Wei et al., 2023; Yao et al., 2024).
 111 Among them, Liu et al. (2024) proposes a novel perspective by bridging RL and LLM, which inspires
 112 us to study ICL from an RL aspect. Xie et al. (2022) interprets ICL as Implicit Bayesian Inference,
 113 while Dai et al. (2023) believes that ICL is performing implicit Gradient Descent. All of these imply
 114 the importance of the content in ICL, an area that remains relatively understudied. To this end, we
 115 propose Hindsight Planner as an exploration.

117 2.2 ADAPTATION MODULE IN POMDP

119 In a Partially Observable Markov Decision Process (POMDP), planners are presented with observ-
 120 able states, while the latent states are invisible to the planner. Making decisions with incomplete
 121 information is challenging; therefore, a component to map the observable state into the latent space is
 122 crucial (Lee et al., 2023). Adaptation modules have been proven effective in legged robots (Kumar
 123 et al., 2021; Zhou et al., 2019; Peng et al., 2020). These modules aim to bridge the gap between
 124 the simulator and the real world. They are often trained to predict crucial information that a robot
 125 can sense in the simulator but not through its sensors in the actual world, such as surface friction or
 126 payload of the robot. The base policy then makes decisions based on the observed information and
 127 the invisible latent information predicted by adaptation modules. Inspired by this, we propose an
 128 adaptation model that maps the visible object list to the latent, invisible Planning Domain Definition
 129 Language (PDDL) (Chapman, 1987) of ALFRED (Shridhar et al., 2020).

130 Previous work such as Min et al. (2021), trains a BERT (Devlin et al., 2019) to predict the PDDL
 131 arguments and decompose high-level instructions into templated sub-goals. However, our approach
 132 differs from these in two aspects: (1) Previous works predict the arguments at the beginning of a
 133 task, which is equivalent to predicting the latent variables based on the initial observed state. In
 134 contrast, our method predicts the latent arguments at each time before reasoning, allowing predictions
 135 to be adjusted through exploration, which makes our planner more robust. (2) We do not apply the
 136 templated approach directly. The adaptation module is used to reveal the latent information for the
 137 planner and assist the planner in making better decisions. Experiments show that our method achieves
 138 competitive performance even without the assistance of the adaptation model, as demonstrated in
 139 Table 4.

140 2.3 HINDSIGHT IN LLMs

142 Hindsight algorithms (Andrychowicz et al., 2018; Li et al., 2020; Pong et al., 2020) are widely
 143 adopted in the reinforcement learning (RL) area. Generally, the hindsight method aims to reveal
 144 future information after collecting a trajectory and relabel the trajectory to make it more informative
 145 during training process (Furuta et al., 2022; Andrychowicz et al., 2018). Furuta et al. (2022) applies
 146 the hindsight method in training a Transformer model and achieves competitive performance on
 147 several baselines. However, training a model from scratch usually requires a large amount of
 148 data. In contrast, in-context learning, leveraging the reasoning ability of LLMs, allows an agent to
 149 complete complex tasks with only a few samples. Dai et al. (2023) has shown that ICL executes
 150 an implicit parameter update. As a result, we utilize ICL in our proposed method. Intuitively, we
 151 hope hindsight prompts can provide guidance when an out-of-distribution state is encountered. For
 152 example, “Wash a pan and put it away” requires the agent to wash a *Pan* and put it on the *DiningTable*.
 153 The trajectory from a planner could be: {(PickupObject, Pan), (PutObject, Sink), (ToggleObjectOn,
 154 Faucet), (PickupObject, Pan), (PutObject, CoffeeMachine)}. Note that in this example, the agent
 155 fails to place the pan in the correct location, does not turn off the faucet, and thus the trajectory
 156 from the planner is suboptimal. Our hindsight method proposes a novel relabeling process that
 157 appends actions to the suboptimal trajectory, aiming to complete the task. In the above example, the
 158 corrected trajectory should be: {(PickupObject, Pan), (PutObject, Sink), (ToggleObjectOn, Faucet),
 159 (PickupObject, Pan), (PutObject, CoffeeMachine), (ToggleObjectOff, Faucet), (PickupObject, Pan),
 160 (PutObject, DiningTable)}. This approach enables us to guide the planner in addressing unknown
 161 states resulting from incorrect actions. Consequently, during the deployment phase, when the planner
 encounters a similar state, it can learn from suboptimal trajectories and subsequently take correct
 actions to correct previous mistakes.

We also analyze our method in comparison to previous hindsight methods (Andrychowicz et al., 2018; Ghosh et al., 2019) following the framework proposed by Furuta et al. (2022). We demonstrate that while previous methods are effective, they alter the distribution of a crucial variable in multi-task RL problems. In contrast, our method optimize the same objective while maintaining the distribution. The detailed discussion can be found in Section 4.3.

3 PRELIMINARIES

3.1 DEFINITION IN POMDP

In a POMDP \mathcal{M} , consider an action space \mathcal{A} , latent state space \mathcal{X} , observation space \mathcal{Y} , transition probability function $p(x'|x, a)$, emission function $o(y|x)$, reward function $r(x, a)$ and discount factor $\gamma \in [0, 1)$. The policy $\pi_\theta(\cdot|y)$ maps the latent state space to the action space, where θ represents its parameter. The goal of RL is to train a policy such that

$$\pi_\theta = \arg \max_{\pi} \frac{1}{1 - \gamma} \mathbb{E}_{x \sim \rho^\pi(x), y \sim o(\cdot|x), a \sim \pi(\cdot|y)} [r(x, a)], \quad (3.1)$$

where $\rho^\pi(x) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t p_t^\pi(x_t = x)$,

$$p_t^\pi(x) = \int_{x_{0:t}, y_{0:t-1}, a_{0:t-1}} \prod_{i=1}^t p(x_i|x_{i-1}, a_{i-1}) o(y_{i-1}|x_{i-1}) \pi(a_{i-1}|y_{i-1}).$$

Given a parameterized reward function $r_z(x, a)$, $z \in \mathcal{Z}$ is a variable that indicates the goal for the agent. A conditional policy $\pi(a|y, z)$ is trained with different values of z where $z \sim p(z)$. The goal in Equation (3.1) becomes

$$\pi_\theta = \arg \max_{\pi} \frac{1}{1 - \gamma} \mathbb{E}_{z \sim p(z), x \sim \rho_z^\pi(x), y \sim o_z(\cdot|x), a \sim \pi(\cdot|y, z)} [r_z(x, a)]. \quad (3.2)$$

Equation (3.2) can be considered as the multi-task RL objective to optimize, which is the core of EIF.

3.2 INFORMATION MATCHING

Following Furuta et al. (2022), we define the information matching (IM) problem as training a policy π_θ that satisfies

$$\pi_\theta = \arg \min_{\pi_\theta} \mathbb{E}_{z \sim p(z), \tau \sim \rho_z^\pi(\tau)} [D(I(\tau), z)], \quad (3.3)$$

where $I(\tau)$ is *information statistic* that can be any function that captures the desired information from a partially observed trajectory $\tau_t = \{y_0, a_0, y_1, y_1, \dots, y_t\}$ and D is a divergence measure such as Kullback-Leibler (KL) divergence or some f -divergences. Competitive results have been achieved with this optimization objective (Lee et al., 2020; Hazan et al., 2019).

Furuta et al. (2022) demonstrates that previous hindsight methods (Andrychowicz et al., 2018; Eysenbach et al., 2020; Guo et al., 2021) utilize various *information statistics* and minimize the divergence $D = 0$ by setting $\hat{z} = I(\tau)$. This allows trajectories to be better used to train a policy $\pi(a|s, z)$. For instance, in HER (Andrychowicz et al., 2018), an MDP trajectory $\tau_t^s = \{s_0, a_0, s_1, \dots, s_t\}$ is collected. The information statistic is set as the final state of the agent, where $I(\tau_t^s) = s_t$, and the relabeling process in HER is equivalent to setting $\hat{z} = I(\tau_t^s)$.

4 HINDSIGHT PLANNER

4.1 OVERVIEW

The Hindsight Planner outputs a sub-goal based on the observed objects and natural language instructions. During the collection phase, suboptimal trajectories are collected, and we apply our

Algorithm 1 Hindsight Planner

```

1: Input: An LLM-planner  $LLM-PL$ , an adaptation module  $Adapter$  and the task instruction  $I$ .
2: Set: Observed Objects  $O \leftarrow \emptyset$ , the sub-goal history  $G \leftarrow \emptyset$ , the current sub-goal  $S \leftarrow \emptyset$ , the time step  $t \leftarrow 0$  and sub-goal index  $k \leftarrow 0$ .
3: Get sample pool  $\mathcal{D}$  and initialize  $Actor_{\theta}$ ,  $Critic$ ,  $Adapter$  from  $\mathcal{D}$ , for any  $\theta \in \{gt, hind\}$  (e.g. Algorithm 3 in Appendix A).
4: while Not Finished do
5:   Get PDDL arguments  $P \leftarrow Adapter(I, O)$ .
6:   Plan and get sub-goal  $S_k \leftarrow LLM-PL(Actor_{gt}, Actor_{hind}, Critic, P, I, O, G)$  (e.g. Algorithm 2 in Appendix A).
7:   Set  $S_k$  as sub-goal for  $LOW-PL$ .
8:   while  $S_k$  not Finished and not Failed do
9:     Invoke  $LOW-PL$  to plan and execute  $a_t$  and update  $O$ .
10:    Set  $t \leftarrow t + 1$ 
11:    if  $S_k$  Finished then
12:      Append  $S_k$  to  $G$ .
13:      Set  $k \leftarrow k + 1$ .
14:    end if
15:  end while
16: end while

```

hindsight method to generate \mathcal{D}_{hind} . The complete dataset $\mathcal{D} = \mathcal{D}_{hind} \cup \mathcal{D}_{gt}$, where \mathcal{D}_{gt} is constructed from training data. In the deployment phase, we initiate hindsight actor $Actor_{hind}$, ground truth actor $Actor_{gt}$, and $Critic$ from \mathcal{D} .

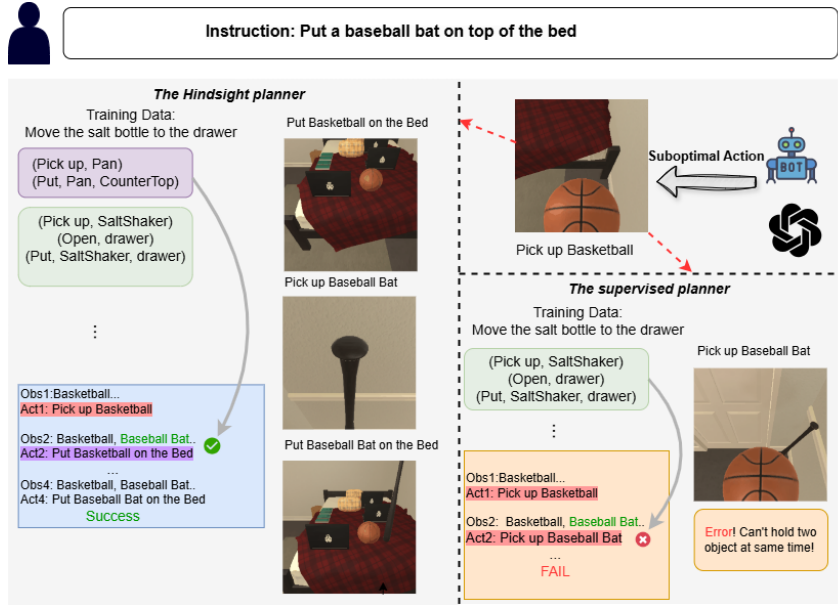


Figure 2: A comparison of Hindsight Planner and previous supervised methods when taking a suboptimal action. The agent initially picks up the incorrect object (“Basketball”). In the supervised method, the planner fails to handle this situation, which leads to task failure. In contrast, the Hindsight Planner can adjust after the incorrect action and successfully complete the task.

At time step t , the planner receives an observed object list y_t from observation functions (Blukis et al., 2021). We then apply the Adaptation module to predict the latent PDDL arguments P based on y_t . The whole state x_t is constructed by y_t and P . With x_t , we invoke the actor-critic task planner $LLM-PL$ to generate a future trajectory over a long horizon and return the sub-goal S_k . To ensure the output from the planner meets the requirements, a frozen BERT (Devlin et al., 2019) is used to map the output to the legal space. The proposed sub-goal will be executed by a low-level controller $LOW-PL$ (Blukis et al., 2021). When a sub-goal is completed or fails, the planner reinvokes the reasoning process to replan another future trajectory from the new state. The complete algorithm is presented in Algorithm 1, and Figure 3 provides an example of the entire process.

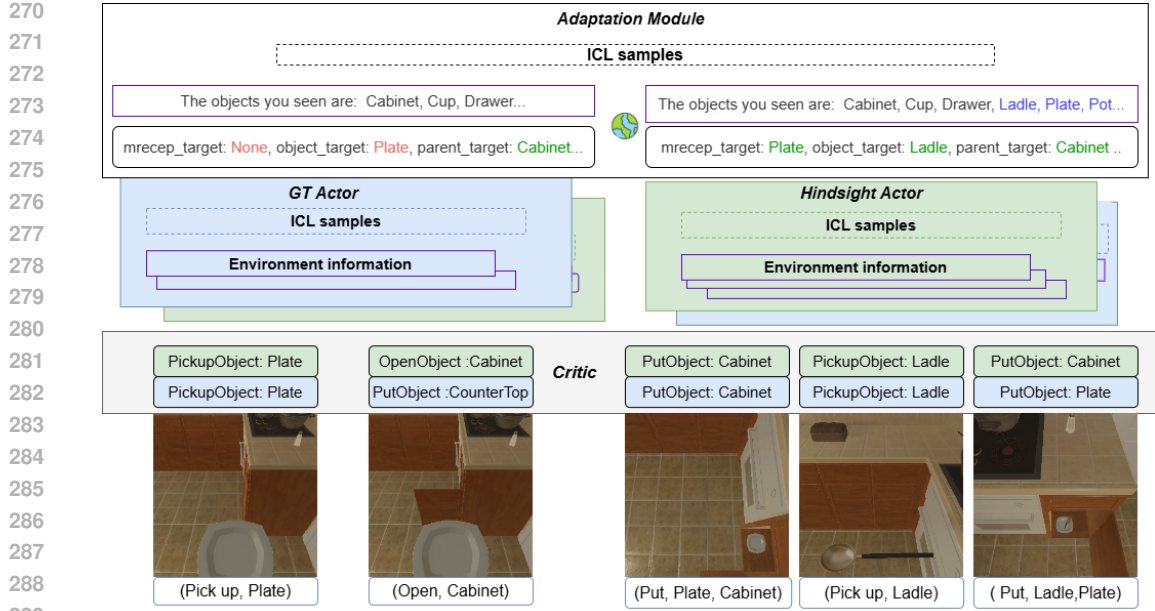


Figure 3: The entire process of the Hindsight Planner is as follows: At the start of the task, which is to “Place a plate with a ladle on it in a cabinet,” the Adapter mistakenly identifies the task as picking up a plate and placing it into a cabinet. $\text{Actor}_{\text{hind}}$ and Actor_{gt} make decisions separately. Critic then selects the best action as its output. Upon further exploration, the agent detects more objects, and the Adapter adjusts its output, recognizing the task as stacking a ladle onto a plate and then placing them into a cabinet. The Actors and Critic subsequently make decisions based on the revised predictions.

4.2 PROMPT DESIGN

All components follow a similar design. The prompt begins with an intuitive explanation of the task and a role description of the LLM. A frozen BERT is then used as a kNN retriever, encoding the task description and selecting K examples with the closest Euclidean distance from the sample pool as in-context samples (Song et al., 2023). Intuitively, the planner would make similar suboptimal actions in similar tasks. For instance, if in an in-context sample “Place two spray bottles into the cabinet,” the planner fails to open the cabinet when putting the second spray bottle into it. In the current task “Putting two candles in a cabinet”, the planner would know to avoid a similar mistake. The detailed prompts for each process can be viewed in Appendix B.

4.3 HINDSIGHT METHOD

In Section 3, we gain a coherent framework to describe previous hindsight methods. However, we find that such methods can lead to the policy π being suboptimal, particularly when the number of samples is insufficient. To illustrate this better, we consider the optimization objective in Equation (3.2). It aims to learn a policy under different values of z where $z \sim p(z)$. During the collection phase, the agent’s trajectory is usually suboptimal and random. Assume the distribution of $I(\tau) \sim q$. The training objective after relabeling is to train a policy $\hat{\pi}$ satisfies that

$$\hat{\pi} = \arg \max_{\pi} \frac{1}{1-\gamma} \mathbb{E}_{z \sim q(z), x \sim \rho_z^{\pi}(x), y \sim o(\cdot|x), a \sim \pi(\cdot|y,z)} [r_z(x, a)]. \quad (4.1)$$

Define the π^* as the oracle. It is easy to see that

$$\begin{aligned} & \frac{1}{1-\gamma} \mathbb{E}_{z \sim p(z), x \sim \rho_z^{\pi^*}(x), y \sim o(\cdot|x), a \sim \pi^*(\cdot|x,z)} [r_z(x, a)] \\ & < \frac{1}{1-\gamma} \mathbb{E}_{z \sim p(z), x \sim \rho_z^{\hat{\pi}}(x), y \sim o(\cdot|x), a \sim \hat{\pi}(\cdot|x,z)} [r_z(x, a)], \end{aligned} \quad (4.2)$$

as the distribution of z is shifted from p to q .

Based on such discovery, we propose a new method of hindsight. Assume that τ^* is the ground truth rollout from the oracle π^* , we can rewrite $z = I(\gamma^*)$, Equation (3.3) then becomes

$$\min_{\pi} \mathbb{E}_{z \sim p(z), \tau \sim \rho_z^{\pi}(\tau)} [D(I(\tau), I(\tau^*))]. \quad (4.3)$$

Our method utilizes LLMs to relabel $\hat{\tau} = \tau_T + \{a_T, y_{T+1}, a_{T+2}, \dots\}$ in such a way that $I(\hat{\tau}) = I(\tau^*) = z$. Thus, we minimize the divergence in Equation (4.3) while keeping the distribution of z unshifted. Intuitively, Equation (4.2) shows that relabeling z alters the distribution of tasks that are truly relevant to our daily lives. This is especially crucial in the reasoning process of EIF.

In practice, our hindsight method consists of two main parts: the collection phase and the deployment phase. During the collection phase, the planner executes tasks and retrieves K examples from a small set of ground truth samples. At each task, the planner generates a possibly suboptimal trajectory τ and relabels them. The algorithm is summarized in Algorithm 3 of Appendix A. During the deployment phase, the `Actorgt` is prompted with ground truth samples while the `Actorhind` and the `Critic` are prompted with relabeled samples. Intuitively, we hope that the `Actorgt` can provide the correct action to complete the task along the shortest path. However, when an incorrect action—which is often unavoidable—is executed, the `Actorhind` and the `Critic` should be able to correct it. The relabeling process utilizes the reasoning ability of LLMs to fit suboptimal trajectories into correct rollouts. The CoT (Wei et al., 2023) method is utilized in the relabeling process. We first prompt the LLM to generate a *Think* about the suboptimal rollout and then prompt it to complete the suboptimal rollout based on the *Think*. A comparison of the hindsight method with the supervised methods is shown in Figure 2, while the right half of Figure 1 illustrates an example of the relabeling process.

4.4 ADAPTATION MODULE

In a POMDP, the adaptation module is used to predict the latent variables from the observed environment y_t (Lee et al., 2023; Kumar et al., 2021) and construct the whole state $x_t = (\text{Adapter}(y_t), y_t)$. In practice, we utilize an LLM as the adaptation module and set PDDL arguments as the prediction target for it. The input prompt for the adaptation module begins with an intuitive explanation of ALFRED, followed by several in-context samples. At the end of the prompt is the current task and the object list. At each step, the object list is updated as the agent explores the environment.

The output from the adaptation module varies depending on the task description. Inspired by PDDL (Chapman, 1987; Silver et al., 2023) of ALFRED, the adaptation module needs to predict the following arguments at each step: (1) *object_target*: The specific object to be interacted with during the task. (2) *parent_target*: The final place for the object in the task. (3) *mrecep_target*: The container or vessel necessary for the task. (4) *toggle_target*: The device that needs to be toggled in the task. (5) *object_state*: Indicates whether the target object needs to be cleaned, heated, or cooled. (6) *object_sliced*: Determines if the object must be sliced. (7) *two_object*: Specifies whether the task involves handling and placing two objects. The adaptation module predicts these arguments at each time before reasoning. Then, the arguments are processed into a specific format to assist the task planner to sense the environment better.

4.5 TASK PLANNER

We adopt an actor-critic planner (Liu et al., 2024). At each time step t , the planner receives x_t from the environment and the adaptation module. We initiate two Actors: `Actorgt` and `Actorhind`, with different samples from the sample pool \mathcal{D} . For each state, we prompt each Actor to generate $\frac{W}{2}$ actions. The `Critic` then selects the top B actions. A generator ψ generates the next state based on each action. In this way, we map `Actors` and `Critic` to B future trajectories and select the best future trajectory (x_t, a_t^*, \dots) through `Critic`. a_t^* is then returned as the sub-goal for `Low-PL`. The left half of Figure 1 shows the reasoning process of the planner.

5 EXPERIMENT

5.1 SETUPS

We validate our framework using the ALFRED benchmark (Shridhar et al., 2020). This benchmark assesses the agent’s capability to execute a series of actions for long-horizon household tasks based

Model	n-shot	Test Seen		Test Unseen	
		SR	GC	SR	GC
HiTUT (Zhang & Chai, 2021)	full	13.63	21.11	11.12	17.89
HLSM (Blukis et al., 2021)	full	25.11	<u>35.79</u>	<u>20.27</u>	<u>27.24</u>
FILM (Min et al., 2021)	full	<u>28.83</u>	39.55	27.80	38.52
MCR-Agent (Bhambri et al., 2024)	full	30.13	-	17.04	-
FILM (low inst.) (Min et al., 2021)	few	0.00	4.23	0.20	6.71
LLM-Planner (Song et al., 2023)	few	15.33	24.57	13.41	22.8
LLM-Planner (low inst.) (Song et al., 2023)	few	<u>18.80</u>	<u>26.77</u>	<u>16.42</u>	<u>23.37</u>
Socratic-Planner (Shin et al., 2024)	few	13.24	21.51	10.66	19.53
Hindsight planner (ours)	few	25.51	34.74	18.77	28.29

Table 1: **Comparison with the state-of-the-art methods on SR and GC in the test set.** Bold symbols in numbers denote the highest accuracy, while underlined symbols indicate the runner-up for each experiment setting. “low inst.” refers to the use of step-by-step instructions.

on natural language task descriptions and egocentric vision. The ALFRED dataset consists of 25k annotations, 108 distinct objects, 7 types of tasks, and 120 scenes. The dataset is divided into training, validation, and testing splits. The validation and test splits contain “seen” subsets, which are part of the training fold, and “unseen” subsets, which are distinct from it. The evaluation is based on Success Rate (SR) and Goal Condition (GC). Given the inherent noise in natural language instructions and the complexities of long-horizon task planning, the ALFRED benchmark presents significant challenges for embodied agents in formulating robust and precise plans.

Similar to previous work (Song et al., 2023; Shin et al., 2024), we only utilize a few examples from the 21k training set annotations. For each of the 7 task types, we randomly select 20 trajectories as the initial sample pool. At the collection phase, we run our planner on the 140 trajectories and collect sub-optimal trajectories. During collection, the same task is not included as in-context samples.

We then give a detailed discussion of the relabeling process. Directly applying the task description from ALFRED may lead to unsatisfactory results, as the task description is often vague. For example, the task “Put a chilled potato on the small black table” requires the planner to put the *potato* on a *SideTable*. If the task description is applied directly, LLMs might focus incorrectly on the *Black Table* and return an incorrect action “PutObject BlackTable”. If the task description is not included in the prompt, it could lead LLMs to imitate the ground truth trajectory. However, planners usually have multiple ways to complete a certain task. For instance, in a task requiring the planner to slice an apple, after slicing the apple, the planner could put the *Knife* on the *DiningTable* or *CounterTop*. To address this issue, we relabel the task based on the latent PDDL arguments. The task description “Put a chilled potato on the small black table” becomes “Pick up one cooled potato and put it on the *SideTable*”. This approach helps clarify the task for the planner and reduces the ambiguity in instructions.

For the kNN retriever, we use a frozen BERT from Wolf et al. (2020). We employ GPT-4 Turbo as the target LLM and set temperature to 0. For the `Adapter`, 5 in-context examples are retrieved from the sample pool through the kNN retriever. For the `Actors` and `Critic` modules, 2 in-context examples are retrieved. The task planner uses beam search with a depth and width of 2. To preserve the few-shot assumption and ensure a fair comparison, we directly adopt the pretrained modules for navigation, perception, and low-level control from HLSM (Blukis et al., 2021).

5.2 MAIN RESULTS

We initially compare our method to other few-shot methods, as shown in Table 1. It is evident that our method achieves a 10.18 and 5.36 higher success rate in “Test Seen” and “Test Unseen” categories, respectively, compared to the previous state-of-the-art method (LLM-Planner) that uses high-level instructions only. Moreover, even when compared to methods utilizing low-level, step-by-step instructions, our method still demonstrates superior performance.

We also compare our method to the other approaches under the same low-level controller (Blukis et al., 2021) in Table 2. The results indicate that our method not only significantly outperforms

Model	n-shot	Valid Seen		Valid Unseen		Test Seen		Test Unseen	
		SR	GC	SR	GC	SR	GC	SR	GC
HLSM (Blukis et al., 2021)	full	29.63	38.74	18.28	31.24	<u>25.11</u>	35.79	20.27	<u>27.24</u>
LLM-Planner (Song et al., 2023)	few	13.53	28.28	12.92	25.35	15.33	24.57	13.41	22.8
Socratic-Planner (Shin et al., 2024)	few	14.88	25.47	13.40	24.91	13.24	21.51	10.66	19.53
Hindsight planner (ours)	few	<u>25.61</u>	<u>34.95</u>	19.00	<u>29.90</u>	25.51	<u>34.74</u>	<u>18.77</u>	28.29

Table 2: **Comparison with the same lower-controller.** Bold symbols in numbers denote the highest accuracy, while underlined symbols indicate the runner-up for each experiment setting.

Task Type	Examine	Pick	Clean	Stack	Pick Two	Heat	Cool
Base Method	40.42	50	15.18	9.56	30.65	7.48	21.43
w.o. hindsight prompt	39.36	49.29	16.96	7.82	29.84	7.47	10.31
w.o. adaptation module	35.1	47.1	8.93	6.09	32.25	9.34	18.26

Table 3: Ablation study on the success rate of different type of tasks in “Valid Seen” split.

previous few-shot LLM planners but also, for the first time, a few-shot LLM method (with around 100 examples) nearly matches and even surpasses (SR in “Valid Unseen”, “Test Seen”, and GC in “Test Unseen”) fully supervised (around 21k samples) methods.

5.3 ABLATION STUDY

We conduct ablation studies to understand the effectiveness of the components in our framework. First, we ablate the adaptation module `Adapter`, which requires the planner to make decisions based solely on the partially observed information. The results show that this causes a drop of -2.44 and -4.01 in the success rates for the “Valid Seen” and “Valid Unseen” splits.

Then, we remove the hindsight prompts. For a fair comparison, the original planner requires both `Actorgt` and `Actorhind` to generate one action per state. We also ablate by prompting `Actorgt` to output two actions for each state. Table 4 shows that the success rates drop by -2.08 and -2.68 in the “Valid Seen” and “Valid Unseen” splits.

For a more comprehensive analysis, we report the success rates for each task type in the “Valid Seen” split, as shown in Table 3. Concurrently, we also present the average sub-goal lengths in the right table. This analysis reveals that hindsight prompting is especially crucial in relatively long-horizon tasks, such as “Cool Object” and “Heat Object”. This is likely because, in long-horizon tasks, planners will output suboptimal actions with a higher probability. On the other hand, the adaptation module can assist the planner in better sensing the environment, leading to a general improvement across nearly all areas.

Model	Valid Seen		Valid Unseen	
	SR	GC	SR	GC
W.O. Adaptation module	23.17	33.28	14.99	27.36
W.O. Hindsight Prompt	23.53	32.76	16.32	28.06
Base Method	25.61	34.95	19.00	29.90

Table 4: Ablation on “Valid Seen”, “Valid Unseen” splits.

Task Type	Avg. Sub-Goal Len.
Examine	2.07
Pick	2.48
Pick Two	5.70
Stack	5.63
Clean	7.25
Cool	10.36
Heat	12.78

6 CONCLUSION

This paper explores an effective few-shot framework for Embodied Instruction Following. We approach the task as a POMDP and design a closed-loop Hindsight Planner equipped with an adaptation module to enhance the agent’s environmental sensing capabilities. Compared to previous open-loop, supervised methods, our approach is more robust and performs better. Furthermore, the planner incorporates a novel hindsight method that enables it to learn from suboptimal trajectories. We hope our work inspires future research in this area.

REFERENCES

- 486
487
488 Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea
489 Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine
490 Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally
491 Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee,
492 Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka
493 Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander
494 Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and
495 Andy Zeng. Do as i can, not as i say: Grounding language in robotic affordances, 2022. URL
496 <https://arxiv.org/abs/2204.01691>.
- 497 Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob
498 McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay, 2018.
499 URL <https://arxiv.org/abs/1707.01495>.
- 500 Suvaansh Bhambri, Byeonghwi Kim, and Jonghyun Choi. Multi-level compositional reasoning for
501 interactive instruction following, 2024. URL <https://arxiv.org/abs/2308.09387>.
- 502 Valts Blukis, Chris Paxton, Dieter Fox, Animesh Garg, and Yoav Artzi. A persistent spatial se-
503 mantic representation for high-level natural language instruction execution. *Cornell University -*
504 *arXiv,Cornell University - arXiv*, Jul 2021.
- 505 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
506 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
507 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler,
508 Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott
509 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya
510 Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.
- 511 David Chapman. Planning for conjunctive goals. *Artif. Intell.*, 32:333–377, 1987. URL <https://api.semanticscholar.org/CorpusID:1525549>.
- 512
513
514
515
516 Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. Why can gpt
517 learn in-context? language models implicitly perform gradient descent as meta-optimizers, 2023.
518 URL <https://arxiv.org/abs/2212.10559>.
- 519 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
520 bidirectional transformers for language understanding, 2019. URL <https://arxiv.org/abs/1810.04805>.
- 521
522 Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu,
523 Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. A survey on in-context learning,
524 2024. URL <https://arxiv.org/abs/2301.00234>.
- 525 Benjamin Eysenbach, Xinyang Geng, Sergey Levine, and Ruslan Salakhutdinov. Rewriting history
526 with inverse rl: Hindsight inference for policy improvement, 2020. URL <https://arxiv.org/abs/2002.11089>.
- 527
528
529 Hiroki Furuta, Yutaka Matsuo, and Shixiang Shane Gu. Generalized decision transformer for offline
530 hindsight information matching, 2022. URL <https://arxiv.org/abs/2111.10364>.
- 531
532 Dibya Ghosh, Abhishek Gupta, Justin Fu, Ashwin Reddy, Coline Devin, Benjamin Eysenbach, and
533 Sergey Levine. Learning to reach goals without reinforcement learning. *ArXiv*, abs/1912.06088,
534 2019.
- 535 Jiaming Guo, Rui Zhang, Xishan Zhang, Shaohui Peng, Qi Yi, Zidong Du, Xing Hu, Qi Guo, and
536 Yunji Chen. Hindsight value function for variance reduction in stochastic dynamic environment,
537 2021. URL <https://arxiv.org/abs/2107.12216>.
- 538
539 Elad Hazan, Sham M. Kakade, Karan Singh, and Abby Van Soest. Provably efficient maximum
entropy exploration, 2019. URL <https://arxiv.org/abs/1812.02690>.

- 540 Byeonghwi Kim, Jinyeon Kim, Yuyeong Kim, Cheolhong Min, and Jonghyun Choi. Context-aware
541 planning and environment-aware memory for instruction following embodied agents, 2024. URL
542 <https://arxiv.org/abs/2308.07241>.
- 543 Ashish Kumar, Zipeng Fu, Deepak Pathak, and Jitendra Malik. Rma: Rapid motor adaptation for
544 legged robots, 2021. URL <https://arxiv.org/abs/2107.04034>.
- 545 Jonathan N. Lee, Alekh Agarwal, Christoph Dann, and Tong Zhang. Learning in pomdps is sample-
546 efficient with hindsight observability, 2023. URL <https://arxiv.org/abs/2301.13857>.
- 547 Lisa Lee, Benjamin Eysenbach, Emilio Parisotto, Eric Xing, Sergey Levine, and Ruslan Salakhutdinov.
548 Efficient exploration via state marginal matching, 2020. URL <https://arxiv.org/abs/1906.05274>.
- 549 Alexander C. Li, Lerrel Pinto, and Pieter Abbeel. Generalized hindsight for reinforcement learning,
550 2020. URL <https://arxiv.org/abs/2002.11708>.
- 551 Zhihan Liu, Hao Hu, Shenao Zhang, Hongyi Guo, Shuqi Ke, Boyi Liu, and Zhaoran Wang. Reason
552 for future, act for now: A principled framework for autonomous llm agents with provable sample
553 efficiency, 2024. URL <https://arxiv.org/abs/2309.17382>.
- 554 SoYeon Min, DevendraSingh Chaplot, Pradeep Ravikumar, Yonatan Bisk, and Ruslan Salakhutdinov.
555 Film: Following instructions in language with modular methods. *Learning, Learning*, Oct 2021.
- 556 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
557 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor
558 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian,
559 Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny
560 Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks,
561 Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea
562 Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen,
563 Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung,
564 Dave Cummings, Jeremiah Carrier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch,
565 Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty
566 Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte,
567 Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel
568 Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua
569 Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike
570 Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon
571 Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne
572 Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo
573 Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar,
574 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik
575 Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich,
576 Aris Konstantinidis, Kyle Kopic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy
577 Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie
578 Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini,
579 Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne,
580 Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David
581 Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie
582 Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély,
583 Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo
584 Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano,
585 Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng,
586 Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto,
587 Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power,
588 Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis
589 Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted
590 Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel
591 Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon
592 Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,
593

- 594 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie
595 Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,
596 Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun
597 Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang,
598 Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian
599 Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren
600 Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming
601 Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao
602 Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL
603 <https://arxiv.org/abs/2303.08774>.
- 604 Xue Bin Peng, Erwin Coumans, Tingnan Zhang, Tsang-Wei Lee, Jie Tan, and Sergey Levine.
605 Learning agile robotic locomotion skills by imitating animals, 2020. URL <https://arxiv.org/abs/2004.00784>.
- 606
607 Vitchyr Pong, Shixiang Gu, Murtaza Dalal, and Sergey Levine. Temporal difference models: Model-
608 free deep rl for model-based control, 2020. URL <https://arxiv.org/abs/1802.09081>.
- 609
610 Suyeon Shin, Sujin jeon, Junghyun Kim, Gi-Cheon Kang, and Byoung-Tak Zhang. Socratic planner:
611 Inquiry-based zero-shot planning for embodied instruction following, 2024. URL <https://arxiv.org/abs/2404.15190>.
- 612
613 Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi,
614 Luke Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions
615 for everyday tasks. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
616 *(CVPR)*, Jun 2020. doi: 10.1109/cvpr42600.2020.01075. URL <http://dx.doi.org/10.1109/cvpr42600.2020.01075>.
- 617
618 Tom Silver, Soham Dan, Kavitha Srinivas, Joshua B. Tenenbaum, Leslie Pack Kaelbling, and Michael
619 Katz. Generalized planning in pddl domains with pretrained large language models, 2023. URL
620 <https://arxiv.org/abs/2305.11014>.
- 621
622 Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M. Sadler, Wei-Lun Chao, and Yu Su.
623 Llm-planner: Few-shot grounded planning for embodied agents with large language models, 2023.
624 URL <https://arxiv.org/abs/2212.04088>.
- 625
626 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
627 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand
628 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language
629 models, 2023. URL <https://arxiv.org/abs/2302.13971>.
- 630
631 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz
632 Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL <https://arxiv.org/abs/1706.03762>.
- 633
634 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
635 Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals,
636 Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models, 2022.
URL <https://arxiv.org/abs/2206.07682>.
- 637
638 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le,
639 and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.
URL <https://arxiv.org/abs/2201.11903>.
- 640
641 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
642 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von
643 Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama
644 Drame, Quentin Lhoest, and Alexander M. Rush. Huggingface’s transformers: State-of-the-art
645 natural language processing, 2020. URL <https://arxiv.org/abs/1910.03771>.
- 646
647 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-
context learning as implicit bayesian inference, 2022. URL <https://arxiv.org/abs/2111.02080>.

648 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik
649 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023.
650 URL <https://arxiv.org/abs/2305.10601>.
651

652 Yao Yao, Zuchao Li, and Hai Zhao. Beyond chain-of-thought, effective graph-of-thought reasoning
653 in language models, 2024. URL <https://arxiv.org/abs/2305.16582>.
654

655 Yichi Zhang and Joyce Chai. Hierarchical task learning from language instructions with unified
656 transformers and self-monitoring, 2021. URL <https://arxiv.org/abs/2106.03427>.
657

658 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao,
659 Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi.
660 Siren’s song in the ai ocean: A survey on hallucination in large language models, 2023. URL
<https://arxiv.org/abs/2309.01219>.

661 Wenxuan Zhou, Lerrel Pinto, and Abhinav Gupta. Environment probing interaction policies, 2019.
662 Publisher Copyright: © 7th International Conference on Learning Representations, ICLR 2019.
663 All Rights Reserved.; 7th International Conference on Learning Representations, ICLR 2019 ;
664 Conference date: 06-05-2019 Through 09-05-2019.
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A MORE ALGORITHM

In algorithm 2, we present a beam search example of a hindsight planner. During the collection phase, one Actor prompted from the ground truth sample pool is required to output W actions for each state, and Critic is used to retain the best B actions for the next round of planning. When the search depth U is reached, the best rollout is selected, and the first action from it is returned. At the deployment phase, two Actors are prompted with hindsight prompts and ground truth samples. Each Actor is required to generate $\frac{W}{2}$ actions.

algorithm 3 outlines the algorithm for the collection phase. To preserve the few-shot assumption, the planner collects suboptimal trajectories from \mathcal{D}_{gt} . During the execution of the current task, this task is specifically excluded from being used as an ICL sample to the planner. We employ a prompt generator ϕ to relabel tasks and mitigate ambiguity in the instructions.

Algorithm 2 LLM Planner: A Beam Search example

```

1: Input Actors, Critic, the initial state  $s$ , a generator  $\psi$ , search Breadth  $B$ , proposal width  $W$  and search
   Depth  $U$ 
2: State  $S_0 \leftarrow \{s\}$ 
3: Action array  $A_0 \leftarrow \emptyset$ 
4: Get numbers of Actors  $n \leftarrow \text{len}(\text{Actors})$ 
5: for  $u = 0, \dots, U$  do
6:   for Actor $_i$  in Actors do
7:     For each  $s_u$  in  $S_u$ , invoke Actor $_i$  to propose  $\frac{W}{n}$  candidate actions.
8:   end for
9:   For each  $a_u^{(w)}$  invoke  $\psi$  to generate next state  $s_{u+1}^{(w)}$ 
10:  For each tuple  $(s_u, a_u^{(w)}, s_{u+1}^{(w)})$ , invoke Critic to evaluate the expected cumulative reward  $V_{u+1}^{(w)}$ 
11:  select  $B$  best  $(s_u, a_u^{(w)}, s_{u+1}^{(w)})$  with highest  $V$  and put them into  $S_u \times A_u \times S_{u+1}$ 
12: end for
13: For  $B$  preserved rollouts in  $S_0 \times A_0 \times \dots \times S_{U+1}$ , invoke Critic to evaluate the expected cumulative
   reward  $V_{u+1}^{(b)}$ 
14: Select the best rollout  $(s_0^*, a_0^*, \dots, s_{U+1}^*)$ 
15: return  $a_0^*$ 

```

Algorithm 3 Hindsight Prompt

```

1: input: A ground truth sample pool  $\mathcal{D}_{gt}$ , a prompt generator  $\phi$ .
2: Initialize initiate Agent from  $\mathcal{D}_{gt}$ , set  $\mathcal{D}_{hind} \leftarrow \emptyset$ .
3: for sample  $s$  in  $\mathcal{D}_{gt}$  do
4:   Extract ground truth rollout  $R$ , task description  $I$ , PDDL arguments  $P$  from  $s$ .
5:   Initialize environment  $E$  with  $s$ .
6:   Collect suboptimal trajectories  $\text{traj} \leftarrow \text{Agent}(I, E, \mathcal{D}_{gt}/\{s\})$  (e.g. algorithm 2 of appendix A).
7:   Rename task description  $\tilde{I} \leftarrow \phi(P)$ .
8:   Get reflection Think  $\leftarrow \text{LLM}(\tilde{I}, \text{traj}, R)$ .
9:   Relabel trajectory prompt $_{actor} \leftarrow \text{LLM}(\tilde{I}, \text{traj}, R, \text{Think})$ .
10:  Generate critic from suboptimal trajectory prompt $_{critic} \leftarrow \text{LLM}(\tilde{I}, \text{traj}, R)$ .
11:  Append prompt $_{actor}$ , prompt $_{critic}$  to  $\mathcal{D}_{hind}$ .
12: end for
13: Build hindsight sample pool  $\mathcal{D} = \mathcal{D}_{gt} \cup \mathcal{D}_{hind}$ .
14: Initial Actor $_{gt}$ , Adapter from  $\mathcal{D}_{gt}$ , initial Critic, Actor $_{hind}$  from  $\mathcal{D}_{hind}$ .
15: Return Actor $_{\theta}$ , Critic, Adapter for any  $\theta \in \{gt, hind\}$ .

```

B PROMPTS

B.1 PROMPTS FOR PLANNER

Here, we display prompts for various components here. The `<base.info>` defines the role descriptions while the `<samples>` provide in-context examples for Actors, Critic and Adapter.

We first show the role description for Actors, Critic and Adapter.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

`<base_info>` of Actor

Interact with a household to solve a task.
At each step, you will be provided with the previous observations and action pairs.

Important: You **are required** to return an action.

The answer should contain two parts, the action type and a target.

The allowed types of actions are:

OpenObject, CloseObject, PickupObject, PutObject, ToggleObjectOn, ToggleObjectOff, SliceObject, Stop

The target of OpenObject, CloseObject, PickupObject, ToggleObjectOn, ToggleObjectOff, SliceObject is the object agent interacts with, and the target of PutObject is the place to put the object.

Stop should end with NIL. Note if all requirements are satisfied, you just need to output Stop

`<base_info>` of Critic

You are a value critic of states in a household task. You would be given a task description, some observations and actions, you need to give a critic about them. **Note** Your critic should end with format: the value is a/b=...******

The allowed types of actions are: OpenObject, CloseObject, PickupObject, PutObject, ToggleObjectOn, ToggleObjectOff, SliceObject, Explore, Stop

The target of OpenObject, CloseObject, PickupObject, ToggleObjectOn, ToggleObjectOff, SliceObject is the object agent interacts with and the target of PutObject is the place to put the object.

Explore and Stop should be followed with NIL. Note if all requirements are satisfied, you just need to output Stop. You might need to OpenObject so you can see the object you need to interact with.

`<base_info>` of Adapter

Predict the necessary components for the following household task:

Moveable Receptacle (mrecep_target):** Identify any container or vessel required for the task. Return `None` if not applicable.

Object Slicing (object_sliced):** Determine if the object needs to be sliced. Provide a boolean value (`True` for yes, `False` for no).

Object Target (object_target):** Identify the specific object that is the focus of the task and will be interacted with. This could be the item that needs to be moved, cleaned, heated, cooled, sliced or examined.

Parent Target (parent_target):** Specify the final resting place for the object or its parts. Return `None` if there is no designated location.

Toggle Target (toggle_target):** Indicate any appliance or device that must be toggled during the task. Return `None` if no toggling is required.

Object State (object_state):** Indicate whether the target object needs to be clean, heat, or cool. Return `None` if no such action is required.

810 `-**Two Objects (two_object)**: Specify whether the task requires the`
 811 `agent to handle and place two *identical* objects into the parent`
 812 `target location. Set to True if needed, otherwise False. Note that`
 813 `this parameter should be True only when the task demands picking and`
 814 `placing two of the *same* items.`
 815 `-**Note that the objects you need to predict might not been seen yet.`

816
 817
 818 We then present the `<samples>` to the Actors, the Critic, and the Adapter. Since there are
 819 140 samples for each component, we select only 2 samples from each to demonstrate.

820
 821 `<samples> for Adapter`

822 Task: Place a cup in the coffee maker.
 823 The objects you seen are: Bread,ButterKnife,Cabinet,Chair,CoffeeMachine,
 824 e,CounterTop,Cup,DishSponge,Drawer,Fork,Fridge,GarbageCan,Lettuce,Micr
 825 owave,Mirror,Mug,Pan,Plate,Pot,SaltShaker,Sink,SoapBottle,Spatula,Spo
 826 n,StoveBurner,StoveKnob,DiningTable,SideTable,Toaster,Window
 827 Predict: mrecep_target: None
 828 object_sliced: False
 829 object_target: Mug
 830 parent_target: CoffeeMachine
 831 toggle_target: None
 832 object_state: cool
 833 two_object: False
 834 Task: Warm a cup to make coffee
 835 The objects you seen are: Apple,Bread,ButterKnife,Cabinet,CoffeeMachine,
 836 e,CounterTop,Cup,Drawer,Egg,Fork,Fridge,GarbageCan,HousePlant,Kettle,K
 837 nife,Ladle,Lettuce,Microwave,Mirror,Pan,PepperShaker,Pot,Potato,SaltSh
 838 aker,Sink,Spatula,StoveBurner,StoveKnob,Toaster,Tomato,Window
 839 Predict: mrecep_target: None
 840 object_sliced: False
 841 object_target: Mug
 842 parent_target: CoffeeMachine
 843 toggle_target: None
 844 object_state: heat
 845 two_object: False

846
 847 `<samples> for Actorgt`

848 Task:Place a cup in the coffee maker.
 849 The objects you have seen are:Bread,ButterKnife,Cabinet,Chair,CoffeeMa
 850 chine,CounterTop,Cup,DishSponge,Drawer,Fork,Fridge,GarbageCan,Lettuce,
 851 Microwave,Mirror,Mug,Pan,Plate,Pot,SaltShaker,Sink,SoapBottle,Spatula,
 852 Spoon,StoveBurner,StoveKnob,DiningTable,SideTable,Toaster,Window
 853 Act: OpenObject : Cabinet
 854 >OK
 855 Act: PickupObject : Mug
 856 >OK
 857 Act: CloseObject : Cabinet
 858 >OK
 859 Act: OpenObject : Fridge
 860 >OK
 861 Act: PutObject : Fridge
 862 >OK
 863 Act: CloseObject : Fridge
 >OK
 Act: OpenObject : Fridge
 >OK
 Act: PickupObject : Mug
 >OK
 Act: CloseObject : Fridge
 >OK

864 Act: PutObject : CoffeeMachine
865 >OK
866 Act: Stop : NIL
867 >OK
868
869 Task:Pick up tissues, drop them in the tub.
870 The objects you have seen are:Bathtub,Cabinet,Candle,Cloth,CounterTop,
871 Mirror,Plunger,ScrubBrush,Sink,SoapBar,SoapBottle,Toilet,ToiletPaper,
872 ToiletPaperHanger,Towel,TowelHolder,Window
873 Act: PickupObject : Cloth
874 >OK
875 Act: PutObject : Bathtub
876 >OK
877 Act: Stop : NIL
878 >OK

879 <samples> for Actor_{hind}

880 Task:Place a cup in the coffee maker.
881 The objects you have seen are:Bowl, Bread, ButterKnife, Cabinet, Chair,
882 CoffeeMachine, CounterTop, Cup, DishSponge, Drawer, Fridge, GarbageCan,
883 Lettuce, LightSwitch, Microwave, Mirror, Mug, Pan, PepperShaker, Plate,
884 SaltShaker, Sink, SoapBottle, Spatula, Spoon, StoveBurner, StoveKnob,
885 DiningTable, SideTable, Toaster, Window
886 Act: OpenObject : Cabinet
887 >OK
888 Act: PickupObject : Cup
889 >OK
890 Act: CloseObject : Cabinet
891 >OK
892 Act: PutObject : DiningTable
893 >OK
894 Act: OpenObject : Fridge
895 >OK
896 Act: PickupObject:Cup
897 >OK
898 Act: PutObject:Fridge
899 >OK
900 Act: CloseObject:Fridge
901 >OK
902 Act: OpenObject:Fridge
903 >OK
904 Act: PickupObject:Cup
905 >OK
906 Act: CloseObject:Fridge
907 >OK
908 Act: PutObject:CoffeeMachine
909 >OK
910 Act: Stop : NIL
911 >OK
912 Task:Pick up tissues, drop them in the tub.
913 The objects you have seen are:Bathtub, Cabinet, Candle, Cloth,
914 CounterTop, GarbageCan, HandTowel, HandTowelHolder, LightSwitch,
915 Mirror, Painting, ScrubBrush, Shelf, ShowerDoor, ShowerGlass, Sink,
916 SoapBottle, Television, Toilet, ToiletPaperHanger, Towel, TowelHolder,
917 Window
918 Act: PickupObject : TissueBox
919 >OK
920 Act: PutObject:CounterTop
921 >OK
922 Act: PickupObject:Cloth
923 >OK
924 Act: PutObject:Bathtub

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

```
>OK
Act: Stop : NIL
>OK
```

Having demonstrated the `<base_info>` and `<samples>`, we can now present the prompt template for the Actor, Critic, and Adapter. Note that the prompt of `Actorgt` and `Actorhind` differs in `<samples>`. The `<object_list>` indicates the objects the agent has seen in the environment. Meanwhile, the `<PDDL_predicted>` refers to the output of the Adapter, and the `<K>` indicates the number of samples in each component. To facilitate better comprehension by LLMs, we convert the PDDL arguments into a natural language description. The `<previous_history>` includes the previous actions executed by the agent, enabling the planner to make better decisions based on this information. Concurrently, we employ a prompt generator that reviews the `<previous_history>` and outputs `<history_information>` exclusively to assist LLMs in identifying the objects being held and the open/closed status of containers.

Prompt of Adapter

```
<Adapter_base_info>
Here are <K> examples:
<Adapter_samples>
Your task is: <task_inst>
The objects you have seen are: <object_list>
```

Prompt of Critic

```
<Critic_base_info>
Here are <K> examples:
<Critic_samples>
Your task is: <task_inst>
Your knowledge about this task is: <PDDL_predicted>
The objects you have seen are: <object_list>
previous_history
Based on the actions and Your knowledge about this task , write
a Critic.
Critic:
```

Prompt of Actor

```
<Actor_base_info>
Here are <K> examples:
<Actor_samples>
Your task is: <task_inst>
Your knowledge about this task is: <PDDL_predicted>
The objects you have seen are: <object_list>
Your knowledge about the current state is: <history_information>
<previous_history>
Act:
```

B.2 PROMPTS FOR HINDSIGHT

We now present the prompts used to query LLMs in our hindsight method. During the relabeling process for Actor, we first prompt LLMs to generate a `<Think>` for the suboptimal trajectory, and then we query the LLMs to complete the task based on it. For the relabeling process of the Critic, we directly prompt the LLMs to generate a critic for the suboptimal trajectory. We first present the hindsight samples.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

<samples> for Actor Think

```
Task: Put a fork on a table.
groundtruth rollout:
PickupObject:Fork
PutObject:Sink
ToggleObjectOn:Faucet
ToggleObjectOff:Faucet
PickupObject:Fork
PutObject:SideTable
Stop:NIL
the incomplete rollout:
PickupObject:Fork
PutObject:SideTable

Think: According to the groundtruth rollout, in this incomplete
rollout, I don't clean the fork and the fork is on the sidetable, I
need to pick up the fork and use faucet to clean the fork and put it
onto the sidetable.
Task: Put a warmed apple in the fridge.
groundtruth rollout:
PickupObject:Apple
OpenObject:Microwave
PutObject:Microwave
CloseObject:Microwave
ToggleObjectOn:Microwave
ToggleObjectOff:Microwave
OpenObject:Microwave
PickupObject:Apple
CloseObject:Microwave
OpenObject:Fridge
PutObject:Fridge
CloseObject:Fridge
Stop:NIL
the incomplete rollout:
PickupObject : Apple
OpenObject : Fridge
PutObject : Fridge
CloseObject : Fridge
OpenObject : Fridge
PickupObject : Apple
CloseObject : Fridge
OpenObject : Microwave
Think: According to the groundtruth rollout, in this incomplete
rollout, I don't heat the apple and the apple is in the fridge, I need
to open the fridge, pickup the apple and use microwave to heat the
apple, then I should put the apple back into the fridge.
```

<samples> for Actor Complete

```
task: Put a fork on a table.
groundtruth rollout:
PickupObject:Fork
PutObject:Sink
ToggleObjectOn:Faucet
ToggleObjectOff:Faucet
PickupObject:Fork
PutObject:SideTable
Stop:NIL
the incomplete rollout:
PickupObject:Fork
PutObject:SideTable
```

1026 Think: According to the groundtruth rollout, in this incomplete
 1027 rollout, I don't clean the fork and the fork is on the sidetable, I
 1028 need to pick up the fork and use faucet to clean the fork and put it
 1029 onto the sidetable.
 1030 Based on the Think and groundtruth rollout, the new actions append to
 1031 the incomplete rollout are:
 1032 PickupObject : Fork
 1033 PutObject :Sink
 1034 ToggleObjectOn : Faucet
 1035 ToggleObjectOff : Faucet
 1036 PickupObject: Fork
 1037 PutObject:SideTable
 1038 Stop : NIL

1038 task: Put a warmed apple in the fridge.
 1039 groundtruth rollout:
 1040 PickupObject : Apple
 1041 OpenObject : Microwave
 1042 PutObject : Microwave
 1043 CloseObject : Microwave
 1044 ToggleObjectOn : Microwave
 1045 ToggleObjectOff : Microwave
 1046 OpenObject : Microwave
 1047 PickupObject : Apple
 1048 CloseObject : Microwave
 1049 OpenObject : Fridge
 1050 PutObject : Fridge
 1051 CloseObject : Fridge
 1052 Stop:NIL
 1053 the incomplete rollout:
 1054 PickupObject : Apple
 1055 OpenObject : Fridge
 1056 PutObject : Fridge
 1057 CloseObject : Fridge
 1058 OpenObject : Fridge
 1059 PickupObject : Apple
 1060 CloseObject : Fridge
 1061 OpenObject : Microwave
 1062 Think: According to the groundtruth rollout, in this incomplete
 1063 rollout, I don't heat the apple and the apple is in the fridge, I need
 1064 to open the fridge, pickup the apple and use microwave to heat the
 1065 apple, then I should put the apple back into the fridge.
 1066 Based on the Think and groundtruth rollout, the new actions append to
 1067 the incomplete rollout are:
 1068 OpenObject : Fridge
 1069 PickupObject : Apple
 1070 CloseObject : Fridge
 1071 PutObject: Microwave
 1072 CloseObject: Microwave
 1073 ToggleObjectOn: Microwave
 1074 ToggleObjectOff : Microwave
 1075 OpenObject : Microwave
 1076 PickupObject : Apple
 1077 CloseObject : Microwave
 1078 OpenObject : Fridge
 1079 PutObject : Fridge
 1080 CloseObject : Fridge

1075

1076 <samples> for critic generation

1077 Your task is: Put the cooked tomato on the round table
 1078 The rollout by agent is: OpenObject : Fridge
 1079 PickupObject : Tomato
 1080 CloseObject : Fridge

```

1080 OpenObject : Microwave
1081 The **ground truth rollout** is:
1082 PickupObject:Tomato OpenObject:Microwave
1083 PutObject:Microwave
1084 CloseObject:Microwave
1085 ToggleObjectOn:Microwave
1086 ToggleObjectOff:Microwave
1087 OpenObject:Microwave
1088 PickupObject:Tomato
1089 CloseObject:Microwave
1090 PutObject:DiningTable
1091 Stop:NIL
1092 Based on the **ground truth rollout** , write a critic
1093 Critic:In this task, I need to do the following things in order: Pick
1094 the tomato and put it into microwave, use microwave to heat it,pick the
1095 tomato from microwave and put it onto the DiningTable.There are 5
1096 subgoals in orde, I only achieved first of them, the value is 1/5=0.2.
1097
1098 Your task is: Put a chilled mug in the bottom cabinet closest to the
1099 fridge.
1100 The rollout by agent is:
1101 PickupObject : Mug
1102 OpenObject : Fridge
1103 PutObject : Fridge
1104 CloseObject : Fridge
1105 OpenObject : Fridge
1106 PickupObject : Mug
1107 CloseObject : Fridge
1108 PutObject : Cabinet
1109 The **ground truth rollout** is:
1110 PickupObject:Mug
1111 OpenObject:Fridge
1112 PutObject:Fridge
1113 CloseObject:Fridge
1114 OpenObject:Fridge
1115 PickupObject:Mug
1116 CloseObject:Fridge
1117 OpenObject:Cabinet
1118 PutObject:Cabinet
1119 CloseObject:Cabinet
1120 Stop:NIL
1121 Based on the **ground truth rollout** , write a critic
1122 Critic:In this task, I need to do the following things in order: pick
1123 the mug and put it into the fridge, pick the mug from the fridge and
1124 put the mug into the cabinet. There are 3 subgoals in all, I achieved 2
1125 of them, this is because I don't open the cabinet, so I can't put the
1126 mug into it, the value is 2/3=0.66

```

We provide the prompts used for querying Actors and Critic, respectively. The `<reabeled_task>` indicates the task rewritten based on its PDDL, as detailed in section 5.1. The `<gt_rollout>` represents the ground truth rollout, while the `<suboptimal_rollout>` denotes the rollout collected by our agent.

prompt of Actor_Think

```

1127 You are a housework agent, you will be given a task, a ground truth
1128 rollout to complete this task, and an incomplete rollout.
1129 Your goal is to consider what action you need to append to the
1130 incomplete rollout to complete the task.
1131 Important: You should use your knowledge to judge what actions need to
1132 do based on the ground truth rollout and incomplete rollout. eg: if the
1133 agent forget to open the fridge, then the action of put object into
fridge should be counted as failed, so you should open the fridge and
put the object into the fridge.

```

1134 Important: the openable object (fridge, mrcrowave...) are initially
 1135 closed, so you need to open them before put object in it.
 1136 Important: You can hold one object in your hand at once.
 1137
 1138 The allowed types of actions are: OpenObject, CloseObject, PickupObject, PutObject, ToggleObjectOn, ToggleObjectOff, SliceObject, Stop
 1139 The target of actions like OpenObject, CloseObject, PickupObject,
 1140 ToggleObjectOn, ToggleObjectOff, and SliceObject is the object the
 1141 agent interacts with, whereas the target of PutObject is the location
 1142 where the object is to be placed.
 1143 The 'Stop' action should be followed by 'NIL'. Note that if all
 1144 requirements are met, you only need to output 'Stop'. Remember that you
 1145 can only pick up one item at a time, so you must put down the object in
 1146 your hand before picking up a new one.
 1147 Here are k examples:
 1148 <Actor_Think_samples>
 1149 Task: <reabeled_task>
 1150 Ground truth rollout: <gt_rollout>
 1151 The incomplete rollout: <suboptimal_rollout>
 1152 Think:

1153
 1154 prompt of Actor.Complete
 1155 You are a housework agent, you will be given a task, a ground truth
 1156 rollout to complete this task, an incomplete rollout, and a think
 1157 about the incomplete rollout.
 1158 Your goal is to finish the incomplete rollout based on the groundtruth
 1159 rollout and your think.
 1160 Important: You can only output the needed actions, seperated by '
 1161 ', you must not output other things
 1162
 1163 The allowed types of actions are: OpenObject, CloseObject, PickupObject, PutObject, ToggleObjectOn, ToggleObjectOff, SliceObject, Stop
 1164 The target of actions like OpenObject, CloseObject, PickupObject,
 1165 ToggleObjectOn, ToggleObjectOff, and SliceObject is the object the
 1166 agent interacts with, whereas the target of PutObject is the location
 1167 where the object is to be placed.
 1168 The 'Stop' action should be followed by 'NIL'. Note that if all
 1169 requirements are met, you only need to output 'Stop'. Remember that you
 1170 can only pick up one item at a time, so you must put down the object in
 1171 your hand before picking up a new one.
 1172 Here is k examples:
 1173 <Actor_Complete_samples>
 1174 Task: <reabeled_task>
 1175 Ground truth rollout: <gt_rollout>
 1176 The incomplete rollout: <suboptimal_rollout>
 1177 Think: <Think>
 1178 Based on the Think and groundtruth rollout, the new actions append to
 1179 the incomplete rollout are:

1180 critic generation prompt
 1181 You will be provided with a household task roll-out conducted by an
 1182 agent and a ground truth roll-out. Your task is to write a critic of
 1183 the agent's roll-out based on the **ground truth rollout**. The critic
 1184 should follow the form: In this task, I need do the following things in
 1185 order: ... There are ... subgoals I need to achieve, My current state
 1186 achieve ...
 1187 Important: You should use your knowledge to judge how many subgoals are
 achieved. eg: if the agent forget to open the fridge, then the action
 of put object into fridge should not counted.

1188 Important: Your critic should end with "the value is a/b=.." You can
1189 round it into 2 decimal.
1190 Important: You should write your critic based on given format, you
1191 shouldn't output other things.
1192 Important: You shouldn't mention about ground truth rollout in your
1193 critic.
1194 Here are examples: <Critic_samples>
1195 Your task is: <relabelled_task>
1196 The rollout by agent is:: <suboptimal_rollout>
1197 The ****ground truth rollout**** is: <gt_rollout>
1198 Based on the ****ground truth rollout**** , write a critic
1199 Critic:
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241