

# Hindsight Planner: A Closed-Loop Few-Shot Planner for Embodied Instruction Following

Yuxiao Yang<sup>1</sup>   Shenao Zhang<sup>2</sup>   Zhihan Liu<sup>2</sup>   Huaxiu Yao<sup>3</sup>   Zhaoran Wang<sup>2</sup>  
<sup>1</sup>Shanghai Jiao Tong University   <sup>2</sup>Northwestern University   <sup>3</sup>UNC-Chapel Hill

October 21, 2024

## 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 a suboptimal action is taken, 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-goals, the low-level module outputs actions to finish each sub-goal. 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, we prompt one actor with ground truth samples, while the other actor is prompted with hindsight samples. Thus, the relabeled trajectories can guide the planner in adjusting its policy when incorrect actions are proposed and executed.

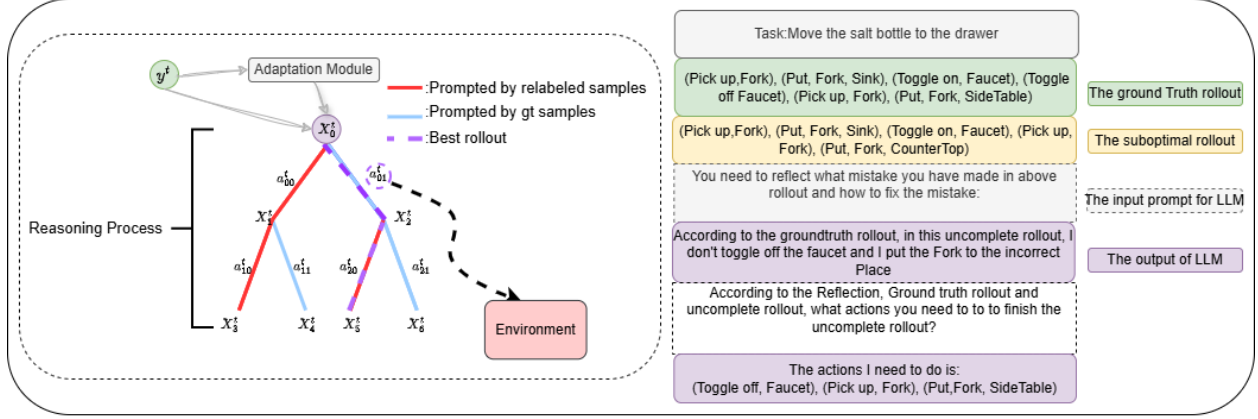


Figure 1: **Left:** The illustration of the Hindsight Planner: at each time step  $t$ , the planner receives a partial observation  $y^t$  from the environment. The adaptation module estimates the latent variable and concatenates it with  $y^t$  to produce the complete state  $x_t$ . **Actor<sub>hind</sub>** and **Actor<sub>gt</sub>** 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 **Actor<sub>hind</sub>**: 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, this approach removes the need for fine-tuning, which can be time-consuming and computationally expensive. To utilize the ICL ability better, many studies propose certain frameworks aimed at enhancing the reasoning capabilities of LLMs (Yao et al., 2023; Wei et al., 2023; Yao et al., 2024). Among them, Liu et al. (2024) proposes a novel perspective by bridging RL and LLM, which inspires us to study ICL from an RL aspect. Xie et al. (2022) interprets ICL as Implicit Bayesian Inference, while Dai et al. (2023) believes that ICL is performing implicit Gradient Descent. All of these imply

the importance of the content in ICL, an area that remains relatively understudied. To this end, we propose Hindsight Planner as an exploration.

## 2.2 Adaptation Module in POMDP

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

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

## 2.3 Hindsight in LLMs

Hindsight algorithms (Andrychowicz et al., 2018; Li et al., 2020; Pong et al., 2020) are widely adopted in the reinforcement learning (RL) area. Generally, the hindsight method aims to reveal future information after collecting a trajectory and relabel the trajectory to make it more informative during training process (Furuta et al., 2022; Andrychowicz et al., 2018). Furuta et al. (2022) applies the hindsight method in training a Transformer model and achieves competitive performance on several baselines. However, training a model from scratch usually requires a large amount of data. In contrast, in-context learning, leveraging the reasoning ability of LLMs, allows an agent to complete complex tasks with only a few samples. Dai et al. (2023) has shown that ICL executes an implicit parameter update. As a result, we utilize ICL in our proposed method. Intuitively, we hope hindsight prompts can provide guidance when an out-of-distribution state is encountered. For example, “Wash a pan and put it away” requires the agent to wash a *Pan* and put it on the *DiningTable*. The trajectory from a planner could be: {(PickupObject, Pan), (PutObject, Sink), (ToggleObjectOn, Faucet), (PickupObject, Pan), (PutObject, CoffeeMachine)}. Note that in this example, the agent

fails to place the pan in the correct location, does not turn off the faucet, and thus the trajectory from the planner is suboptimal. Our hindsight method proposes a novel relabeling process that appends actions to the suboptimal trajectory, aiming to complete the task. In the above example, the corrected trajectory should be:  $\{(\text{PickupObject}, \text{Pan}), (\text{PutObject}, \text{Sink}), (\text{ToggleObjectOn}, \text{Faucet}), (\text{PickupObject}, \text{Pan}), (\text{PutObject}, \text{CoffeeMachine}), (\text{ToggleObjectOff}, \text{Faucet}), (\text{PickupObject}, \text{Pan}), (\text{PutObject}, \text{DiningTable})\}$ . This approach enables us to guide the planner in addressing unknown 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 (the information statistic) in multi-task RL problems. In contrast, our method optimizes 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 trajectory  $\tau$  is defined as  $\tau = \{x_0, y_0, a_0, x_1, y_1, a_1, \dots\}$  and the initial state  $x_0$  is generated through  $x_0 \sim \rho_0(\cdot)$ . The policy  $\pi_\theta(\cdot|y)$  aims to map the observation space into the action space, with  $\theta$  denoting its parameters. The goal of RL is to train a policy such that

$$\pi_\theta = \arg \max_{\pi} \mathbb{E}_{\tau \sim P(\cdot|\pi)} [R(\tau)], \quad (3.1)$$

where  $P(\tau|\pi) = \rho_0(x_0) \prod_{t=0}^T p(x_{t+1}|x_t, a_t) o(y_t|x_t) \pi(a_t|y_t)$ ,  $R(\tau) = \sum_{t=0}^T \gamma^t r(x_t, a_t)$ .

Given a parameterized reward function  $r_z(x, a)$ , where  $z \in \mathcal{Z}$  is a variable indicating the goal for the agent, the conditional policy  $\pi(\cdot|y, z)$  aims to accomplish different goals based on its observations. The goal in Equation (3.1) becomes

$$\pi_\theta = \arg \max_{\pi} \mathbb{E}_{\tau \sim P(\cdot|\pi, z), z \sim p(z)} [R_z(\tau)], \quad (3.2)$$

where  $R_z(\tau) = \sum_{t=0}^{\infty} \gamma^t r_z(x_t, a_t)$ . 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  $\pi_\theta$  that satisfies

$$\pi_\theta = \arg \min_{\pi} \mathbb{E}_{\tau \sim P(\cdot|\pi, z), z \sim p(z)} [\text{KL}(I(\tau), z)], \quad (3.3)$$

where  $I(\tau)$  is *information statistic* that can be any function that captures the desired information from a trajectory  $\tau_t = \{x_0, y_0, a_0, x_1, y_1, a_1, \dots, x_t, y_t\}$  and KL is the Kullback-Leibler divergence. This optimization objective has achieved competitive results in previous studies (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  $\pi(\cdot|x, z)$ . 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 the sub-goal index  $k \leftarrow 0$ .
  - 3: Get sample pool  $\mathcal{D}$  and initialize **Actor** $_{\theta}$ , **Critic**, **Adapter** from  $\mathcal{D}$ , for any  $\theta \in \{\text{gt}, \text{hind}\}$  (e.g. Algorithm 3 in Appendix A). (Hindsight process)
  - 4: **while** Not *Finished* **do**
  - 5:   Get PDDL arguments  $P \leftarrow \text{Adapter}(I, O)$ .
  - 6:   Plan and get sub-goal  $S_k \leftarrow \text{LLM-PL}(\text{Actor}_{\text{gt}}, \text{Actor}_{\text{hind}}, \text{Critic}, P, I, O, G)$  (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  $\mathcal{D}_{\text{hind}}$ . The complete dataset  $\mathcal{D} = \mathcal{D}_{\text{hind}} \cup \mathcal{D}_{\text{gt}}$ , where  $\mathcal{D}_{\text{gt}}$  is constructed from training data. In the deployment phase, we initiate hindsight actor **Actor** $_{\text{hind}}$ , ground truth actor **Actor** $_{\text{gt}}$ , and **Critic** from  $\mathcal{D}$ .

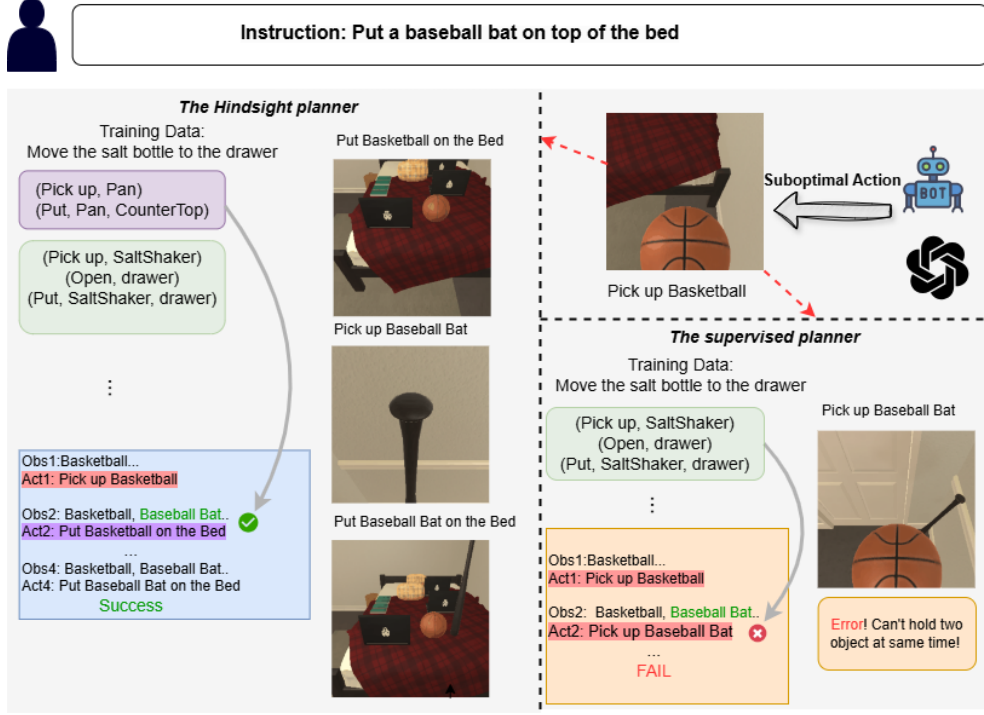


Figure 2: A comparison of Hindsight Planner and previous supervised methods when taking a suboptimal action. The agent initially picks up the incorrect object (“Basketball”). In the supervised method, the planner fails to handle this situation, which leads to task failure. In contrast, the Hindsight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

## 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.

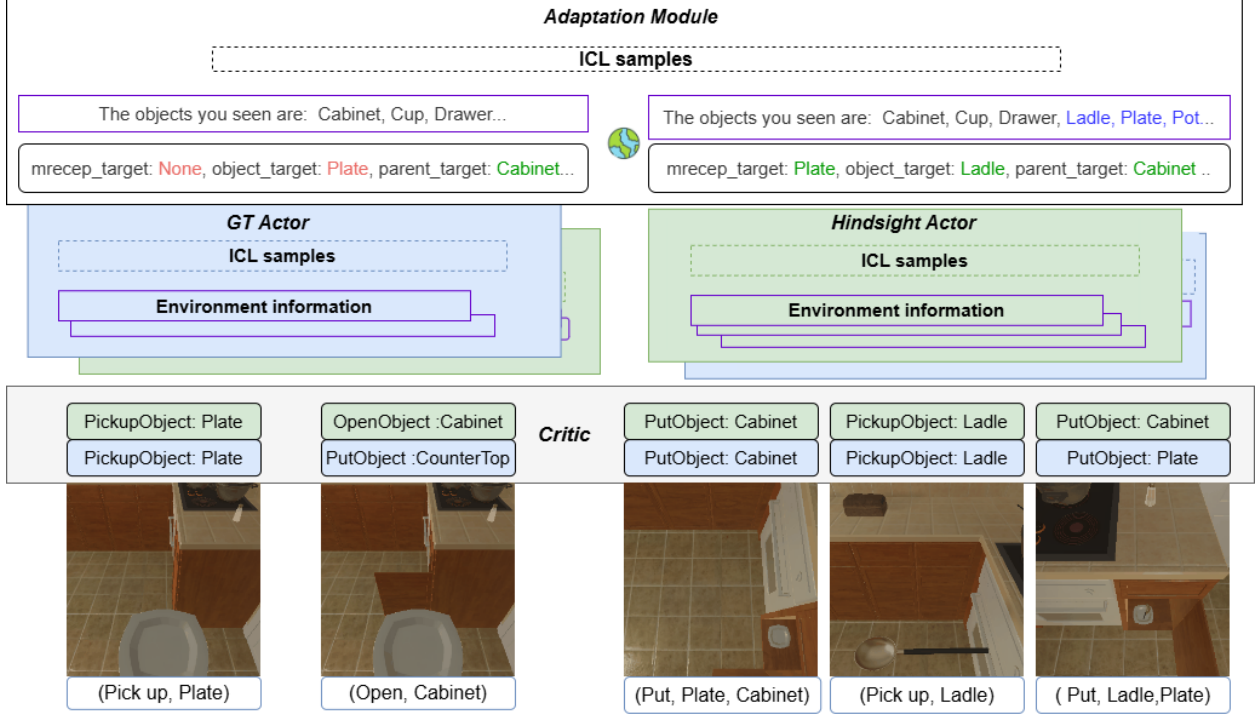


Figure 3: The entire process of the Hindsight Planner is as follows: At the start of the task, which is to “Place a plate with a ladle on it in a cabinet,” the **Adapter** mistakenly identifies the task as picking up a plate and placing it into a cabinet. **Actor<sub>hind</sub>** and **Actor<sub>gt</sub>** 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.

### 4.3 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  $\pi$  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

$$\hat{\pi} = \arg \max_{\pi} \mathbb{E}_{\tau \sim P(\cdot | \pi, z), z \sim q(z)} [R_z(\tau)], \quad (4.1)$$

Define the  $\pi^*$  as the oracle. It is easy to see that

$$\mathbb{E}_{\tau \sim P(\cdot | \hat{\pi}, z), z \sim p(z)} [R_z(\tau)] < \mathbb{E}_{\tau \sim P(\cdot | \pi^*, z), z \sim p(z)} [R_z(\tau)], \quad (4.2)$$

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



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

$$\min_{\pi} \mathbb{E}_{\tau \sim P(\cdot | \pi, z), z \sim p(z)} [\text{KL}(I(\tau), I(\tau^*))]. \quad (4.3)$$

Our method utilizes LLMs to relabel  $\hat{\tau} = \tau_T + \{a_T, x_{T+1}, y_{T+1}, a_{T+2}, \dots\}$  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  $\tau$  and relabels them. The algorithm is summarized in Algorithm 3 of Appendix A. During the deployment phase, the **Actor<sub>gt</sub>** is prompted with ground truth samples while the **Actor<sub>hind</sub>** and the **Critic** are prompted with relabeled samples. Intuitively, we hope that the **Actor<sub>gt</sub>** 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 **Actor<sub>hind</sub>** 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: **Actor<sub>gt</sub>** and **Actor<sub>hind</sub>**, 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  $\psi$  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 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 *BlackTable* 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 (OpenAI et al., 2024) as the target LLM and set temperature to 0. For the **Adapter**, 5 in-context

Model	n-shot	Test Seen		Test Unseen	
		SR	GC	SR	GC
HiTUT (Zhang and 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>	<b>39.55</b>	<b>27.80</b>	<b>38.52</b>
MCR-Agent (Bhambri et al., 2024)	full	<b>30.13</b>	-	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	<b>25.51</b>	<b>34.74</b>	<b>18.77</b>	<b>28.29</b>

Table 1: **Comparison with the state-of-the-art methods on SR and GC in the test set.** Bold numbers represent the highest level of accuracy, whereas underlined numbers signify the second-highest accuracy for each experimental configuration. “low inst.” refers to the use of step-by-step instructions.

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	<b>29.63</b>	<b>38.74</b>	<u>18.28</u>	<b>31.24</b>	<u>25.11</u>	<b>35.79</b>	<b>20.27</b>	<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>	<b>19.00</b>	<u>29.90</u>	<b>25.51</b>	<u>34.74</u>	<u>18.77</u>	<b>28.29</b>

Table 2: **Comparison with the same lower-controller.** Bold numbers represent the highest level of accuracy, whereas underlined numbers signify the second-highest accuracy for each experimental configuration.

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,

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 Method	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.

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 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

Model	Valid Seen		Valid Unseen	
	SR	GC	SR	GC
W.O. Adaptation Module	23.17	33.28	14.99	27.36
W.O. Hindsight Method	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.

rates for the “Valid Seen” and “Valid Unseen” splits. Then, we remove the hindsight prompts. For a fair comparison, the original planner requires both **Actor<sub>gt</sub>** and **Actor<sub>hind</sub>** to generate one action per state. We also ablate by prompting **Actor<sub>gt</sub>** 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. Additionally, we also present the average sub-goal lengths in Table 5. 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 are more likely to output suboptimal actions, allowing the hindsight actor to correct its mistakes. 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.

## 6 Conclusion

This paper explores an effective few-shot framework for Embodied Instruction Following. We approach the task as a POMDP and

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

Table 5: Average sub-goal lengths.

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

- Ahn, M., Brohan, A., Brown, N., Chebotar, Y., Cortes, O., David, B., Finn, C., Fu, C., Gopalakrishnan, K., Hausman, K., Herzog, A., Ho, D., Hsu, J., Ibarz, J., Ichter, B., Irpan, A., Jang, E., Ruano, R. J., Jeffrey, K., Jesmonth, S., Joshi, N. J., Julian, R., Kalashnikov, D., Kuang, Y., Lee, K.-H., Levine, S., Lu, Y., Luu, L., Parada, C., Pastor, P., Quiambao, J., Rao, K., Rettinghouse, J., Reyes, D., Sermanet, P., Sievers, N., Tan, C., Toshev, A., Vanhoucke, V., Xia, F., Xiao, T., Xu, P., Xu, S., Yan, M. and Zeng, A. (2022). Do as i can, not as i say: Grounding language in robotic affordances.  
<https://arxiv.org/abs/2204.01691>
- Andrychowicz, M., Wolski, F., Ray, A., Schneider, J., Fong, R., Welinder, P., McGrew, B., Tobin, J., Abbeel, P. and Zaremba, W. (2018). Hindsight experience replay.  
<https://arxiv.org/abs/1707.01495>
- Bhambri, S., Kim, B. and Choi, J. (2024). Multi-level compositional reasoning for interactive instruction following.  
<https://arxiv.org/abs/2308.09387>
- Blukis, V., Paxton, C., Fox, D., Garg, A. and Artzi, Y. (2021). A persistent spatial semantic representation for high-level natural language instruction execution. *Cornell University - arXiv, Cornell University - arXiv*.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I. and Amodei, D. (2020). Language models are few-shot learners.  
<https://arxiv.org/abs/2005.14165>
- Chapman, D. (1987). Planning for conjunctive goals. *Artif. Intell.*, **32** 333–377.  
<https://api.semanticscholar.org/CorpusID:1525549>
- Dai, D., Sun, Y., Dong, L., Hao, Y., Ma, S., Sui, Z. and Wei, F. (2023). Why can gpt learn in-context? language models implicitly perform gradient descent as meta-optimizers.  
<https://arxiv.org/abs/2212.10559>

- Devlin, J., Chang, M.-W., Lee, K. and Toutanova, K. (2019). Bert: Pre-training of deep bidirectional transformers for language understanding.  
<https://arxiv.org/abs/1810.04805>
- Dong, Q., Li, L., Dai, D., Zheng, C., Ma, J., Li, R., Xia, H., Xu, J., Wu, Z., Chang, B., Sun, X., Li, L. and Sui, Z. (2024). A survey on in-context learning.  
<https://arxiv.org/abs/2301.00234>
- Eysenbach, B., Geng, X., Levine, S. and Salakhutdinov, R. (2020). Rewriting history with inverse rl: Hindsight inference for policy improvement.  
<https://arxiv.org/abs/2002.11089>
- Furuta, H., Matsuo, Y. and Gu, S. S. (2022). Generalized decision transformer for offline hindsight information matching.  
<https://arxiv.org/abs/2111.10364>
- Ghosh, D., Gupta, A., Fu, J., Reddy, A., Devin, C., Eysenbach, B. and Levine, S. (2019). Learning to reach goals without reinforcement learning. *ArXiv*, **abs/1912.06088**.
- Guo, J., Zhang, R., Zhang, X., Peng, S., Yi, Q., Du, Z., Hu, X., Guo, Q. and Chen, Y. (2021). Hindsight value function for variance reduction in stochastic dynamic environment.  
<https://arxiv.org/abs/2107.12216>
- Hazan, E., Kakade, S. M., Singh, K. and Soest, A. V. (2019). Provably efficient maximum entropy exploration.  
<https://arxiv.org/abs/1812.02690>
- Kim, B., Kim, J., Kim, Y., Min, C. and Choi, J. (2024). Context-aware planning and environment-aware memory for instruction following embodied agents.  
<https://arxiv.org/abs/2308.07241>
- Kumar, A., Fu, Z., Pathak, D. and Malik, J. (2021). Rma: Rapid motor adaptation for legged robots.  
<https://arxiv.org/abs/2107.04034>
- Lee, J. N., Agarwal, A., Dann, C. and Zhang, T. (2023). Learning in pomdps is sample-efficient with hindsight observability.  
<https://arxiv.org/abs/2301.13857>
- Lee, L., Eysenbach, B., Parisotto, E., Xing, E., Levine, S. and Salakhutdinov, R. (2020). Efficient exploration via state marginal matching.  
<https://arxiv.org/abs/1906.05274>
- Li, A. C., Pinto, L. and Abbeel, P. (2020). Generalized hindsight for reinforcement learning.  
<https://arxiv.org/abs/2002.11708>

Liu, Z., Hu, H., Zhang, S., Guo, H., Ke, S., Liu, B. and Wang, Z. (2024). Reason for future, act for now: A principled framework for autonomous llm agents with provable sample efficiency.

<https://arxiv.org/abs/2309.17382>

Min, S., Chaplot, D., Ravikumar, P., Bisk, Y. and Salakhutdinov, R. (2021). Film: Following instructions in language with modular methods. *Learning, Learning*.

OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J., Bello, I., Berdine, J., Bernadett-Shapiro, G., Berner, C., Bogdonoff, L., Boiko, O., Boyd, M., Brakman, A.-L., Brockman, G., Brooks, T., Brundage, M., Button, K., Cai, T., Campbell, R., Cann, A., Carey, B., Carlson, C., Carmichael, R., Chan, B., Chang, C., Chantzis, F., Chen, D., Chen, S., Chen, R., Chen, J., Chen, M., Chess, B., Cho, C., Chu, C., Chung, H. W., Cummings, D., Currier, J., Dai, Y., Decareaux, C., Degry, T., Deutsch, N., Deville, D., Dhar, A., Dohan, D., Dowling, S., Dunning, S., Ecoffet, A., Eleti, A., Eloundou, T., Farhi, D., Fedus, L., Felix, N., Fishman, S. P., Forte, J., Fulford, I., Gao, L., Georges, E., Gibson, C., Goel, V., Gogineni, T., Goh, G., Gontijo-Lopes, R., Gordon, J., Grafstein, M., Gray, S., Greene, R., Gross, J., Gu, S. S., Guo, Y., Hallacy, C., Han, J., Harris, J., He, Y., Heaton, M., Heidecke, J., Hesse, C., Hickey, A., Hickey, W., Hoeschele, P., Houghton, B., Hsu, K., Hu, S., Hu, X., Huizinga, J., Jain, S., Jain, S., Jang, J., Jiang, A., Jiang, R., Jin, H., Jin, D., Jomoto, S., Jonn, B., Jun, H., Kaftan, T., Łukasz Kaiser, Kamali, A., Kanitscheider, I., Keskar, N. S., Khan, T., Kilpatrick, L., Kim, J. W., Kim, C., Kim, Y., Kirchner, J. H., Kiros, J., Knight, M., Kokotajlo, D., Łukasz Kondraciuk, Kondrich, A., Konstantinidis, A., Kosic, K., Krueger, G., Kuo, V., Lampe, M., Lan, I., Lee, T., Leike, J., Leung, J., Levy, D., Li, C. M., Lim, R., Lin, M., Lin, S., Litwin, M., Lopez, T., Lowe, R., Lue, P., Makanju, A., Malfacini, K., Manning, S., Markov, T., Markovski, Y., Martin, B., Mayer, K., Mayne, A., McGrew, B., McKinney, S. M., McLeavey, C., McMillan, P., McNeil, J., Medina, D., Mehta, A., Menick, J., Metz, L., Mishchenko, A., Mishkin, P., Monaco, V., Morikawa, E., Mossing, D., Mu, T., Murati, M., Murk, O., Mély, D., Nair, A., Nakano, R., Nayak, R., Neelakantan, A., Ngo, R., Noh, H., Ouyang, L., O’Keefe, C., Pachocki, J., Paino, A., Palermo, J., Pantuliano, A., Parascandolo, G., Parish, J., Parparita, E., Passos, A., Pavlov, M., Peng, A., Perelman, A., de Avila Belbute Peres, F., Petrov, M., de Oliveira Pinto, H. P., Michael, Pokorny, Pokrass, M., Pong, V. H., Powell, T., Power, A., Power, B., Proehl, E., Puri, R., Radford, A., Rae, J., Ramesh, A., Raymond, C., Real, F., Rimbach, K., Ross, C., Rotsted, B., Roussez, H., Ryder, N., Saltarelli, M., Sanders, T., Santurkar, S., Sastry, G., Schmidt, H., Schnurr, D., Schulman, J., Selsam, D., Sheppard, K., Sherbakov, T., Shieh, J., Shoker, S., Shyam, P., Sidor, S., Sigler, E., Simens, M., Sitkin, J., Slama, K., Sohl, I., Sokolowsky, B., Song, Y., Staudacher, N., Such, F. P., Summers, N., Sutskever, I., Tang, J., Tezak, N., Thompson, M. B., Tillet, P., Tootoonchian, A., Tseng, E., Tuggle, P., Turley, N., Tworek, J., Uribe, J. F. C., Vallone, A., Vijayvergiya, A., Voss, C., Wainwright, C., Wang, J. J., Wang, A., Wang, B., Ward, J., Wei, J., Weinmann, C., Welihinda, A., Welinder, P., Weng, J., Weng, L., Wiethoff, M., Willner, D., Winter, C.,



- Wolrich, S., Wong, H., Workman, L., Wu, S., Wu, J., Wu, M., Xiao, K., Xu, T., Yoo, S., Yu, K., Yuan, Q., Zaremba, W., Zellers, R., Zhang, C., Zhang, M., Zhao, S., Zheng, T., Zhuang, J., Zhuk, W. and Zoph, B. (2024). Gpt-4 technical report.  
<https://arxiv.org/abs/2303.08774>
- Peng, X. B., Coumans, E., Zhang, T., Lee, T.-W., Tan, J. and Levine, S. (2020). Learning agile robotic locomotion skills by imitating animals.  
<https://arxiv.org/abs/2004.00784>
- Pong, V., Gu, S., Dalal, M. and Levine, S. (2020). Temporal difference models: Model-free deep rl for model-based control.  
<https://arxiv.org/abs/1802.09081>
- Shin, S., jeon, S., Kim, J., Kang, G.-C. and Zhang, B.-T. (2024). Socratic planner: Inquiry-based zero-shot planning for embodied instruction following.  
<https://arxiv.org/abs/2404.15190>
- Shridhar, M., Thomason, J., Gordon, D., Bisk, Y., Han, W., Mottaghi, R., Zettlemoyer, L. and Fox, D. (2020). Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.  
<http://dx.doi.org/10.1109/cvpr42600.2020.01075>
- Silver, T., Dan, S., Srinivas, K., Tenenbaum, J. B., Kaelbling, L. P. and Katz, M. (2023). Generalized planning in pddl domains with pretrained large language models.  
<https://arxiv.org/abs/2305.11014>
- Song, C. H., Wu, J., Washington, C., Sadler, B. M., Chao, W.-L. and Su, Y. (2023). Llm-planner: Few-shot grounded planning for embodied agents with large language models.  
<https://arxiv.org/abs/2212.04088>
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E. and Lample, G. (2023). Llama: Open and efficient foundation language models.  
<https://arxiv.org/abs/2302.13971>
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. and Polosukhin, I. (2023). Attention is all you need.  
<https://arxiv.org/abs/1706.03762>
- Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., Chi, E. H., Hashimoto, T., Vinyals, O., Liang, P., Dean, J. and Fedus, W. (2022). Emergent abilities of large language models.  
<https://arxiv.org/abs/2206.07682>
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q. and Zhou, D. (2023). Chain-of-thought prompting elicits reasoning in large language models.  
<https://arxiv.org/abs/2201.11903>

- Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Scao, T. L., Gugger, S., Drame, M., Lhoest, Q. and Rush, A. M. (2020). Huggingface’s transformers: State-of-the-art natural language processing.  
<https://arxiv.org/abs/1910.03771>
- Xie, S. M., Raghunathan, A., Liang, P. and Ma, T. (2022). An explanation of in-context learning as implicit bayesian inference.  
<https://arxiv.org/abs/2111.02080>
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T. L., Cao, Y. and Narasimhan, K. (2023). Tree of thoughts: Deliberate problem solving with large language models.  
<https://arxiv.org/abs/2305.10601>
- Yao, Y., Li, Z. and Zhao, H. (2024). Beyond chain-of-thought, effective graph-of-thought reasoning in language models.  
<https://arxiv.org/abs/2305.16582>
- Zhang, Y. and Chai, J. (2021). Hierarchical task learning from language instructions with unified transformers and self-monitoring.  
<https://arxiv.org/abs/2106.03427>
- Zhang, Y., Li, Y., Cui, L., Cai, D., Liu, L., Fu, T., Huang, X., Zhao, E., Zhang, Y., Chen, Y., Wang, L., Luu, A. T., Bi, W., Shi, F. and Shi, S. (2023). Siren’s song in the ai ocean: A survey on hallucination in large language models.  
<https://arxiv.org/abs/2309.01219>
- Zhou, W., Pinto, L. and Gupta, A. (2019). Environment probing interaction policies. Publisher Copyright: © 7th International Conference on Learning Representations, ICLR 2019. All Rights Reserved.; 7th International Conference on Learning Representations, ICLR 2019 ; Conference date: 06-05-2019 Through 09-05-2019.

## 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  $\mathcal{D}_{\text{gt}}$ . 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  $\phi$  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  $\psi$ , the search Breadth  $B$ , the proposal width  $W$  and the search Depth  $U$ .
  - 2: Set State  $S_0 \leftarrow \{s\}$ .
  - 3: Set 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** <sub>$i$</sub>  in **Actors** **do**
  - 7:     For each  $s_u$  in  $S_u$ , invoke **Actor** <sub>$i$</sub>  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  $(s_u, a_u^{(w)}, s_{u+1}^{(w)})$ , use **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 write them to  $S_u \times A_u \times S_{u+1}$ .
  - 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  $(s_0^*, a_0^*, \dots, s_{U+1}^*)$  and return  $a_0^*$ .
- 

## B Prompts

### B.1 Prompts for Planner

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

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

<base_info> of Actor
<p>Interact with a household to solve a task.</p> <p>At each step, you will be provided with the previous observations and action pairs.</p> <p>Important: You <b>**are required**</b> to return an action.</p>

---

**Algorithm 3** Hindsight Prompt

---

- 1: **Input:** A ground truth sample pool  $\mathcal{D}_{gt}$ , a prompt generator  $\phi$ .
  - 2: **Initialize** Initialize **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  $traj \leftarrow \text{Agent}(I, E, \mathcal{D}_{gt}/\{s\})$  (e.g. Algorithm 2 of Appendix A).
  - 7:   Rename the task description  $\tilde{I} \leftarrow \phi(P)$ .
  - 8:   Get reflection Think  $\leftarrow \text{LLM}(\tilde{I}, traj, R)$ .
  - 9:   Relabel trajectory prompt<sub>actor</sub>  $\leftarrow \text{LLM}(\tilde{I}, traj, R, \text{Think})$ .
  - 10:   Generate critic from suboptimal trajectory prompt<sub>critic</sub>  $\leftarrow \text{LLM}(\tilde{I}, traj, R)$ .
  - 11:   Append prompt<sub>actor</sub>, prompt<sub>critic</sub> to  $\mathcal{D}_{hind}$ .
  - 12: **end for**
  - 13: Build hindsight sample pool  $\mathcal{D} = \mathcal{D}_{gt} \cup \mathcal{D}_{hind}$ .
  - 14: Initialize Actor<sub>gt</sub>, Adapter from  $\mathcal{D}_{gt}$ , initial Critic, Actor<sub>hind</sub> from  $\mathcal{D}_{hind}$ .
  - 15: **Return** Actor <sub>$\theta$</sub> , Critic, Adapter for any  $\theta \in \{\text{gt}, \text{hind}\}$ .
- 

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

We then show the `<samples>` for the `Actors`, the `Critic`, and the `Adapter`. Since there are 140 samples for each component, we select just 2 samples from each to demonstrate.

<samples> for Adapter
-----------------------

```
Task: Place a cup in the coffee maker.
The objects you seen are: Bread,ButterKnife,Cabinet,Chair,CoffeeMachine,CounterTop,Cup,DishSpoon,Drawer,Fork,Fridge,GarbageCan,Lettuce,Microwave,Mirror,Mug,Pan,Plate,Pot,SaltShaker,Sink,SoupBottle,Spatula,Spoon,StoveBurner,StoveKnob,DiningTable,SideTable,Toaster,Window
Predict:
mrecep_target: None
object_sliced: False
object_target: Mug
parent_target: CoffeeMachine
toggle_target: None
object_state: cool
two_object: False
Task: Warm a cup to make coffee
The objects you seen are: Apple,Bread,ButterKnife,Cabinet,CoffeeMachine,CounterTop,Cup,Drawer,Egg,Fork,Fridge,GarbageCan,HousePlant,Kettle,Knife,Ladle,Lettuce,Microwave,Mirror,Pan,PepperShaker,Pot,Potato,SaltShaker,Sink,Spatula,StoveBurner,StoveKnob,Toaster,Tomato,Window
Predict:
mrecep_target: None
object_sliced: False
object_target: Mug
parent_target: CoffeeMachine
toggle_target: None
object_state: heat
```

two\_object: False

<samples> for Actor<sub>gt</sub>

Task:Place a cup in the coffee maker.

The objects you have seen are:Bread,ButterKnife,Cabinet,Chair,CoffeeMachine,CounterTop,Cup,DishSponge,Drawer,Fork,Fridge,GarbageCan,Lettuce,Microwave,Mirror,Mug,Pan,Plate,Pot,SaltShaker,Sink,SoapBottle,Spatula,Spoon,StoveBurner,StoveKnob,DiningTable,SideTable,Toaster,Window

Act: OpenObject : Cabinet

>OK

Act: PickupObject : Mug

>OK

Act: CloseObject : Cabinet

>OK

Act: OpenObject : Fridge

>OK

Act: PutObject : Fridge

>OK

Act: CloseObject : Fridge

>OK

Act: OpenObject : Fridge

>OK

Act: PickupObject : Mug

>OK

Act: CloseObject : Fridge

>OK

Act: PutObject : CoffeeMachine

>OK

Act: Stop : NIL

>OK

Task:Pick up tissues, drop them in the tub.

The objects you have seen are:Bathtub,Cabinet,Candle,Cloth,CounterTop,Mirror,Plunger,ScrubBrush,Sink,SoapBar,SoapBottle,Toilet,ToiletPaper,ToiletPaperHanger,Towel,TowelHolder,Window

Act: PickupObject : Cloth

>OK

Act: PutObject : Bathtub

>OK

Act: Stop : NIL

>OK

<samples> for Actor<sub>hind</sub>

Task:Place a cup in the coffee maker.

The objects you have seen are:Bowl, Bread, ButterKnife, Cabinet, Chair, CoffeeMachine, CounterTop, Cup, DishSponge, Drawer, Fridge, GarbageCan, Lettuce, LightSwitch, Microwave, Mirror, Mug, Pan, PepperShaker, Plate, SaltShaker, Sink, SoapBottle, Spatula, Spoon, StoveBurner, StoveKnob, DiningTable, SideTable, Toaster, Window

```

Act: OpenObject : Cabinet
>OK
Act: PickupObject : Cup
>OK
Act: CloseObject : Cabinet
>OK
Act: PutObject : DiningTable
>OK
Act: OpenObject : Fridge
>OK
Act: PickupObject:Cup
>OK
Act: PutObject:Fridge
>OK
Act: CloseObject:Fridge
>OK
Act: OpenObject:Fridge
>OK
Act: PickupObject:Cup
>OK
Act: CloseObject:Fridge
>OK
Act: PutObject:CoffeeMachine
>OK
Act: Stop : NIL
>OK

```

Task:Pick up tissues, drop them in the tub.

The objects you have seen are:Bathtub, Cabinet, Candle, Cloth, CounterTop, GarbageCan, HandTowel, HandTowelHolder, LightSwitch, Mirror, Painting, ScrubBrush, Shelf, ShowerDoor, ShowerGlass, Sink, SoapBottle, Television, Toilet, ToiletPaperHanger, Towel, TowelHolder, Window

```

Act: PickupObject : TissueBox
>OK
Act: PutObject:CounterTop
>OK
Act: PickupObject:Cloth
>OK
Act: PutObject:Bathtub
>OK
Act: Stop : NIL
>OK

```

Having demonstrated the `<base_info>` and `<samples>`, we can now present the prompt template for the `Actors`, the `Critic`, and the `Adapter`. Note that the prompt of `Actorgt` and `Actorhind` differ in `<samples>`. The `<object_list>` indicates the objects the agent has seen in the environment. Meanwhile, the `<PDDL_predicted>` represents to the output of the `Adapter`, and the `<K>` indicates the number of samples in each component. To facilitate better comprehension by LLMs, the PDDL



arguments are converted 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, a prompt generator reviews the `<previous_history>` and outputs `<history_information>` to assist LLMs in identifying 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 the **Actor**, we first prompt the LLMs to generate a `<Think>` for the suboptimal trajectory, and then query them to complete the task based on it. For the relabeling process of the **Critic**, we directly prompt the LLMs to generate an evaluation for the suboptimal trajectory. We first present the hindsight samples for clarity.

<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
```

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.

Based on the Think and groundtruth rollout, the new actions append to the incomplete rollout are:

```
PickupObject : Fork
PutObject :Sink
ToggleObjectOn : Faucet
ToggleObjectOff : Faucet
PickupObject: Fork
PutObject:SideTable
Stop : NIL
```

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.

Based on the Think and groundtruth rollout, the new actions append to the incomplete rollout are:

OpenObject : Fridge  
PickupObject : Apple  
CloseObject : Fridge  
PutObject: Microwave  
CloseObject: Microwave  
ToggleObjectOn: Microwave  
ToggleObjectOff : Microwave  
OpenObject : Microwave  
PickupObject : Apple  
CloseObject : Microwave  
OpenObject : Fridge  
PutObject : Fridge  
CloseObject : Fridge

<samples> for critic generation

Your task is: Put the cooked tomato on the round table

The rollout by agent is: OpenObject : Fridge

PickupObject : Tomato

CloseObject : Fridge

OpenObject : Microwave

The **ground truth rollout** is:

PickupObject:Tomato OpenObject:Microwave

PutObject:Microwave

CloseObject:Microwave

ToggleObjectOn:Microwave

ToggleObjectOff:Microwave

OpenObject:Microwave

PickupObject:Tomato

CloseObject:Microwave

PutObject:DiningTable

Stop:NIL

Based on the **ground truth rollout** , write a critic

Critic:In this task, I need to do the following things in order: Pick the tomato and put it into microwave, use microwave to heat it,pick the tomato from microwave and put it onto the DiningTable.There are 5 subgoals in orde, I only achieved first of them, the value is  $1/5=0.2$ .

Your task is: Put a chilled mug in the bottom cabinet closest to the fridge.

The rollout by agent is:

PickupObject : Mug

OpenObject : Fridge

PutObject : Fridge

CloseObject : Fridge

OpenObject : Fridge

PickupObject : Mug

CloseObject : Fridge

PutObject : Cabinet

The **ground truth rollout** is:

```

PickupObject:Mug
OpenObject:Fridge
PutObject:Fridge
CloseObject:Fridge
OpenObject:Fridge
PickupObject:Mug
CloseObject:Fridge
OpenObject:Cabinet
PutObject:Cabinet
CloseObject:Cabinet
Stop:NIL
Based on the **ground truth rollout** , write a critic
Critic:In this task, I need to do the following things in order: pick the mug and put it into
the fridge, pick the mug from the fridge and put the mug into the cabinet. There are 3 subgoals
in all, I achieved 2 of them, this is because I don't open the cabinet, so I can't put the mug
into it, the value is  $2/3=0.66$ 

```

We provide the prompts used for querying Actors and Critic, respectively. The `<relabeled_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

You are a housework agent, you will be given a task, a ground truth rollout to complete this task, and an incomplete rollout.

Your goal is to consider what action you need to append to the incomplete rollout to complete the task.

Important: You should use your knowledge to judge what actions need to do based on the ground truth rollout and incomplete rollout. eg: if the 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.

Important: the openable object (fridge,mrcrowave...) are initially closed, so you need to open them before put object in it.

Important: You can hold one object in your hand at once.

The allowed types of actions are:  
 OpenObject,CloseObject,PickupObject,PutObject,ToggleObjectOn,ToggleObjectOff,SliceObject,Stop  
 The target of actions like OpenObject, CloseObject, PickupObject, ToggleObjectOn, ToggleObjectOff, and SliceObject is the object the agent interacts with, whereas the target of PutObject is the location where the object is to be placed.  
 The 'Stop' action should be followed by 'NIL'. Note that if all requirements are met, you only need to output 'Stop'. Remember that you can only pick up one item at a time, so you must put down the object in your hand before picking up a new one.

Here are k examples:  
`<Actor_Think_samples>`  
 Task: `<relabeled_task>`  
 Ground truth rollout: `<gt_rollout>`

The incomplete rollout: <suboptimal\_rollout>  
Think:

prompt of Actor\_Complete

You are a housework agent, you will be given a task, a ground truth rollout to complete this task, an incomplete rollout, and a think about the incomplete rollout. Your goal is to finish the incomplete rollout based on the groundtruth rollout and your think. Important: You can only output the needed actions,seperated by ' ', you must not output other things

The allowed types of actions are:

OpenObject,CloseObject,PickupObject,PutObject,ToggleObjectOn,ToggleObjectOff,SliceObject,Stop  
The target of actions like OpenObject, CloseObject, PickupObject, ToggleObjectOn, ToggleObjectOff, and SliceObject is the object the agent interacts with, whereas the target of PutObject is the location where the object is to be placed.  
The 'Stop' action should be followed by 'NIL'. Note that if all requirements are met, you only need to output 'Stop'. Remember that you can only pick up one item at a time, so you must put down the object in your hand before picking up a new one.

Here are k examples:

<Actor\_Complete\_samples>

Task: <relabeled\_task>

Ground truth rollout: <gt\_rollout>

The incomplete rollout: <suboptimal\_rollout>

Think: <Think>

Based on the Think and groundtruth rollout, the new actions append to the incomplete rollout are:

critic generation prompt

You will be provided with a household task roll-out conducted by an agent and a ground truth roll-out. Your task is to write a critic of the agent's roll-out based on the **ground truth rollout** The critic should follow the form:In this task, I need do the follwing things in order:... There are ... subgoals I need to achieve,My current state achieve ...  
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.  
Important: Your critic should end with "the value is a/b=.." You can round it into 2 decimal.  
Important: You should write your critic based on given format, you should't output other things.  
Important: You shouldn't mention about ground truth rollout in your critic.

Here are k examples: <Critic\_samples>

Your task is: <relabelled\_task>

The rollout by agent is: <suboptimal\_rollout>

The **ground truth rollout** is: <gt\_rollout>

Based on the **ground truth rollout** , write a critic

Critic: