

Turn-level Multiscale Density Ratio Estimation for LLM Agents

Anonymous ACL submission

Abstract

With the rapid development of Large language model (LLM), agent systems enhanced by LLMs show huge potential in being able to deal with complex tasks, especially involving multi-step thinking or interaction with tools. For applying LLM techniques with a well-designed agent paradigm, post-training of LLM on multiple agent scenarios is necessary to achieve better performance. Among the variable post-training techniques, alignment methods such as PPO, DPO, DIL and GRPO, become popular since many papers show the significant positive impact on the model’s performance by introducing the negative samples to be punished while keeping acceptable training complexity. However, most alignment methods address simple single-turn tasks, and there remains room for improvement for complex multi-turn tasks. We propose turn-level multiscale Density Ratio Estimation (**tlm-DRE**), which assigns different weights on corresponding turns and proposes asymmetric token-level training based on the positive-negative space gaps across multiple turns of tasks. The results of the experiment on a wide range of agent benchmarks show that the proposed method performs competitively compared to traditional alignment methods. The proposed training method enables LLMs to perform robust in multi-turn reasoning tasks with both in-domain and out-of-domain conditions.

1 Introduction

Enhancing agents’ capabilities to tackle diverse complex tasks, which often involve interacting with a sophisticated environment equipped with a bunch of tools, has attracted considerable attention. For example, such tasks include complex social dialogue (Wang et al., 2023; Park et al., 2023), scientific experiment (Wang et al., 2022), embodied housework (Shridhar et al., 2020; Li et al., 2024), multi-hop question answering (Yang et al., 2018; Ho et al., 2020), etc.

To accomplish these tasks, LLM-based agents must interact with the environment step by step, decomposing the final goal into sub-goals, and then plan the trajectories based on feedback from the environment. Research on LLM-based agents initially focused on directly generating trajectories using large language models. Most studies employ prompt engineering to enhance the trajectory generation capabilities of large language models, such as CoT (Wei et al., 2022), ReAct (Yao et al., 2022), and Reflexion (Shinn et al., 2023). Subsequent research focused on trajectory tuning to further enhance the agent’s trajectory planning capabilities (Chen et al., 2023; Yin et al., 2023).

Meanwhile, reinforcement learning of Large Language Models (Achiam et al., 2023; Touvron et al., 2023; Bai et al., 2023) becomes one of the most important post-training approaches (Kumar et al., 2022) to tune LLMs more applicable to overcome shortcomings, such as hallucinations and logical consistency. Several efficient alignment training methods have been proposed in recent research discourse (Ouyang et al., 2022; Rafailov et al., 2023; Shao et al., 2024). Such post-training paradigms have also been applied to LLM-based agent tasks recently, trying to overcome the drawbacks of simply utilizing the zero-shot LLMs, which neglect agent training.

More specifically, LLM-based agent tasks typically employ the heuristic model(e.g., GPT-4) to generate a group of expert trajectories with a certain CoT form via a set of sampling strategies as a filter. Further supervised fine-tuning (SFT) is then launched to enhance the model’s reasoning and planning adaptation to certain domains. Driven by the SFT-trained reference policy, more trajectories are sampled Song et al.; Shi et al.; Xiong et al.; Kong et al. and rewarded step by step to evaluate the capability gaps of the reference model through feedback from the simulated environment. Subsequently, an alignment approach such as RHLF

*Corresponding author.

(Ouyang et al., 2022), DPO (Rafailov et al., 2023), DRE (Xiao et al., 2025), GRPO (Shao et al., 2024) will be applied as a key role of post-training in those tasks to calibrate the sampling distribution by penalizing low-quality trajectories while preserving the original probability mass over high-quality samples.

However, existing methods of alignment on agent tasks still exhibit discrepancies for future development. For example, the aforementioned RL approaches perform alignment directly at the trajectory level while being lack of attention to turn-level details, resulting in suboptimal overall alignment performance. Furthermore, even though a new LLM alignment paradigm has been put forward that views the process as a typical imitation learning under the framework of density ratio estimation, studies on the performance of the "imitation type" of alignment in agent-based tasks remain relatively scarce. Aimed to address these challenges, we propose Turn-level Multiscale Density Ratio Estimation (**tlm-DRE**) for LLM Agents, which applies imitation learning on agent-based tasks with turn-level alignment efficiency.

In particular, we assign lower weights to the turns of the samples with higher policy confidence, since there is no need to further align between chosen and rejected sample pairs. On the other hand, we assign higher weights to the turns with lower policy confidence since those turns are more crucial to be calibrated in the language space. Then, we design a novel density ratio estimation definition under turn-level varying scales based on the turn weights derived from the reference policy. Afterwards, we launch the alignment over extensive agent-based tasks in the framework of imitation learning.

Our contributions are as follows: (1) We propose a novel **tlm-DRE** method that systematically applies imitation learning to agent-based tasks in a turn-level multiscale probability space. (2) We conduct extensive experiments on several agent-based tasks (Shridhar et al., 2020; Wang et al., 2022; Yang et al., 2018) to show that our method surpasses current methods with stable and robust performance. (3) We present a comprehensive analysis to support the efficacy of our method from various points of view.

2 Related Work

LLM application in agent-based tasks The development of LLM has inspired the intelligent agents designing complex tasks involved with a dynamic environment and a complicated real-world toolbox. The main motivation to utilize LLM in those tasks is that the agents need strong abilities of reasoning(both externally-oriented and reflectional) and action planning. CoT (Wei et al., 2022) enables LLM to articulate its own thought processes, enhancing its reasoning capabilities and laying the groundwork for subsequent agent reasoning frameworks. ReAct (Yao et al., 2022) integrates feedback from the environment into its reasoning process, allowing the model to think while taking actions and interacting with the environment, and adjust subsequent actions based on that feedback. Reflexion (Shinn et al., 2023) builds upon ReAct by incorporating LLM self-reflection, allowing the model to autonomously correct its previous erroneous actions. There are even more complicated multi-agent designs (Park et al., 2023) that simulate believable human behavior in the simulated sandbox.

Alignment process: Reinforcement Learning

Many researchers focus on using reinforcement learning to further enhance agent performance. ETO (Song et al., 2024) introduced Direct Preference Optimization for post alignment. They used the SFT model as a reference policy to generate bad cases, which were then paired with expert trajectories to form sample pairs. DMPO (Shi et al., 2024) has noted that training directly with DPO leads to mismatching trajectory step lengths between positive and negative samples, and length regularization has been proposed. In contrast, SDPO (Kong et al., 2025) directly utilizes GPT to capture positive and negative samples of equal step lengths, avoiding the issue of mismatched step length.

Some research approaches agent training using online RL, such as GRPO (Shao et al., 2024), GSPO (Zheng et al., 2025). Such kind of methods are group-based RL approaches that are critic-free, making training simple and stable. GiGPO (Feng et al., 2025) groups objects that share identical environmental interaction states, employing this hierarchical Group-in-Group structure to train GRPO. These methods have also been shown to be highly effective in post-alignment.

Alignment process: Imitation Learning Other research focuses on applying imitation learning(IR)

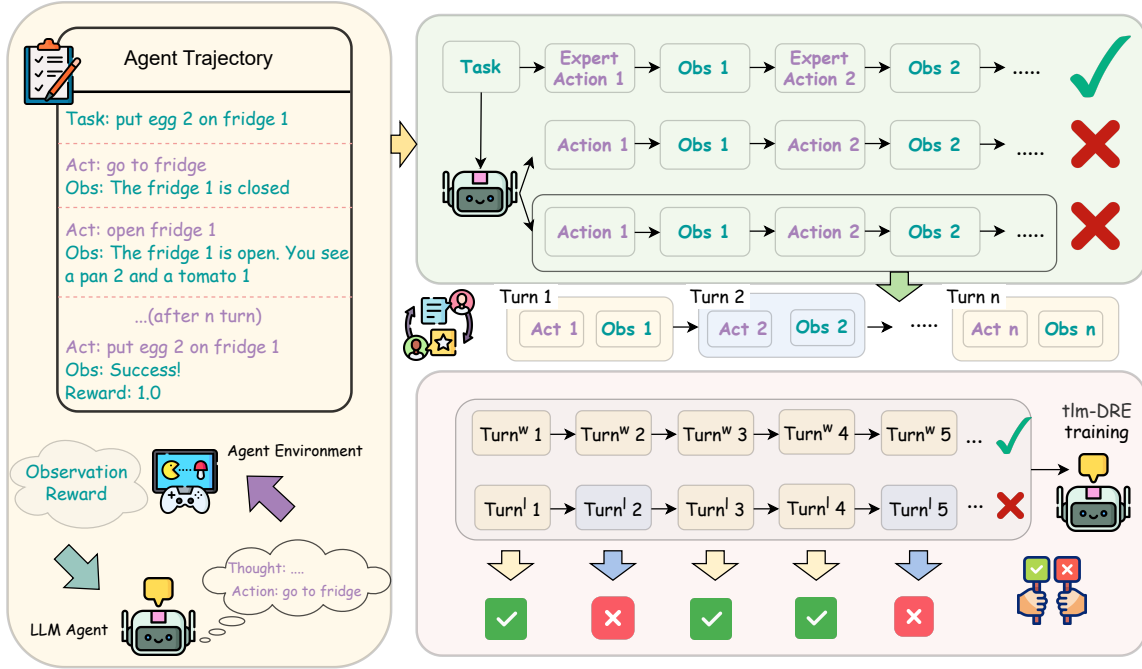


Figure 1: The overall architecture of **tlm-DRE** in a single iteration. First, the agent initialized with the SFT policy samples trajectory paths and collects failure trajectories. Then, within these failure trajectories, we identify low-confidence turns—i.e., turns where the model exhibits low self-confidence. Finally, we train the agent using Multiscale Density Ratio Estimation, which explicitly upweights these identified failure turns during optimization.

182 in the artificial intelligence domain. The key idea
 183 of IR is to directly extract knowledge from demon-
 184 strations by human experts or artificially created
 185 agents by replicating their behavior. Inverse Rein-
 186 forcement Learning has been proposed to recover
 187 the reward function of empirical resources from
 188 the uncertain environment (Russell, 1998; Ng et al.,
 189 2000; Sun and van der Schaar, 2024). Behavioral
 190 cloning (Pomerleau, 1991; Ross et al., 2011) effi-
 191 ciently bridges the domain gap by supervised fine-
 192 tuning on expert trajectories. The adversarial Imi-
 193 tation Learning method, such as GAIL (Ho and
 194 Ermon, 2016), AIRL (Fu et al., 2018), leverages
 195 adversarial training with online interaction with the
 196 environment.

197 Density Ratio Estimation (Sugiyama et al., 2012;
 198 Amari and Cichocki, 2010) can be viewed as a typi-
 199 cal way to mimic the behavior of experts under a cer-
 200 tain distribution space. DRE can be measured by
 201 optimizing the discrepancy between the learnable
 202 policy and the ideal distribution equipped with the
 203 Bregman divergence. Some methods have already
 204 tried to transplant this idea into the post training of
 205 LLM based tasks, such as DIL (Xiao et al., 2025),
 206 GSIL (Xiao et al., 2024). Furthermore, DIL (Xiao
 207 et al., 2025) also shows that traditional reinforce-

208 ment learning methods such as RLHF and DPO are
 209 actually a form of imitation learning.

210 3 Notations and Preliminaries

211 3.1 Task Formulation

212 Consider one agent which is designed to deal with a
 213 certain bunch of complex tasks interacting with the
 214 real world or simulation environment. Let $o_i \in \mathcal{O}$
 215 represent the observation of the agent for the i
 216 turn of interaction with the environment. When
 217 $i = 1$, o_i represents the initial description of a
 218 given task and the initial status of the agent. Based
 219 on the agent’s observation from the environment,
 220 as well as the previous context or trajectory $c_t =$
 221 $(o_1, a_1, \dots, o_{t-1}, a_{t-1}, o_t)$, a type of action from the
 222 action space \mathcal{A} will be chosen after a reasoning
 223 process. To avoid redundancy, we might as well
 224 view the thinking process and the action that the
 225 agent takes as a whole, named $a_t \in \mathcal{A}$ for t turns,
 226 which follows a parameterized policy $\pi_\theta(a_t|c_t)$.

227 For simplicity, the potential stochastic effect of
 228 how the action is being executed is neglected. Fur-
 229 thermore, if the policy π is driven by a certain large
 230 language model, then the action space \mathcal{A} is equiva-
 231 lent to a language space \mathcal{L} . In detail, the whole
 232 trajectory c_T can be viewed as (x, y_T) , where

$\mathbf{x} = [x_1, x_2, \dots]$ belongs to L while x_i is the language token. Similarly, $\mathbf{y}_T = [\mathbf{y}_{a_1}, \mathbf{y}_{o_2}, \mathbf{y}_{a_2}, \dots]$, $\mathbf{y}_{o_i}, \mathbf{y}_{a_i}$ also belong to the language space \mathcal{L} .

Denote π_{ref} as the reference policy. Based on the reference policy and the given input language sequences \mathbf{x} , several trajectories \mathbf{y} will be sampled following a sample strategy. Denote \mathbf{y}_w as the preferred trajectory which is marginally superior to another trajectory \mathbf{y}_l , based on the final behavior and the evaluation of the whole process. Suppose that the ideal optimized policy of the entire policy space is π_c . The object is to find the optimized projected policy π_θ for a parameterized policy space

3.2 Directed Preference Optimization

For standard reinforcement learning from human feedback (RLHF), a reward model is trained to evaluate whether the policy behaves well enough, such as PPO. (Rafailov et al., 2023) proposes a way to directly use π_θ/π_{ref} as the reward function $r(\mathbf{y}|\mathbf{x})$ and interpret the alignment process as a task to optimize the margin between the pair-wise samples

Bradley-Terry Model The PPO method applies a trained reward model to Within the framework of probability theory, the Bradley-Terry model is applied to measure the partial order relation of a sample pair $(\mathbf{y}_w, \mathbf{y}_l)$ stipulating the human preference distribution p^* as

$$p^*(\mathbf{y}_w \succ \mathbf{y}_l) = \sigma(r(\mathbf{x}, \mathbf{y}_w) - r(\mathbf{x}, \mathbf{y}_l)), \quad (1)$$

where σ is the sigmoid function.

Assuming a static dataset of pair-wise samples $\mathcal{D} = \{\mathbf{x}^i, \mathbf{y}_w^i, \mathbf{y}_l^i\}_{i=1}^N$ exists, then the DPO method will try to find a reward model r_θ in a parametrized model space that minimizes the following negative log-likelihood loss

$$\mathcal{L}_{DPO} = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l \sim \mathcal{D})} [\log \sigma(r_\theta(\mathbf{y}_w) - r_\theta(\mathbf{y}_l))] + \beta \mathbb{D}_{kl}[\pi_\theta \| \pi_{ref}]. \quad (2)$$

3.3 Density Ratio Estimation

Bregman Divergence Denote the density ratio as r_θ as $\frac{\pi_\theta}{\pi_{ref}}$, the true density ratio is r^* as $\frac{\pi_c}{\pi_{ref}}$. To employ the Bregman divergence (BR) for estimating the density ratio according to (Sugiyama et al., 2012), let f be a differentiable and strictly convex function on the density-ratio model space. the discrepancy from the true density ratio r^* to the density-ratio model r_θ can be measured as

$$BR'_f(r^* \| r_\theta) = (f(r^*) - f(r_\theta) - \partial f(r_\theta)(r^* - r_\theta)). \quad (3)$$

Fig. 4 in Appendix A illustrates how the current point r_θ can slide towards the target point r^* iteratively by optimizing the Bregman divergence. Several kernel functions f for Bregman divergence can be found in Appendix A.

Directly Imitation Learning Under the sample assumption that we have datasets \mathcal{D} defined above. The DIL method (Xiao et al., 2025) tries to imitate the optimal density ratio by minimizing the discrepancy under the Bregman divergence in a parameterized constriction

$$\min_{\theta} D_h(r^* \| r_\theta) = \sum_{\mathbf{y}} \pi_{ref}(f(r^*) - f(r_\theta) - \partial f(r_\theta)(r^* - r_\theta)). \quad (4)$$

Removing the θ irrelevant constant and under the assumption that $\mathbf{y}_l \sim \pi_{ref}(\mathbf{y}_l | \mathbf{x})$, then we get the loss function of the DIL method

$$\mathcal{L}_{DIL} = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l \sim \mathcal{D})} \{ \partial f(r_\theta(\mathbf{y}_w)) - [\partial f(r_\theta(\mathbf{y}_l)) r_\theta(\mathbf{y}_l) - f(r_\theta(\mathbf{y}_l))] \}. \quad (5)$$

4 Methodology

Due to varying degrees of training exposure of tokens in the embedding space within the reference model, one daunting challenge is how to train the policy over the candidate tokens more efficiently. This problem warrants more attention for the tasks that need multi-step thinking and reasoning. (Liu et al., 2025) introduced a way to allocate token weights to each token according to a certain importance measurement for the DPO. However, this idea cannot be trivially generalized to turn-level optimization for the DRE loss due to the fact that most of the Bregman divergence functions show a high degree of asymmetry and nonlinearity. We will introduce a new density ratio function that will be used to solve this problem. Our method framework is illustrated in Figure 1.

4.1 Multiscale Density Ratio Representation

Since the density ratio is designed to estimate how close an approximation to the optimal solution can be achieved over the output language space, we can only focus on the linguistic action results \mathbf{y}_{a_i} under the assumption that the observation results are determined by omitting the stochastic effect from the simulation environment. More specifically, the

density ratio can be expressed as

$$\begin{aligned} r(\mathbf{y}_t|\mathbf{x}) &= \frac{\pi(\mathbf{y}_t|\mathbf{x})}{\pi_{ref}(\mathbf{y}_t|\mathbf{x})} \\ &= \frac{\pi(\mathbf{y}_{a_t}|\mathbf{x}, \mathbf{c}_{t-1}) \dots \pi(\mathbf{y}_{a_1}|\mathbf{x})}{\pi_{ref}(\mathbf{y}_{a_t}|\mathbf{x}, \mathbf{c}_{t-1}) \dots \pi_{ref}(\mathbf{y}_{a_1}|\mathbf{x})}, \end{aligned} \quad (6)$$

where $\mathbf{c}_{t-1} = [\mathbf{y}_{a_1}, \dots, \mathbf{y}_{o_{t-1}}]$

The main drawback of the traditional expression of the DRE is the inefficiency of training the underfitting, while well-trained tokens will possibly suffer the risk of overfitting. Meanwhile, as demonstrated in (Liu et al., 2025), the importance of each token in latent embedding spaces is not uniformly distributed, implying heterogeneous training demands. In the agent-related tasks with multi-step reasoning, the importance weighting with non-uniformity becomes more pronounced from the turn-level perspective.

Suppose that a turn-level weighted function

$$\omega_t = \omega(\mathbf{y}_{a_t}|\mathbf{x}, \mathbf{c}_{t-1}), \quad (7)$$

that measures how important and how well-trained each turn in the agentic reasoning process is. Then the density ratio is to be measured as

$$\begin{aligned} r(\mathbf{y}_t|\mathbf{x}) &= \frac{\hat{\pi}(\mathbf{y}_t|\mathbf{x})}{\hat{\pi}_{ref}(\mathbf{y}_t|\mathbf{x})} \\ &= \frac{\pi^{\omega_t}(\mathbf{y}_{a_t}|\mathbf{x}, \mathbf{c}_{t-1}) \dots \pi^{\omega_1}(\mathbf{y}_{a_1}|\mathbf{x})}{\pi_{ref}^{\omega_t}(\mathbf{y}_{a_t}|\mathbf{x}, \mathbf{c}_{t-1}) \dots \pi_{ref}^{\omega_1}(\mathbf{y}_{a_1}|\mathbf{x})}. \end{aligned} \quad (8)$$

4.2 Turn-level Weighted Measurement

We define the adequacy estimate of the i -th turn output \mathbf{y}_i as p_i :

$$\begin{aligned} \log p_i &= \frac{1}{|\mathbf{y}_i|} \log \pi_{ref}(\mathbf{y}_{a_i}|\mathbf{x}, \mathbf{c}_{i-1}) \\ &= \frac{1}{|\mathbf{y}_i|} \sum_{j=1}^{|\mathbf{y}_i|} \log \pi_{ref}(y_{i,j}|\mathbf{x}, \mathbf{c}_{i-1}, y_{i,<j}), \end{aligned} \quad (9)$$

where $|\mathbf{y}_i|$ is the total number of tokens of the \mathbf{y}_i .

In some previous searches (Nagumo and Fujisawa, 2024), they view the statistical possibility of the output as an outlier measurement. However, we argue that, under the assumption that previous training in the agent-based task is inefficient and imbalanced, the adequacy estimation actually shows more weight on how well trained the i th turn of the trajectory based on the reference policy is, rather than the outliers measurement. Hence, a natural

point is to make sure that there is a negative correlation between the adequacy estimation and the turn-weights applied on the following alignment.

We define the turn-level weights ω_i for the i -th turn of the trajectory as

$$\omega_i = \begin{cases} \omega_L & p_i > p_{pivot}, \\ \omega_U & p_i \leq p_{pivot}. \end{cases} \quad (10)$$

4.3 Turn-level Weighted DRE Optimization

Applying the turn-level weights version of density ratio as defined in Equation (8) on the DRE optimization, Equation (4) becomes

$$\begin{aligned} \min_{\theta} D_{\omega, f}(\hat{r}^*|\hat{r}_{\theta}) &= \sum_{\mathbf{y}} \pi_{ref}(f(\hat{r}^*) \\ &\quad - f(\hat{r}_{\theta}) - \partial f(\hat{r}_{\theta})(\hat{r}^* - \hat{r}_{\theta})). \end{aligned} \quad (11)$$

Subtracting the constant $\sum_{\mathbf{y}} \pi_{ref} f(\hat{r}^*)$, we obtain

$$\begin{aligned} \min_{\theta} D_{\omega, f}(\hat{r}^*|\hat{r}_{\theta}) &= - \sum_{\mathbf{y}} \pi_{ref} \partial f(\hat{r}_{\theta}) \hat{r}^* \\ &\quad + \sum_{\mathbf{y}} \pi_{ref} (\partial f(\hat{r}_{\theta}) \hat{r}_{\theta} - f(\hat{r}_{\theta})). \end{aligned} \quad (12)$$

The last term above $\pi_{ref} \partial f(\hat{r}_{\theta}) \hat{r}^*$ can be expanded as

$$\begin{aligned} &\pi_{ref} \partial f(\hat{r}_{\theta}) \hat{r}^* \\ &= \pi_{ref} \partial f(\hat{r}_{\theta}) \prod_{i=1}^t \frac{\pi_c^{\omega_i}(\mathbf{y}_{a_i}|\mathbf{x}, \mathbf{c}_{i-1})}{\pi_{ref}^{\omega_i}(\mathbf{y}_{a_i}|\mathbf{x}, \mathbf{c}_{i-1})} \\ &= \partial f(\hat{r}_{\theta}) \prod_{i=1}^t \frac{\pi_c^{\omega_i}(\mathbf{y}_{a_i}|\mathbf{x}, \mathbf{c}_{i-1})}{\pi_{ref}^{\omega_i-1}(\mathbf{y}_{a_i}|\mathbf{x}, \mathbf{c}_{i-1})}. \end{aligned} \quad (13)$$

Due to the fact that the possibility of the un-chosen trajectory by policy π_c is vanishingly small and $\omega \in [\omega_L, \omega_U]$, Substituting (13) into (12) the equation yields

$$\begin{aligned} &\min_{\theta} D_{\omega, f}(\hat{r}^*|\hat{r}_{\theta}) \\ &= \sum_{\mathbf{y}} \pi_{ref} (\partial f(\hat{r}_{\theta}) \hat{r}_{\theta} - f(\hat{r}_{\theta})) \\ &\quad - \sum_{\mathbf{y}_w} \pi_c \partial f(\hat{r}_{\theta}) \prod_{i=1}^t \frac{\pi_c^{\omega_i-1}(\mathbf{y}_{a_i}|\mathbf{x}, \mathbf{c}_{i-1})}{\pi_{ref}^{\omega_i-1}(\mathbf{y}_{a_i}|\mathbf{x}, \mathbf{c}_{i-1})}, \end{aligned} \quad (14)$$

under the same assumption that $\mathbf{y}_l \sim \pi_{ref}(\mathbf{y}_l|\mathbf{x})$ the loss function of (tlm-DRE) is

$$\begin{aligned} &\mathcal{L}_{tlm-DRE} \\ &= \int_{\mathbf{y}_l} [\partial f(\hat{r}_{\theta}(\mathbf{y}_l)) \hat{r}_{\theta}(\mathbf{y}_l) - f(\hat{r}_{\theta}(\mathbf{y}_l))] \\ &\quad - \int_{\mathbf{y}_w} \partial f(\hat{r}_{\theta}) \prod_{i=1}^t \frac{\pi_c^{\omega_i-1}(\mathbf{y}_{a_i}|\mathbf{x}, \mathbf{c}_{i-1})}{\pi_{ref}^{\omega_i-1}(\mathbf{y}_{a_i}|\mathbf{x}, \mathbf{c}_{i-1})}. \end{aligned} \quad (15)$$

Since the distribution shows significant concentration on the candidate tokens from the chosen trajectory over the language space, we can assume that $\pi_c(\mathbf{y}_{a_i} \parallel \mathbf{x}, \mathbf{c}_{i-1}) \sim 1$ uniformly, leading to the final **tlm-DRE** loss for training

$$\begin{aligned} \mathcal{L}_{tlm-DRE} &= \mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l \sim \mathcal{D})} \{ [\partial f(\hat{r}_\theta(\mathbf{y}_l)) \hat{r}_\theta(\mathbf{y}_l) - f(\hat{r}_\theta(\mathbf{y}_l))] \\ &\quad - \partial f(\hat{r}_\theta(\mathbf{y}_w)) \prod_{i=1}^t \frac{\pi_c^{\omega_i-1}(\mathbf{y}_w, \mathbf{a}_i \parallel \mathbf{x}, \mathbf{c}_{w,i-1})}{\pi_{ref}^{\omega_i-1}(\mathbf{y}_w, \mathbf{a}_i \parallel \mathbf{x}, \mathbf{c}_{w,i-1})} \} \\ &\approx \mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l \sim \mathcal{D})} \{ [\partial f(\hat{r}_\theta(\mathbf{y}_l)) \hat{r}_\theta(\mathbf{y}_l) - f(\hat{r}_\theta(\mathbf{y}_l))] \\ &\quad - \partial f(\hat{r}_\theta(\mathbf{y}_w)) \prod_{i=1}^t \frac{1}{\pi_{ref}^{\omega_i-1}(\mathbf{y}_w, \mathbf{a}_i \parallel \mathbf{x}, \mathbf{c}_{w,i-1})} \}. \end{aligned} \quad (16)$$

In this paper, we use UKL (Nguyen et al., 2010) defined in Appendix A as the kernel function for Bregman divergence. Then equation (16) will be as

$$\begin{aligned} \mathcal{L}_{tlm-DRE} &= \mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l \sim \mathcal{D})} \left\{ \prod_i \frac{\pi_\theta^{\omega_{l,i}}(\mathbf{y}_{l,i})}{\pi_{ref}^{\omega_{l,i}}(\mathbf{y}_{l,i})} - \right. \\ &\quad \left. \sum_i \frac{\omega_{w,i}}{\prod_{i=1}^t \pi_{ref}^{\omega_i-1}(\mathbf{y}_w, \mathbf{a}_i)} \log \frac{\pi_\theta(\mathbf{y}_w, i)}{\pi_{ref}(\mathbf{y}_w, i)} \right\}. \end{aligned} \quad (17)$$

5 Experiments

In this section, we demonstrate the performance of **tlm-DRE** across various agent based tasks, provide detailed experimental procedures, and introduce other related baselines.

5.1 Experimental Settings

Datasets We perform experiments across a diverse set of environments to evaluate the capabilities of our agent. Specifically, we use SciWorld (Wang et al., 2022) under the Apache-2.0 license for grounded scientific experimentation in simulated labs, ALFWorld under the MIT License (Shridhar et al., 2020) for embodied household tasks requiring object manipulation and planning in 3D environments. In addition, we also include HotpotQA (Yang et al., 2018) under the CC BY-SA 4.0 license as the benchmark of multi-hop reasoning in open-domain settings, a task that requires the agent to retrieve and integrate information from multiple sources to answer complex questions. SciWorld provides dense final rewards on a continuous scale from 0 to 1, whereas ALFWorld offers only sparse binary rewards that indicate whether the task

was completed or not. For multi-hop QA tasks, we measure reasoning performance using EM and F1 scores against ground-truth answers.

Baselines We compare our approach with a series of benchmarks: (1) Zero-shot using LLMs such as GPT-4o, Qwen2.5-7B, applying the ReAct prompting paradigm, which represents the state-of-the-art zero-shot capabilities of LLMs. (2) SFT (Supervised Fine-Tuning) conducts behavioral cloning on expert trajectories. (3) PPO (Proximal Policy Optimization) as an actor-critic reinforcement learning algorithm to directly optimize the SFT-initialized policy. (4) ETO (Song et al., 2024) uses successful and failure trajectories as sample pairs for DPO training. (5) DMPO (Shi et al., 2024) adds length regularization to ETO to eliminate noise caused by failure trajectories with excessive steps. (6) GiGPO (Feng et al., 2025) groups objects with identical environmental interaction states, employing the hierarchical group structure to train GRPO.

All experiments were conducted on 8x Nvidia A100 GPUs, each with 80GB of memory, implemented using PyTorch in Python. Details of the experiment setup can be found in Appendix B.

5.2 Main Results

Results on Interactive Tasks Table 1 demonstrates the strong performance of **tlm-DRE** on both ALFWorld and SciWorld. As shown, prompt-only ReAct baselines achieve only moderate results: GPT-4 attains 38.1 average reward on ALFWorld (unseen) and 64.4 on SciWorld (unseen), while GPT-3.5 lags significantly behind. Qwen2.5-7B performs modestly in all settings, with an average reward of around 25, indicating limited effectiveness without further alignment or training.

For SFT and RL training, most prior work adopts LLaMA2-7B as the backbone. To ensure a maximally fair comparison, we also report results with LLaMA2-7B. We further evaluate a stronger and more recent model, Qwen2.5-7B. ETO directly applies DPO to agent tasks, avoiding the need for a critic network as in PPO and thus offering a lighter and simpler training pipeline. With LLaMA2-7B, ETO achieves 72.4 score on ALFWorld (*unseen*) and 61.1 on SciWorld (*unseen*), though it is weaker in the seen setting. DMPO performs particularly well on SciWorld with LLaMA2-7B, reaching 72.4 on SciWorld (*seen*), but degrades on ALFWorld. Based on the same base model, our method attains the best results in the seen tasks on both datasets;

Paradigm	Models	ALFWorld		ScienceWorld	
		Seen	Unseen	Seen	Unseen
Prompt-based	GPT-3.5 (Ouyang et al., 2022)	7.9	10.5	16.5	13.0
	GPT-4 (Achiam et al., 2023)	42.9	38.1	64.8	64.4
	Qwen2.5-7B (Team et al., 2024)	25.1	28.4	26.8	25.2
<i>Llama-2-7B-Chat</i>					
SFT & RL	SFT (Chen et al., 2023)	60.0	67.2	56.8	56.0
	PPO (Trung et al., 2024)	22.1	29.1	59.4	51.7
	RFT (Zhang et al., 2023)	62.9	66.4	71.6	54.3
	ETO (Song et al., 2024)	68.6	72.4	68.5	61.1
	DMPO (Shi et al., 2024)	43.3	55.0	72.4	61.7
	tlm-DRE (ours)	70.7	72.4	73.0	61.2
<i>Qwen2.5-7B-instruct</i>					
RL-training	SFT (Chen et al., 2023)	70.7	83.6	71.8	61.8
	ETO (Song et al., 2024)	75.0	86.6	69.1	62.8
	tlm-DRE w/o DRE	75.0	86.6	72.1	63.3
	tlm-DRE w/o tlm	75.7	88.8	71.9	63.7
	tlm-DRE (ours)	75.7	90.3	73.1	63.8

Table 1: Performance of different methods on ALFWorld and SciWorld, reported as average reward. "Seen" denotes the held-out test set containing task types observed during training, while "Unseen" refers to test tasks with critical unseen variations (e.g., novel object or goal). "tlm-DRE w/o DRE" denotes the linearly turn-weighted DPO training. For fair comparison, all methods use the same base models: LLaMA2-7B and Qwen2.5-7B.

Method	ALFWorld (all)
RLOO (Ahmadian et al., 2024)	75.5
GRPO (Shao et al., 2024)	77.6
GiGPO w/ std (Feng et al., 2025)	90.8
GiGPO w/o std (Feng et al., 2025)	90.2
tlm-DRE (ours)	92.4

Table 2: For fair comparison, results are evaluated in the GiGPO (Feng et al., 2025) test environment (with longer exploration steps and different test categories). All methods use the same base models: Qwen2.5-7B-Instruct

Method	HotpotQA	
	EM	F1
SFT (Chen et al., 2023)	27.80	36.45
CoH (Liu et al., 2023)	28.60	39.53
PPO (Trung et al., 2024)	28.20	36.47
DPO (Song et al., 2024)	26.40	34.83
NAT (Wang et al., 2025)	29.60	42.50
tlm-DRE (ours)	33.8	43.74

Table 3: Overall results on open-domain question answering tasks. We measure the performance using Exact Match and F1 score. For fair comparison, all methods use the same base models: LLaMA2-7B

461 moreover, when switching to Qwen2.5-7B, it further
462 boosts ALFWorld (*unseen*) to 90.3, far surpassing
463 the LLaMA2-7B counterpart.

464 GiGPO uses a different evaluation protocol with
465 longer exploration steps. For a fair comparison,
466 we also evaluated our method in the GiGPO en-
467 vironment. As shown in Table 2, our approach
468 outperforms GiGPO by approximately three points
469 in overall success reward, demonstrating strong
470 effectiveness in this more challenging setting.

471 **Results on Multi-Hop QA Tasks** As shown in
472 Table 3, **tlm-DRE** consistently improves perfor-
473 mance on multi-turn search-augmented QA tasks,

474 achieving an Exact Match (EM) of 33.8 and an F1
475 score of 43.74 on HotpotQA, substantially outper-
476 forming strong baselines such as NAT. Although
477 search-augmented QA uses a different set of tools
478 and typically requires fewer exploration steps than
479 traditional agent tasks (e.g., ALFWorld), it implies
480 generalization of our method.

481 **Ablation Study** We conduct ablation studies to
482 evaluate the effectiveness of each component: (1)
483 Density Ratio Estimation (DRE) with UKL-kernel
484 for Bregman divergence. (2) A variant version of
485 DRE proposed with turn-level multiscale. Table 1
486 shows that both components yield clear improve-

ments over standard SFT and DPO. Specifically, turn-level weights and DRE each contribute independently, but exhibit different strengths across evaluation splits. We observe that the turn-level weights tend to provide larger gains in the *seen* setting on both datasets, suggesting that emphasizing critical turns helps better fit in-distribution interaction patterns. Moreover, **tlm-DRE** brings more consistent improvements in *unseen* tasks and achieves the best overall performance.

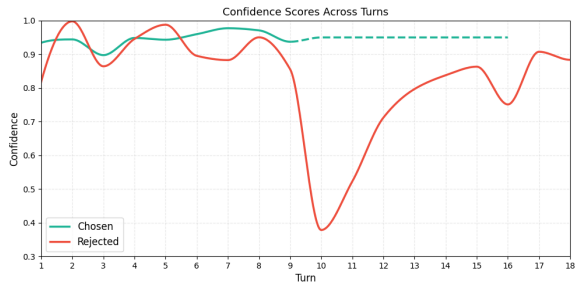


Figure 2: Illustration of confidence scores for different turns computed using the SFT policy, shown for both the chosen and rejected trajectories.

6 Analysis

6.1 Illustration of turn weights

Previous methods, such as ETO and DMPO, often view the entire trajectory as a unified object, thereby neglecting the contribution gaps of different turns within the trajectory. For example, in a failed trajectory, not every turn is necessarily incorrect. Therefore, as shown in Eq. (9), we use the SFT policy to recompute the confidence scores of the model in different turns for both positive and negative examples. As illustrated in Fig. 2, we present the confidence scores of the SFT policy on different turns for both chosen expert trajectories and rejected trajectories. We observe that, in the rejected trajectories, the model exhibits notably low confidence at Turns 10 and 11. By examining the chosen expert trajectories and specific case studies of Fig. 5, we find that these two turns are indeed critical steps that lead the entire trajectory astray. Consequently, in the subsequent alignment stage, we place particular emphasis on such uncertain turns, and our experimental results confirm that this idea is effective.

6.2 Margin analysis

Fig. 3 shows the reward margin between the chosen and rejected trajectories during training for tlm-

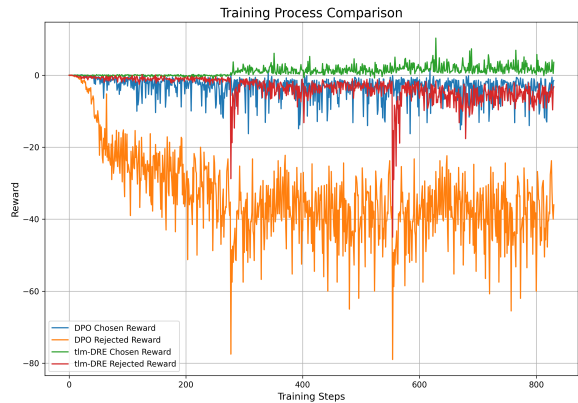


Figure 3: Illustration of the reward margin between the chosen and rejected trajectories during training for tlm-DRE and DPO

DRE and DPO. We argue that DPO separates chosen and rejected trajectories in a coarse-grained manner, treating a failed trajectory as entirely negative and neglecting that most turns within it can still be correct. As a result, many correct turns from rejected trajectories may be pushed toward the negative region, which can degrade overall performance. In practice, the failure trajectory is often caused by only a few critical erroneous steps.

To address this issue, our method up-weights critical turns and implicitly down-weights correct turns to be punished. Concretely, we only decrease the reward for the key erroneous turns while keeping the weights of other turns almost unchanged. As illustrated in Fig. 3, our method produces a chosen-rejected margin smaller than the standard DPO, reflecting a finer-grained distinction between positive and negative samples. The empirical results further support the effectiveness of this design.

7 Conclusions

In this paper, Turn-level Multiscale Density Ratio Estimation (**tlm-DRE**) is proposed to enhance LLM performance in multi-turn agent tasks. We theoretically design a multiscale density ratio representation that fully leverages the contrastive information between positive and negative samples. Our approach assigns turn-specific weights under the framework of DRE to make the alignment of multi-turn tasks more flexible and robust. Our approach has demonstrated strong performance across multiple multi-turn agent tasks, including ALFWorld, SciWorld, and HotpotQA. In the future, we will explore the generalizability of this method across more multi-turn scenarios.

8 Limitations

Our Turn-level Multiscale Density Ratio Estimation (**tlm-DRE**) employs pair-level asymmetric training with turn weights assigned to key steps, which adapts well to multi-turn agent tasks with long trajectories. However, this study is still limited from several perspectives, pointing to promising directions for future research. First, **tlm-DRE** is conducted mainly based on the offline sampling strategy, traditionally adopted by PPO or DPO. Whether the gain is still maintained when we combine the proposed alignment method with the dynamic sampling strategy as used in GRPO or GSPO deserves careful investigation. Second, while the UKL kernel of Bregman divergence shows advances in several multi-turn agent tasks, systematic analysis of the impact of varying the choice of different function kernels needs to be further conducted. Moreover, with the rapid development of the AI/LLM agent area, there are more agent-based tasks emerging in recent years, which need more complex task decomposition and tool-using strategy, as well as more sophisticated environmental interaction. Extending the experiment of our method to more relevant tasks also warrants an in-depth analysis. We hope our work will inspire researchers to explore multi-step agent training in this field.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. 2024. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv preprint arXiv:2402.14740*.

Shun-ichi Amari and Andrzej Cichocki. 2010. Information geometry of divergence functions. *Bulletin of the polish academy of sciences. Technical sciences*, 58(1):183–195.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, and 1 others. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

Ayanendranath Basu, Ian R Harris, Nils L Hjort, and MC Jones. 1998. Robust and efficient estimation by minimising a density power divergence. *Biometrika*, 85(3):549–559.

Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao. 2023. Fireact: Toward language agent fine-tuning. *arXiv preprint arXiv:2310.05915*.

Lang Feng, Zhenghai Xue, Tingcong Liu, and Bo An. 2025. Group-in-group policy optimization for llm agent training. *arXiv preprint arXiv:2505.10978*.

Justin Fu, Katie Luo, and Sergey Levine. 2018. [Learning robust rewards with adversarial inverse reinforcement learning](#). In *International Conference on Learning Representations*.

Trevor Hastie. 2009. The elements of statistical learning: data mining, inference, and prediction.

Jonathan Ho and Stefano Ermon. 2016. Generative adversarial imitation learning. *Advances in neural information processing systems*, 29.

Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. [Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Takafumi Kanamori, Shohei Hido, and Masashi Sugiyama. 2009. A least-squares approach to direct importance estimation. *The Journal of Machine Learning Research*, 10:1391–1445.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Aobo Kong, Wentao Ma, Shiwan Zhao, Yongbin Li, Yuchuan Wu, Ke Wang, Xiaoqian Liu, Qicheng Li, Yong Qin, and Fei Huang. 2025. Sdpo: Segment-level direct preference optimization for social agents. *arXiv preprint arXiv:2501.01821*.

Komal Kumar, Tajamul Ashraf, Omkar Thawakar, Rao Muhammad Anwer, Hisham Cholakkal, Mubarak Shah, Ming-Hsuan Yang, Phillip HS Torr, Fahad Shahbaz Khan, and Salman Khan. 2022. Llm post-training: a deep dive into reasoning large language models (2025). *URL https://arxiv.org/abs/2502.21321*, 3(7).

Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen, Sanjana Srivastava, Roberto Martín-Martín, Chen Wang, Gabriel Levine, Wensi Ai, Benjamin Martinez, and 1 others. 2024. Behavior-1k: A human-centered, embodied ai benchmark with 1,000 everyday activities and realistic simulation. *arXiv preprint arXiv:2403.09227*.

Aiwei Liu, Haoping Bai, Zhiyun Lu, Yanchao Sun, Xiang Kong, Xiaoming Wang, Jiulong Shan, Albin Madappally Jose, Xiaojiang Liu, Lijie Wen, Philip Yu, and Meng Cao. 2025. [Tis-dpo: Token-level importance sampling for direct preference optimization with estimated weights](#). In *International Conference*

771 and 1 others. 2023. Rolellm: Benchmarking, elic-
772 iting, and enhancing role-playing abilities of large
773 language models. *arXiv preprint arXiv:2310.00746*.

774 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
775 Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,
776 and 1 others. 2022. Chain-of-thought prompting elic-
777 its reasoning in large language models. *Advances*
778 *in neural information processing systems*, 35:24824–
779 24837.

780 Teng Xiao, Mingxiao Li, Yige Yuan, Huaisheng Zhu,
781 Chao Cui, and Vasant G Honavar. 2024. How to
782 leverage demonstration data in alignment for large
783 language model? a self-imitation learning perspec-
784 tive. *arXiv preprint arXiv:2410.10093*.

785 Teng Xiao, Yige Yuan, Mingxiao Li, Zhengyu Chen,
786 and Vasant G Honavar. 2025. [On a connection be-
787 tween imitation learning and RLHF](#). In *The Thir-
788 teenth International Conference on Learning Repre-
789 sentations*.

790 Weimin Xiong, Yifan Song, Xiutian Zhao, Wenhao Wu,
791 Xun Wang, Ke Wang, Cheng Li, Wei Peng, and Su-
792 jian Li. 2024. Watch every step! llm agent learning
793 via iterative step-level process refinement. *arXiv*
794 *preprint arXiv:2406.11176*.

795 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,
796 William Cohen, Ruslan Salakhutdinov, and Christo-
797 pher D Manning. 2018. Hotpotqa: A dataset for
798 diverse, explainable multi-hop question answering.
799 In *Proceedings of the 2018 conference on empiri-
800 cal methods in natural language processing*, pages
801 2369–2380.

802 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak
803 Shafran, Karthik R Narasimhan, and Yuan Cao. 2022.
804 React: Synergizing reasoning and acting in language
805 models. In *The eleventh international conference on*
806 *learning representations*.

807 Da Yin, Faeze Brahman, Abhilasha Ravichander, Khy-
808 athi Chandu, Kai-Wei Chang, Yejin Choi, and
809 Bill Yuchen Lin. 2023. Agent lumos: Unified and
810 modular training for open-source language agents.
811 *arXiv preprint arXiv:2311.05657*.

812 Yifan Zhang, Jingqin Yang, Yang Yuan, and An-
813 drew Chi-Chih Yao. 2023. Cumulative reason-
814 ing with large language models. *arXiv preprint*
815 *arXiv:2308.04371*.

816 Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui
817 Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong
818 Liu, Rui Men, An Yang, and 1 others. 2025.
819 Group sequence policy optimization. *arXiv preprint*
820 *arXiv:2507.18071*.

A Details of Density Ratio Estimation

Fig. 4 shows how optimizing the Bregman divergence gradually drives the point toward the target point.

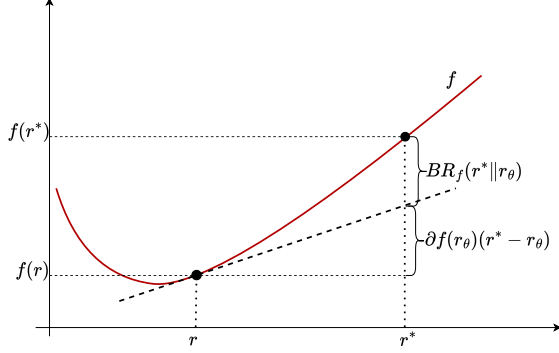


Figure 4: Illustration of DRE with Bregman divergence.

Several kernel functions for Bregman divergence have been proposed in the past discourse, such as LSIF (Kanamori et al., 2009), UKL (Nguyen et al., 2010), BCE (Hastie, 2009), and Basu’s power divergence (Basu et al., 1998). We list the details of those functions below.

LSIF

$$f(r) = \frac{1}{2}(r - 1)^2. \quad (18)$$

Bregman divergence (BR) defined in Equation 3 is reduced to the squared (SQ) distance:

$$SQ'(r^* || r_\theta) = \frac{1}{2}(r^* - r)^2. \quad (19)$$

UKL

$$f(r) = r \log r - r. \quad (20)$$

BR is reduced to the unnormalized Kullback–Leibler (UKL) divergence:

$$UKL'(r^* || r_\theta) = r^* \log \frac{r^*}{r} - r^* + r. \quad (21)$$

BCE/BKL

$$f(r) = r \log r - (r + 1) \log(r + 1). \quad (22)$$

BR is reduced to the binary Kullback–Leibler (BKL) divergence:

$$BKL'(r^* || r_\theta) = (1 + r^*) \log \frac{1 + r}{1 + r^*} + r^* \log \frac{r}{r^*}. \quad (23)$$

Basu’s power For $\alpha > 0$

$$f(r) = \frac{r^{1+\alpha} - r}{\alpha}. \quad (24)$$

Then BR is reduced to the BA divergence:

$$BA'_\alpha(r^* || r_\theta) = r^\alpha(r - r^*) - \frac{r^* r^\alpha - (r^*)^{1+\alpha}}{\alpha}. \quad (25)$$

B Experiment Setup

We mainly select Qwen2.5-7B-Instruct (Team et al., 2024) for experiment. For a comprehensive comparison, we also selected Llama2-7B-Chat (Touvron et al., 2023). Our model is fully fine-tuned (not PEFT) in two stages: 3 epochs of SFT followed by 1 epoch of alignment training, optimizing with AdamW (Kingma and Ba, 2014). For SFT, initial learning rate is 1×10^{-6} for Alfworld, Sciworld and 2×10^{-5} for HotpotQA. For alignment, initial learning rate is 7×10^{-7} for Alfworld and 1×10^{-6} for Sciworld, HotpotQA. Training uses a batch size of 4. In the alignment stage, policy sampling uses temperature 1; Inference testing uses temperature 0. Contrastive pairs are filtered by score reward with threshold $\tau = 0.4$. The pivot value in Equation 10 is 0.9. The lower bound and upper bound of the turn weight in Equation 10 are 0.2 and 1.0.

C Trajectory Sample

Fig. 5 shows how different turns in a trajectory vary in their importance and completeness of training.

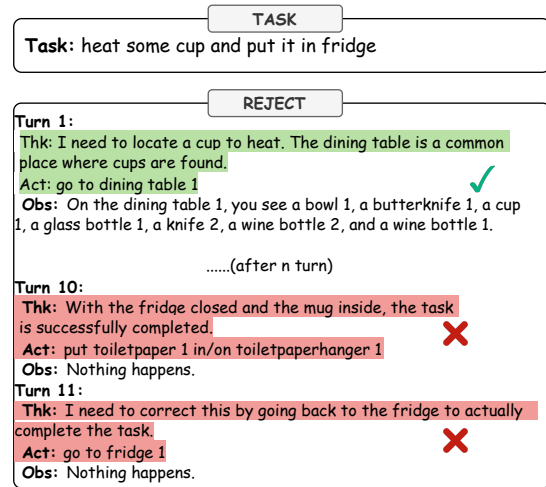


Figure 5: Case analysis across different turns.

Fig. 6 and 7 show how the key turn improved after alignment plays an important role for the task.

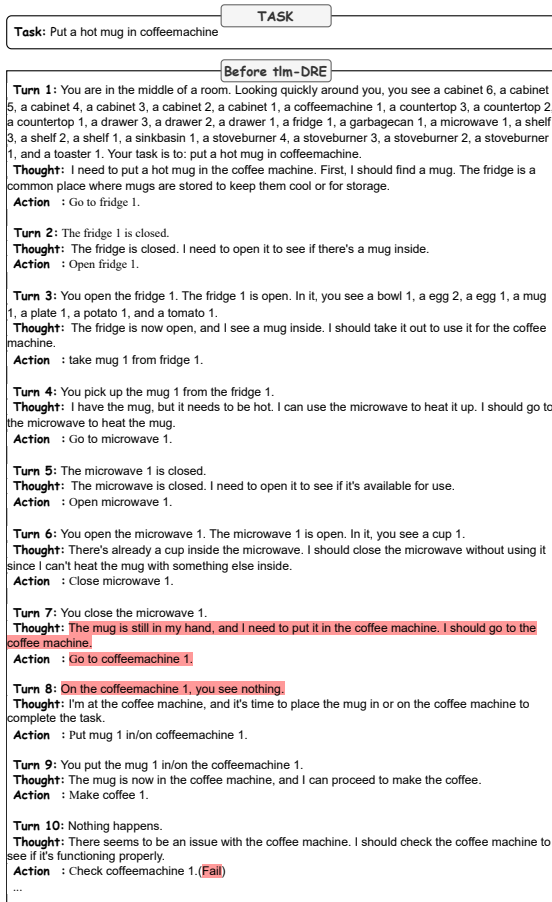


Figure 6: Trajectory sample before tlm-DRE.

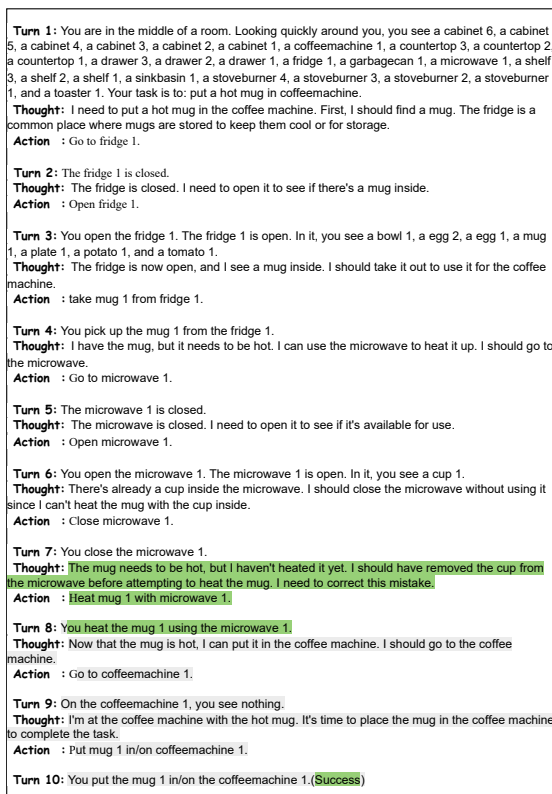


Figure 7: Trajectory sample after tlm-DRE.