

Unleashing Interactive Agent Planning Ability in LLMs via Multi-Trajectory Reinforcement Learning

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have demonstrated remarkable capabilities across diverse domains, yet they struggle with agent task planning in dynamic environments requiring continuous observation and sequential decision-making. Current methods generate static action sequences from pre-trained knowledge without learning from environmental feedback, limiting their effectiveness in partially observable settings. We present **Interactive Planner-R1**, a novel trajectory-level reinforcement learning framework that enables LLMs to develop interactive planning capabilities through autonomous environmental exploration. Our approach addresses three key challenges: (1) limited exploration diversity by introducing multi-trajectory autonomous exploration through parallel group rollouts, (2) sparse reward signals by developing a completion-driven reward architecture that promotes genuine environmental understanding, and (3) single-step optimization constraints by proposing Interactive Policy Optimization (IPO) that extends group-relative policy optimization for multi-step trajectory learning. Extensive experiments on ALFWorld and ScienceWorld demonstrate that Interactive Planner-R1 achieves substantial improvements over existing approaches, reaching 97.55% completion rate on ALFWorld and 79.92% on ScienceWorld, with strong generalization exhibiting only 3.33% performance gap in unseen environments. Our work establishes a new paradigm for LLM-based interactive planning through trajectory-level policy learning.

1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities across a wide range of tasks, from open-domain dialogue to complex reasoning (Brown et al., 2020; OpenAI, 2023; DeepSeek-AI et al., 2025; Shao et al., 2024; Team et al., 2024; Yang et al., 2024a). However, translating these capabilities into effective planning in in-

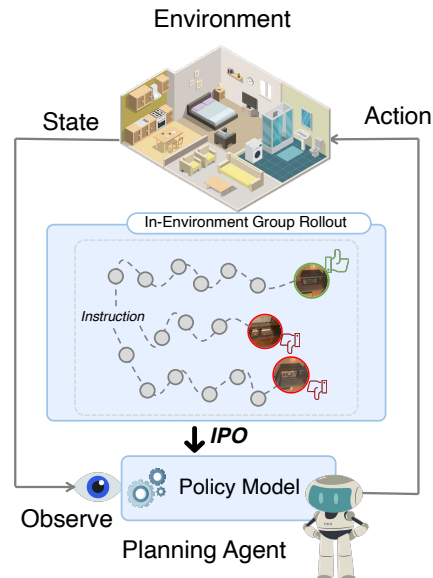


Figure 1: Overview of Interactive Planner-R1, which encourages agents to autonomously explore environments and trains them with binary rewards based on task completion status.

teractive environments remains a fundamental challenge. Planning agents must continually perceive, reason about, and act upon their environment, forming a tightly coupled decision-making loop. Current LLMs still lack the situational and interactive awareness required for robust planning. In contrast, humans actively probe their surroundings, learning how actions transform the environment and refining both perception and world understanding through exploration (Gibson, 1979). While intuitive for humans, such exploration-driven learning remains challenging for artificial intelligence—especially when applying LLMs to interactive planning domains, such as robotics and embodied agents.

Traditional LLM-based planners treat the model as a static decision-maker (Chen et al., 2023b; Min et al., 2022), generating fixed action sequences from pre-trained knowledge and ignoring real-time

feedback. While Yao et al. (2023) improved task performance by introducing reasoning processes into planning, these methods still primarily perform offline rollouts based on prior knowledge, failing to establish crucial causal links between actions and environmental feedback during interaction. This limitation is particularly problematic in partially observable or unfamiliar environments, where agents must continually update their world models to succeed. Recent work (Song et al., 2024; Wang et al., 2024b; Chen et al., 2023a; Zeng et al., 2024a) has attempted to bridge this gap by incorporating human-curated demonstrations or hand-crafted knowledge. However, heavy reliance on human priors restricts scalability and limits the overall capability of LLMs. This raises a fundamental question: Can we instead encourage autonomous exploration with minimal supervision, enabling LLMs to acquire robust interaction and generalization abilities?

Inspired by human cognition, we introduce **Interactive Planner-R1**, a reinforcement learning framework that enables LLMs to learn interactive planning through environmental exploration, requiring only minimal task-specific supervision.

Extensive experiments demonstrate the effectiveness of Interactive Planner-R1, achieving 97.55% and 79.92% cumulative reward on the ALFWorld and ScienceWorld benchmarks, respectively, significantly outperforming existing approaches. These results validate the framework’s capability to foster genuine environmental understanding and adaptive behavior in interactive planning tasks.

In summary, our main contributions are as follows:

- We introduce Interactive Planner-R1, a novel trajectory-level reinforcement learning framework that enables LLMs to learn interactive planning through autonomous environmental exploration with minimal supervision.
- We propose three key technical innovations: (1) parallel group rollouts for diverse trajectory collection, (2) completion-driven sparse rewards that prevent reward hacking while encouraging genuine environmental understanding, and (3) Interactive Policy Optimization (IPO) that extends group-relative optimization to multi-step sequential planning.
- We achieve state-of-the-art performance on ALFWorld (97.55%) and ScienceWorld

(79.92%) with superior generalization (only 3.33% performance gap on unseen environments), demonstrating the effectiveness of autonomous exploration over human-supervised approaches.

2 Preliminaries

Following established researches (Qiao et al., 2024), we formalize embodied planning tasks within the framework of a Partially Observable Markov Decision Process (POMDP) due to the agent’s inability to directly observe the environment’s complete state. Within this framework, the agent interacts within a partially observable environment and develops plans based on environmental feedback. The POMDP is defined by a 7-tuple:

$$(S, A, \Psi, \mathcal{M}, O, R, \gamma), \quad (1)$$

where S represents the state space defined by the environment, A denotes the action space available to the agent, O constitutes the observation space of the agent, $\Psi : S \times A \rightarrow S$ is the state transition function determining how actions will have effect on environmental states, $\mathcal{M} : S \rightarrow O$ is the observation function of the agent. Due to partial observability constraints, the agent cannot directly access states $s \in S$, but rather obtains partial observations of the current environment via the observation function, such that $o = \mathcal{M}(s)$, where $o \in O$ and $s \in S$. Additionally, R represents the reward function, and γ is the discount factor.

In the context of planning tasks, the agent receives a task q , typically expressed as a natural language instruction. The agent’s objective is to complete this instruction through environmental interaction. During this interaction process, the agent initially obtains an observation o_0 from the environment and initializes a trajectory $\tau_0 = (q, o_0)$. At each time step t , the policy model π_θ generates the current action based on the historical trajectory:

$$a_t = \pi_\theta(\tau_t) = \pi_\theta(q, o_0, a_0, o_1, \dots, a_{t-1}, o_t). \quad (2)$$

When performed, the action will lead to a change to the environment, resulting in a new state $s_{t+1} = \Psi(s_t, a_t)$. Following this transition, the agent receives a new partial observation $o_{t+1} = \mathcal{M}(s_{t+1})$. Subsequently, the trajectory is updated to $\tau_{t+1} = (q, o_0, a_0, o_1, \dots, a_t, o_{t+1})$, incorporating all past actions and observations with the new observation. The reward for the agent is defined with a binary signal, where the agent receives a reward of 1 if the

new state s_{t+1} meets the predefined goal conditions specified by q , otherwise the reward is 0.

3 Interactive Planner-R1

3.1 Overview

In this section, we introduce **Interactive Planner-R1**, an end-to-end training framework designed for long-term planning. Interactive Planner-R1 learns solely from outcome-based reward signals without requiring additional human supervision. Our key insight is to unleash the latent reasoning and planning capabilities of large language models (LLMs) through direct interaction with the environment, allowing them to self-evolve via reinforcement learning. To better model the interaction process, we formulate it as **ReAct-style trajectories**, where progress is made through continuous ‘*Observation* \rightarrow *Think* \rightarrow *Action*’ cycles until task completion. Interactive Planner-R1 consists of three tightly integrated core components: Firstly, to address the problem of sparse rewards, we develop a parallel sampling strategy that simultaneously generates multiple trajectories forming trajectory groups to establish stable statistical baselines, which we call **Group Rollout with In-Environment Interaction**. Secondly, we design a pure **Completion-based Sparse Reward** structure that encourages development of genuine environmental understanding through autonomous interaction. Finally, and most importantly, we introduce the **Interactive Policy Optimization (IPO)** algorithm, which is specifically designed for ReAct-style multi-turn interactions as its POMDP-oriented extension.

3.2 ReAct Paradigm

As described above, at time step t , the agent generates an action based on a trajectory that includes environmental observations and previous actions. However, mapping directly from such a trajectory to action demands complex reasoning capabilities, like task goal decomposition, common sense knowledge application, and extraction of key information from observations is exceptionally challenging. Following Yao et al. (2023), we adopt the ReAct paradigm that introduces the thought process (Wei et al., 2022) into LLM-based planning trajectories, as demonstrated in Figure 3. Unlike traditional approaches, with the historical trajectory as input, the agent generates a thought context:

$$\tau_{t+1} = (q, o_0, m_0, o_1, \dots, m_t, o_{t+1}), \quad (3)$$

where m_t denotes the model’s response for t -th step, it consists of thoughts and a text-form action:

$$m_t = (\phi_t, a_t). \quad (4)$$

We extract action a_t from m_t , which is then used to interact with the environment to obtain the new observation o_{t+1} .

3.3 Group Rollout with In-Environment Interaction

To effectively train an LLM for solving interactive environment tasks, we introduce a group rollout procedure that improves exploration efficiency while preserving a simple and well-defined interaction protocol.

Episode-level Parallel Rollout. Our group rollout mechanism should be distinguished from step-level branching or tree-search-based rollouts. Instead of expanding multiple actions from a shared intermediate environment state at each timestep, we perform episode-level parallel rollouts. Specifically, for each sampled task–environment pair, we initialize n independent and identical environment replicas at the beginning of an episode. Each replica maintains its own environment state, observation history, and trajectory buffer throughout the rollout.

A single frozen policy $\pi_{\theta_{\text{old}}}$ is used as the decision-maker for all replicas. At timestep t , the action for replica i is sampled from $\pi_{\theta_{\text{old}}}(\cdot \mid h_i^t)$, where h_i^t denotes the unique history of that replica up to step t . Due to policy stochasticity and environment dynamics, trajectories naturally diverge over time. This procedure constitutes a set of independent episode-level simulations rather than a branching rollout from shared intermediate states.

Parallel Rollout Procedure. At each training iteration, we sample K task–environment pairs from the distribution \mathcal{D} and create n parallel replicas for each pair. Each replica is initialized with a reset environment and an empty trajectory buffer. For up to max_steps interaction steps, the frozen policy $\pi_{\theta_{\text{old}}}$ generates a ReAct-style response conditioned on the task instruction and the current trajectory history. From each response, we parse an executable action.

To allow explicit termination, we introduce a special done action. The policy may generate done when it determines that the task has been completed or is unsolvable. If done is generated before the

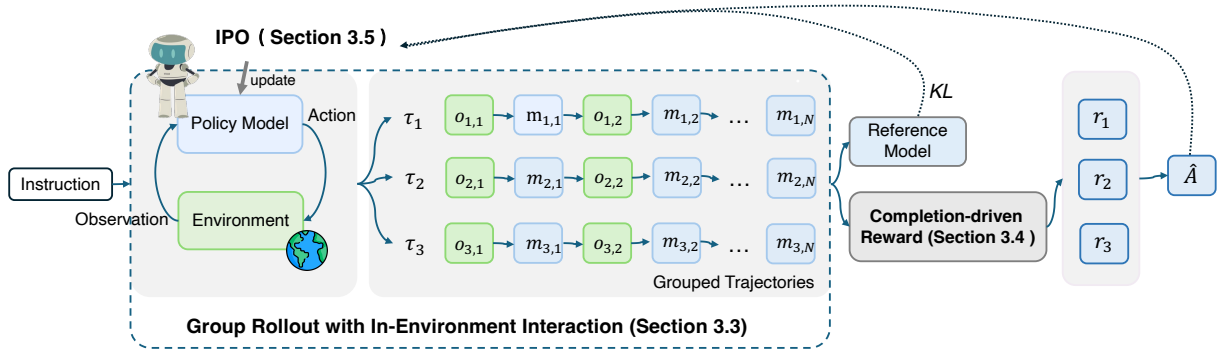


Figure 2: The Interactive Planner-R1 consists of three tightly integrated core components: 1) Group Rollout with Environment Interaction, a parallel sampling strategy that simultaneously generates multiple trajectories to form trajectory groups; 2) Completion-driven Sparse Reward; 3) Interactive Policy Optimization algorithm that redesigns probability ratios and advantage allocation, and updates the policy function based on group relative rewards. The KL divergence between policy and reference model serves as a regularization term to prevent over-optimization.

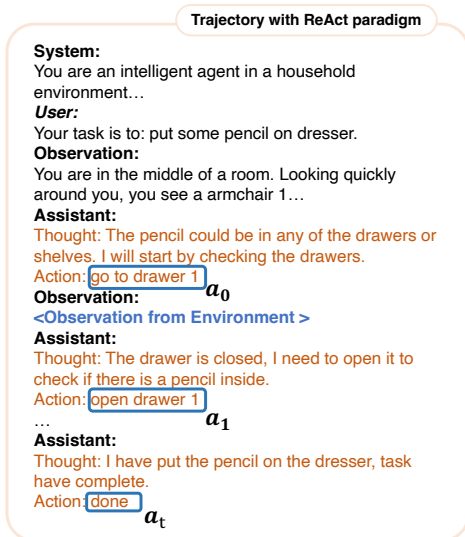


Figure 3: We employ the ReAct paradigm during generation and extract the action from the response m_t .

environment confirms task completion, the corresponding trajectory receives zero reward.

The rollout for a task terminates when all replicas either emit done or reach the maximum step limit. All collected trajectories are then aggregated into a global buffer \mathcal{T} , from which trajectory-level returns are computed and used to update the policy according to Eqs. (7)–(8). We provide the algorithm pseudo code in Appendix B and an illustrative example of the parallel rollout process is provided in Appendix C.

3.4 Completion-driven Reward

To encourage autonomous exploration, we use only a rule-based completion reward without further re-

striction to optimize the model:

$$r = \begin{cases} 0 & \text{if the task is incomplete,} \\ 1 & \text{if the task is completed.} \end{cases} \quad (5)$$

Unlike previous approaches that introduced format rewards or length penalty, we do not explicitly provide such constraints. On the one hand, we found that LLMs can maintain a clear response structure with guidance from the system prompt, and only a parsable and valid response can interact with the environment to complete the task. On the other hand, we argue that reducing human priors in reward design can lower the risk of reward hacking and encourage the model to explore autonomously, generating more flexible and diverse trajectories.

3.5 Interactive Policy Optimization (IPO)

To enable efficient environment interaction and planning while substantially reducing training cost, we introduce Interactive Policy Optimization (IPO). IPO relies on the group rollout with In-Environment Interaction to make the policy model π_θ generate multiple trajectories. We employ the group-normalized advantage estimator, which has no learnable critic model, but its probability ratio is calculated over the designated trajectories tokens prefix up to step t , including all preceding thought ϕ and action a tokens:

$$Pr_t(\theta) = \frac{\pi_\theta(\phi_t, a_t | \tau_{i, < t})}{\pi_{\theta_{\text{old}}}(\phi_t, a_t | \tau_{i, < t})}. \quad (6)$$

We optimize the policy model by maximizing

the following objective function:

$$\mathcal{J}_{\text{IPO}}(\theta) = \mathbb{E}_{(q,\epsilon) \sim \mathcal{D}, \{\tau_i\}_{i=1}^n \sim \pi_{\theta, \text{old}}(\mathcal{T}(q,\epsilon))} \left\{ \frac{1}{n} \sum_{i=1}^n \frac{1}{|\tau_i|} \sum_{t=1}^{|\tau_i|} \left[\min \left[Pr_i(\theta) \hat{A}_{i,t}, \text{clip} \left(Pr_i(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t} \right] - \beta \mathbb{D}_{\text{KL}}[\pi_{\theta} \parallel \pi_{\text{ref}}] \right] \right\}, \quad (7)$$

where $\hat{A}_{i,t}$ is the advantage of i -th trajectory in the group for time step t , which is calculated as follows:

$$\hat{A}_{i,t} = \frac{r_i - \mu_r}{\sigma_r}, \quad (8)$$

where t represents the t -th step of the trajectories, μ_r is the mean reward of the group, and σ_r is the reward deviation of the whole group. Additionally, the objective of IPO integrates a KL divergence penalty term between the training policy and the reference policy to prevent the model from deviating too much (Eq.7).

4 Experiments

4.1 Implementation

Evaluation To assess model capabilities in interactive planning tasks, we conducted systematic evaluations using two text-based embodied world simulator benchmarks: ALFWorld (Shridhar et al., 2021) and ScienceWorld (Wang et al., 2022). ALFWorld encompasses six categories of planning tasks set primarily in home environments, covering not only basic object manipulation (such as “pick and place”) but also tasks requiring complex interaction sequences. ScienceWorld presents a more challenging benchmark, requiring models to complete scientific experiments in a highly interactive environment. Each task was evaluated only once, with interaction steps limited to 30 to ensure evaluation efficiency and consistency. More evaluation and benchmark details will be introduced in Appendix D.1.

Baseline We compared our Interactive Planner-R1 against three state-of-the-art training-based methods: SFT (Zeng et al., 2024a), NAT (Wang et al., 2024b), and ETO (Song et al., 2024). We also evaluated larger foundation models including GPT-3.5 Turbo (Brown et al., 2020)(gpt-3.5-turbo-0125), GPT-4 (OpenAI, 2023)(gpt-4o-2024-08-06), DeepSeek-V3 (DeepSeek-AI et al., 2024)(DeepSeek-V3-0324), and InternLM 2.5 20B (Team et al., 2024). All training-based methods are fine-tuned on Qwen2.5-7B-Instruct (Yang et al., 2024a). Detailed baseline descriptions are provided in Appendix D.2.

Training Details We implemented our approach using the verl library and trained on 8 NVIDIA A100 80GB GPUs. Each training step sampled 128 tasks with 5 trajectories per task. Additional implementation details are provided in Appendix D.3.

4.2 Main results

As shown in Table 1, we conducted comprehensive evaluations of prompt-based and training-based approaches on two challenging interactive environments, ALFWorld and ScienceWorld. The experimental results reveal several key findings:

First, current state-of-the-art foundational LLMs demonstrate notable limitations in interactive planning capabilities. While GPT-4 achieved the best performance (55.41% on ALFWorld and 58.18% on ScienceWorld), it still lags behind training-based approaches, indicating the challenges large language models face in environmental interaction tasks.

Compared with other baselines, Interactive Planner-R1 demonstrates consistent improvements across all evaluation metrics. On the ALFWorld benchmark, it achieved accuracy rates of 98.57% and 96.52% in seen and unseen scenarios, respectively, surpassing the ETO by 14.28% and 16.67%. On the ScienceWorld dataset, Interactive Planner-R1 achieved accuracy rates of 83.66% and 76.18% in seen and unseen scenarios, outperforming competitive methods by 28.58% and 18.86%, respectively. The performance improvement from the Qwen2.5-7B-Instruct (averaging 31.05% on ALFWorld and 22.05% on ScienceWorld) to Interactive Planner-R1 demonstrates the effectiveness of our reinforcement learning framework.

We evaluate three trajectory-based baselines in ALFWorld that leverage human prior knowledge to construct training datasets and observe their performance clustering around 80%, suggesting a performance ceiling imposed by the reliance on human-annotated data. In contrast, Interactive Planner-R1 easily surpasses this ceiling, achieving over 95% by leveraging reinforcement learning.

4.3 Ablation Study

To comprehensively understand the contribution of each component in Interactive Planner-R1, we conduct systematic ablation studies across three key dimensions: policy learning architecture, ReAct paradigm integration, group rollout mechanism. All ablation experiments are performed on ALFWorld to ensure consistent evaluation conditions.

Method	ALFWorld			ScienceWorld		
	Seen	Unseen	Avg.	Seen	Unseen	Avg.
GPT-3.5-Turbo (Brown et al., 2020)	7.14	7.46	7.30	28.03	21.63	24.83
GPT-4o (OpenAI, 2023)	58.57	52.24	55.41	59.6	56.75	58.18
Deepseek-v3 (DeepSeek-AI et al., 2024)	36.43	31.34	33.89	29.28	27.45	28.37
InternLM 2.5 20B Chat (Team et al., 2024)	7.86	11.94	9.90	33.66	31.15	32.41
Qwen2.5-72B-Instruct (Yang et al., 2024a)	25.71	26.12	25.92	29.38	27.96	28.67
Qwen2.5-7B-Instruct (Yang et al., 2024a)	30.00	32.09	31.05	28.01	16.09	22.05
Qwen2.5-7B-Instruct + SFT(Zeng et al., 2024a)	80.71	79.10	79.91	68.81	55.7	62.26
Qwen2.5-7B-Instruct + NAT(Wang et al., 2024b)	63.57	66.42	65.00	58.56	49.75	54.16
Qwen2.5-7B-Instruct + ETO(Song et al., 2024)	84.29	79.85	82.07	55.08	57.32	56.20
Qwen2.5-7B-Instruct + Interactive Planner-R1(Ours)	98.57 ± 0.7	96.52 ± 1.9	97.55	83.66 ± 0.5	76.18 ± 0.9	79.92

Table 1: **Main Results.** The best results are marked in **bold**. We compare our Interactive Planner-R1 with representative baselines including vanilla language models and training-based methods on ALFWorld and ScienceWorld benchmarks, evaluating performance on both seen and unseen task scenarios. See Table 5 for complete comparison with additional baselines.

	Seen	Unseen	Avg.	Δ
Interactive Planner-R1	98.57	96.52	97.55	0
w/o React	96.56	95.71	96.14	-1.41
w/o IPO	21.43	26.11	23.77	-73.78
w/o group		N/A		-

Note: ‘w/o group’ ($k = 1$) is not applicable to IPO, as within-group normalization collapses and the advantage estimator degenerates, resulting in non-learning behavior. We therefore report it as N/A.

Table 2: Ablation study on ALFWorld. Δ denotes the performance difference compared to the full model.

As shown in Table 2, the experimental results reveal several critical insights about our framework design:

Trajectory-level vs. Single-turn Policy Learning

The most significant finding concerns the importance of Interactive Policy Optimization (IPO) for multi-turn planning tasks. To isolate the effect of the optimization objective from increased exploration, we compare IPO with a naive GRPO baseline (denoted as w/o IPO) under *identical rollout and optimization budgets*, including the same group size ($k = 5$), the same number of rollouts per task, and the same number of gradient updates.

Removing IPO and replacing it with this single-turn optimization objective results in a dramatic performance collapse, dropping from 97.55% to merely 23.77%—a substantial decrease of 73.78%. Despite having access to the same amount of interaction data, the naive GRPO baseline fails to learn effective multi-step strategies, demonstrating that performance gains do not arise from collecting more trajectories alone. Instead, effective interactive planning critically depends on trajectory-

level, group-relative optimization that captures long-horizon credit assignment across action sequences.

ReAct Paradigm Contribution

The structured *Observation* \rightarrow *Think* \rightarrow *Action* paradigm shows measurable but moderate impact on performance. Removing the explicit reasoning chains (w/o React) leads to a 1.41% performance decrease, from 97.55% to 96.14%. While this decline appears modest, it represents consistent degradation across both seen and unseen scenarios, indicating that structured reasoning helps the model better understand environmental contexts and maintain performance stability during planning.

Group Rollout Necessity

The group rollout mechanism proves fundamental to our approach’s functionality. As indicated by the missing results for “w/o group” in Table 2, this ablation becomes infeasible because group-normalized advantages require multiple parallel trajectories for statistical baseline computation. Without diverse trajectory groups, the advantage estimation becomes unstable, preventing effective policy learning entirely. This dependency underscores that parallel exploration is not merely beneficial but architecturally essential for our framework.

These ablation results collectively demonstrate that Interactive Planner-R1’s superior performance stems from the synergistic combination of trajectory-level optimization, structured reasoning, and parallel exploration, with IPO serving as the most critical component for achieving effective interactive planning capabilities.

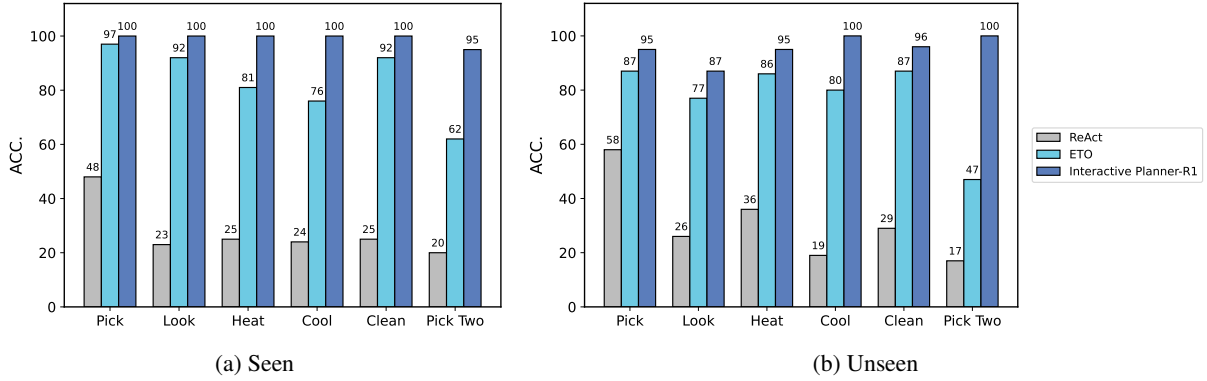


Figure 4: Completion rates of Qwen2.5-7B-Instruct on various ALFWorld tasks, covering both seen and unseen scenarios. We compare the prompt-based method ReAct with the training-based method ETO and Interactive Planner-R1.

	Pick	Look	Heat	Clean	Cool	Pick Two	Avg.
ETO	-9.89	-15.82	+6.95	-5.94	+5.39	-24.72	-7.34
Ours	-4.20	-12.50	-4.55	-3.13	0.00	+4.38	-3.33

Table 3: ETO and Interactive Planner-R1’s relative improvement on unseen tasks compared with seen tasks in ALFWorld, where Interactive Planner-R1 maintains strong generalization ability even on unseen tasks.

4.4 Training Analysis

Interactive Planner-R1 maintains high-level performance even on more difficult problems and unseen environments. Figure 4 compares the performance differences among various methods across different ALFWorld tasks, revealing clear difficulty hierarchies and capability differentials among the methods. As shown in the Figure, the “Pick” single-object task is relatively the most basic, where all training-based methods perform well, while purely prompt-based methods achieve less than 60% task completion rate. The “Pick Two” task is clearly the most challenging, with ReAct performing worst on this task (only 20.8% in seen and 17.6% in unseen), and ETO also showing significant decline (62.5% in seen, 47.05% in unseen). Only Interactive Planner-R1 maintains high-level performance (95.8% in seen, 100% in unseen). This difficulty distribution reflects the challenges drastically increase when tasks involve more object interactions, complex state transitions, or environmental reasoning. However, Interactive Planner-R1 maintains stable performance across the different tasks.

Generalization capabilities in unseen environments. To evaluate the generalization capability of different methods, we define the generalization gap $\Delta(abs)$ as the difference between task completion rates on unseen and seen tasks: $\Delta(abs) =$

$Acc_{unseen} - Acc_{seen}$, as shown in Table 3. Experimental results demonstrate that our method exhibits significant advantages in generalization performance. Overall, our approach achieves a smaller average generalization gap (-3.33%) compared to the baseline method ETO (-7.34%). Notably, in the complex “Pick Two” task, our method not only overcomes the generalization challenge but achieves positive improvement (+4.38%), while ETO shows substantial performance degradation (-24.72%). This superior generalization ability is particularly evident in handling complex tasks. Further analysis reveals that our method also demonstrates stronger stability, with a relatively smaller performance variance across different task types (ranging from -12.50% to +4.38%). These results indicate that Interactive Planner-R1 effectively enhances the model’s adaptability to variations in task scenarios while maintaining cross-task stability.

Efficient planning through invalid action reduction. Figure 5 demonstrates how Interactive Planner-R1 achieves efficient planning by systematically reducing invalid actions during training in the ScienceWorld. The reward curve shows steady improvement as the model learns to generate more valid action sequences (see Figure 5a). Notably, the average response length decreases from approximately 900 tokens to 500 tokens by step 50, reflecting more concise and efficient action planning. This efficiency gain is directly linked to the model’s ability to avoid invalid actions, as shown in Figure 5b. The total number of action steps (red line) reduces from an initial 20 steps to a consistent 12.5 steps, while invalid actions (blue line) are nearly eliminated after 50 training steps. This substantial reduction in invalid actions indicates that Interactive Planner-R1 has developed a robust

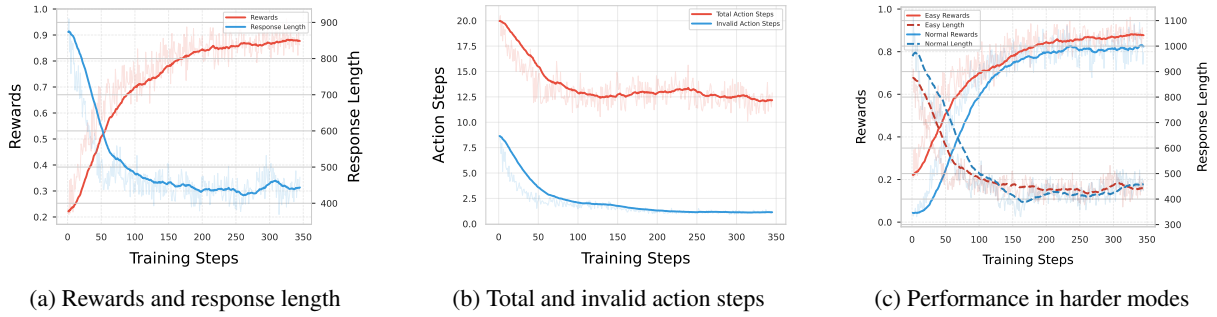


Figure 5: Training progress of Interactive Planner-R1 in ScienceWorld: (a) shows the rewards and total response length in multi-round exploration; (b) illustrates the decreasing trends in both total and invalid action steps during training; (c) reveals the sustained performance trends across varying difficulty modes.

514 understanding of environmental constraints, lead-
 515 ing to more efficient planning strategies and more
 516 efficient plans.

517 **Robust Performance Under Increased Environ-**
 518 **ment Complexity Level.** Figure 5c further vali-
 519 dates this capability by comparing learning per-
 520 formance in ScienceWorld under "Easy" and "Nor-
 521 mal" difficulty modes. While "Normal" mode re-
 522 quires complete environmental interactions (navi-
 523 gation, container manipulation, fundamental ac-
 524 tions), Interactive Planner-R1 demonstrates notable
 525 progress in both settings. Though "Easy" mode
 526 shows faster convergence (0.6 reward at 50 steps
 527 vs 100 steps) and higher final performance (0.88
 528 vs 0.82), the model maintaining flexible explo-
 529 ration strategies and planning capabilities even un-
 530 der challenging environmental constraints.

5 Related Work

531 **LLM-based Long-Horizon Planning.** LLMs
 532 have been widely studied for long-horizon plan-
 533 ning in interactive environments (Liu et al., 2024a;
 534 Xi et al., 2025). Prompt-based methods, such as
 535 ReAct (Yao et al., 2023) and its extensions (Shinn
 536 et al., 2023; Yao et al., 2024), enhance planning
 537 via structured prompting without additional train-
 538 ing, but often struggle in complex or partially ob-
 539 servable environments. Training-based approaches
 540 typically rely on supervised fine-tuning or behav-
 541 ioral cloning on static expert trajectories (Yin et al.,
 542 2024; Song et al., 2024), which limits exploration
 543 diversity and generalization.

544 In contrast, our work learns planning policies
 545 through online environment interaction, enabling
 546 LLM agents to actively explore and adapt their
 547 strategies under partial observability.

548 **Interactive Agents and Reinforcement Learn-**
 549 **ing.** Recent work shows that reinforcement learn-
 550 ing with simple rewards can drive strong self-
 551

552 improvement in LLMs (Lightman et al., 2024;
 553 DeepSeek-AI et al., 2025), but most studies focus
 554 on static or single-turn reasoning tasks. Concur-
 555 rently, interactive agent frameworks such as API-
 556 Gen (Liu et al., 2024b) and LAM Simulator (Hoang
 557 et al., 2025) explore tool use or language-based hu-
 558 man feedback as supervision.

559 Different from these approaches, our method
 560 targets text-based embodied environments and
 561 learns long-horizon interactive planning purely
 562 from sparse, outcome-level environment rewards
 563 via trajectory-level reinforcement learning, without
 564 dense supervision or human feedback.

565 More detailed discussions of related work are
 566 provided in Appendix A.

6 Conclusion

567 In this paper, we propose Interactive Planner-R1,
 568 a reinforcement learning framework that enables
 569 LLMs to achieve self-evolution in multi-turn en-
 570 vironmental interaction planning tasks. Our ap-
 571 proach demonstrates that models can autonomously
 572 develop effective planning strategies through care-
 573 fully designed prompt constraints and outcome-
 574 based reward mechanisms, without relying on ex-
 575 pert trajectories or extensive human priors. Ex-
 576 perimental results across multiple environments
 577 show that Interactive Planner-R1 significantly out-
 578 performs existing baselines, validating the effec-
 579 tiveness of self-evolution via reinforcement learn-
 580 ing in complex interactive scenarios. Notably, our
 581 analysis reveals that multi-turn interactive plan-
 582 ning benefits from shorter context lengths, as per-
 583 formance improves when ineffective actions are re-
 584 duced. This suggests that interactive planning tasks
 585 prioritize efficient action planning over lengthy
 586 deliberation. These findings may provide some
 587 insights for LLM-based interactive planning sys-
 588 tems.
 589

590 Limitations

591 While Interactive Planner-R1 demonstrates strong
592 performance on interactive planning tasks, several
593 limitations merit consideration for future research:

594 **Environment Scope.** Our evaluation focuses ex-
595 clusively on text-based environments (ALFWorld
596 and ScienceWorld). Although the underlying prin-
597 ciples of multi-trajectory reinforcement learning
598 are domain-agnostic, empirical validation in mul-
599 timodal settings—particularly vision-language en-
600 vironments—remains an important direction for
601 establishing broader applicability.

602 **Computational Requirements.** The group roll-
603 out mechanism, while effective for exploration di-
604 versity, introduces substantial computational over-
605 head compared to single-trajectory methods. Train-
606 ing requires 8 A100 GPUs for several hours,
607 which may limit accessibility for researchers with
608 constrained resources. Developing more sample-
609 efficient variants or exploring distributed training
610 strategies could address this scalability concern.

611 References

612 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
613 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
614 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
615 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
616 Gretchen Krueger, Tom Henighan, Rewon Child,
617 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
618 Clemens Winter, and 12 others. 2020. [Language
619 models are few-shot learners](#). In *Advances in Neural
620 Information Processing Systems 33: Annual Confer-
621 ence on Neural Information Processing Systems 2020,
622 NeurIPS 2020, December 6-12, 2020, virtual*.

623 Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Col-
624 lier, Karthik Narasimhan, and Shunyu Yao. 2023a. [Fireact: Toward language agent fine-tuning](#). *CoRR*,
625 abs/2310.05915.
626

627 Yaran Chen, Wenbo Cui, Yuanwen Chen, Mining Tan,
628 Xinyao Zhang, Dongbin Zhao, and He Wang. 2023b. [Robogpt: an intelligent agent of making embodied
629 long-term decisions for daily instruction tasks](#). *CoRR*,
630 abs/2311.15649.
631

632 Sanjiban Choudhury. 2025. [Process reward models
633 for LLM agents: Practical framework and directions](#).
634 *CoRR*, abs/2502.10325.

635 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,
636 Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
637 Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,
638 Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi-
639 hong Shao, Zhuoshu Li, Ziyi Gao, and 81 others.
640 2025. [Deepseek-r1: Incentivizing reasoning capa-
641 bility in llms via reinforcement learning](#). *CoRR*,
642 abs/2501.12948.

DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingx-
uan Wang, Bochao Wu, Chengda Lu, Chenggang
Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,
Damai Dai, Daya Guo, Dejian Yang, Deli Chen,
Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai,
and 81 others. 2024. [Deepseek-v3 technical report](#).
CoRR, abs/2412.19437.

Jiafei Duan, Samson Yu, Hui Li Tan, Hongyuan Zhu,
and Cheston Tan. 2022. [A survey of embodied
AI: from simulators to research tasks](#). *IEEE Trans.
Emerg. Top. Comput. Intell.*, 6(2):230–244.

James J. Gibson. 1979. *The Ecological Approach to Vi-
sual Perception: Classic Edition*. Houghton Mifflin.

Thai Quoc Hoang, Kung-Hsiang Huang, Shirley
Kokane, Jianguo Zhang, Zuxin Liu, Ming Zhu, Jake
Grigsby, Tian Lan, Michael S. Ryoo, Chien-Sheng
Wu, Shelby Heinecke, Huan Wang, Silvio Savarese,
Caiming Xiong, and Juan Carlos Niebles. 2025. [LAM SIMULATOR: advancing data generation for
large action model training via online exploration
and trajectory feedback](#). In *Findings of the Associ-
ation for Computational Linguistics, ACL 2025, Vi-
enna, Austria, July 27 - August 1, 2025*, pages 12921–
12934. Association for Computational Linguistics.

Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xi-
angyu Zhang, and Heung-Yeung Shum. 2025. [Open-
reasoner-zero: An open source approach to scaling
up reinforcement learning on the base model](#). *CoRR*,
abs/2503.24290.

Chengxing Jia, Ziniu Li, Pengyuan Wang, Yi-Chen
Li, Zhenyu Hou, Yuxiao Dong, and Yang Yu. 2025. [Controlling large language model with latent actions](#).
ArXiv, abs/2503.21383.

Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harri-
son Edwards, Bowen Baker, Teddy Lee, Jan Leike,
John Schulman, Ilya Sutskever, and Karl Cobbe.
2024. [Let’s verify step by step](#). In *The Twelfth In-
ternational Conference on Learning Representations,
ICLR 2024, Vienna, Austria, May 7-11, 2024*. Open-
Review.net.

Zongyu Lin, Yao Tang, Xingcheng Yao, Da Yin, Ziniu
Hu, Yizhou Sun, and Kai-Wei Chang. 2025. [Qlass:
Boosting language agent inference via q-guided step-
wise search](#). *ArXiv*, abs/2502.02584.

Yang Liu, Weixing Chen, Yongjie Bai, Guanbin Li, Wen
Gao, and Liang Lin. 2024a. [Aligning cyber space
with physical world: A comprehensive survey on
embodied AI](#). *CoRR*, abs/2407.06886.

Yuqi Liu, Bohao Peng, Zhisheng Zhong, Zihao Yue,
Fanbin Lu, Bei Yu, and Jiaya Jia. 2025. [Seg-zero:
Reasoning-chain guided segmentation via cognitive
reinforcement](#). *CoRR*, abs/2503.06520.

Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian
Lan, Shirley Kokane, Juntao Tan, Weiran Yao, Zhi-
wei Liu, Yihao Feng, Rithesh R. N., Liangwei Yang,
Silvio Savarese, Juan Carlos Niebles, Huan Wang,

699	Shelby Heinecke, and Caiming Xiong. 2024b. Api-gen: Automated pipeline for generating verifiable and diverse function-calling datasets . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. 2024. Trial and error: Exploration-based trajectory optimization for LLM agents . <i>CoRR</i> , abs/2403.02502.	755	
700			756	
701			757	
702			758	
703				
704		Haotian Sun, Yuchen Zhuang, Lingkai Kong, Bo Dai, and Chao Zhang. 2023a. Adaplaner: Adaptive planning from feedback with language models . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	759	
705			760	
706	Chengqi Lyu, Songyang Gao, Yuzhe Gu, Wenwei Zhang, Jianfei Gao, Kuikun Liu, Ziyi Wang, Shuaibin Li, Qian Zhao, Haian Huang, Weihao Cao, Jiangning Liu, Hongwei Liu, Junnan Liu, Songyang Zhang, Dahua Lin, and Kai Chen. 2025. Exploring the limit of outcome reward for learning mathematical reasoning . <i>CoRR</i> , abs/2502.06781.		761	
707			762	
708			763	
709			764	
710			765	
711		Jiankai Sun, Chuanyang Zheng, Enze Xie, Zhengying Liu, Ruihang Chu, Jianing Qiu, Jiaqi Xu, Mingyu Ding, Hongyang Li, Mengzhe Geng, Yue Wu, Wenhai Wang, Junsong Chen, Zhangyue Yin, Xiaozhe Ren, Jie Fu, Junxian He, Wu Yuan, Qi Liu, and 15 others. 2023b. A survey of reasoning with foundation models . <i>CoRR</i> , abs/2312.11562.	766	
712			767	
713	So Yeon Min, Devendra Singh Chaplot, Pradeep Kumar Ravikumar, Yonatan Bisk, and Ruslan Salakhutdinov. 2022. FILM: following instructions in language with modular methods . In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.		768	
714			769	
715			770	
716			771	
717			772	
718		InternLM Team, Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan, Zhaoye Fei, Yang Gao, Jiaye Ge, Chenya Gu, Yuzhe Gu, and 53 others. 2024. Internlm2 technical report . <i>CoRR</i> , abs/2403.17297.	773	
719	Kolby Nottingham, Bodhisattwa Prasad Majumder, Bhavana Dalvi Mishra, Sameer Singh, Peter Clark, and Roy Fox. 2024. Skill set optimization: Reinforcing language model behavior via transferable skills . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.		774	
720			775	
721			776	
722			777	
723			778	
724		Hanlin Wang, Jian Wang, Chak Tou Leong, and Wenjie Li. 2025. Steca: Step-level trajectory calibration for llm agent learning . <i>ArXiv</i> , abs/2502.14276.	779	
725			780	
726	OpenAI. 2023. GPT-4 technical report . <i>CoRR</i> , abs/2303.08774.		781	
727				
728	Shuofei Qiao, Runnan Fang, Ningyu Zhang, Yuqi Zhu, Xiang Chen, Shumin Deng, Yong Jiang, Pengjun Xie, Fei Huang, and Huajun Chen. 2024. Agent planning with world knowledge model . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .		782	
729			783	
730			784	
731			785	
732				
733				
734				
735				
736	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models . <i>CoRR</i> , abs/2402.03300.		786	
737			787	
738			788	
739			789	
740				
741	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: language agents with verbal reinforcement learning . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .		790	
742			791	
743			792	
744			793	
745			794	
746			795	
747			796	
748	Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew J. Hausknecht. 2021. Alfworld: Aligning text and embodied environments for interactive learning . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.		797	
749			798	
750			799	
751			800	
752			801	
753			802	
754				
		Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. 2023. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents . <i>CoRR</i> , abs/2302.01560.	803	
			804	
			805	
			806	
			807	
			808	
			809	
			810	
		Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models . In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .		

811	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, and 9 others. 2025. The rise and potential of large language model based agents: a survey . <i>Sci. China Inf. Sci.</i> , 68(2).	
819	Tian Xie, Zitian Gao, Qingnan Ren, Haoming Luo, Yuqian Hong, Bryan Dai, Joey Zhou, Kai Qiu, Zhirong Wu, and Chong Luo. 2025. Logic-rl: Unleashing LLM reasoning with rule-based reinforcement learning . <i>CoRR</i> , abs/2502.14768.	
824	Weimin Xiong, Yifan Song, Xiutian Zhao, Wenhao Wu, Xun Wang, Ke Wang, Cheng Li, Wei Peng, and Sujian Li. 2024. Watch every step! LLM agent learning via iterative step-level process refinement . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 1556–1572, Miami, Florida, USA. Association for Computational Linguistics.	
832	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jixi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and 22 others. 2024a. Qwen2.5 technical report . <i>CoRR</i> , abs/2412.15115.	
839	Chen Yang, Chenyang Zhao, Quanquan Gu, and Dongruo Zhou. 2024b. Cops: Empowering LLM agents with provable cross-task experience sharing . <i>CoRR</i> , abs/2410.16670.	
843	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
849	Weiran Yao, Shelby Heinecke, Juan Carlos Niebles, Zhiwei Liu, Yihao Feng, Le Xue, Rithesh R. N., Zeyuan Chen, Jianguo Zhang, Devansh Arpit, Ran Xu, Phil Mui, Huan Wang, Caiming Xiong, and Silvio Savarese. 2024. Retroformer: Retrospective large language agents with policy gradient optimization . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	
858	Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Raghavi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. 2024. Agent lumos: Unified and modular training for open-source language agents . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 12380–12403. Association for Computational Linguistics.	
	Siyu Yuan, Zehui Chen, Zhiheng Xi, Junjie Ye, Zhengyin Du, and Jiecao Chen. 2025. Agent-r: Training language model agents to reflect via iterative self-training . <i>CoRR</i> , abs/2501.11425.	867 868 869 870
	Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. 2024a. Agenttuning: Enabling generalized agent abilities for llms . In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 3053–3077. Association for Computational Linguistics.	871 872 873 874 875 876 877
	Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. 2024b. AgentTuning: Enabling generalized agent abilities for LLMs . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 3053–3077, Bangkok, Thailand. Association for Computational Linguistics.	878 879 880 881 882 883
	Zirui Zhao, Wee Sun Lee, and David Hsu. 2023. Large language models as commonsense knowledge for large-scale task planning . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	884 885 886 887 888 889 890
	Yuqi Zhu, Shuofei Qiao, Yixin Ou, Shumin Deng, Shiwei Lyu, Yue Shen, Lei Liang, Jinjie Gu, Hua-jun Chen, and Ningyu Zhang. 2025. KnowAgent: Knowledge-augmented planning for LLM-based agents . In <i>Findings of the Association for Computational Linguistics: NAACL 2025</i> , pages 3709–3732, Albuquerque, New Mexico. Association for Computational Linguistics.	891 892 893 894 895 896 897 898

899 A Related Work

900 A.1 Long-Horizon Planning with Large 901 Language Models

902 Large language models (LLMs) have recently
903 demonstrated strong potential for long-horizon
904 planning in interactive environments (Liu et al.,
905 2024a; Xi et al., 2025; Duan et al., 2022). Existing
906 approaches can be broadly categorized into prompt-
907 based methods and training-based methods.

908 Prompt-based methods aim to enhance planning
909 capabilities through carefully designed prompts
910 or exemplars without additional training. Re-
911 Act (Yao et al., 2023) introduces an explicit think-
912 act paradigm that enables LLMs to interleave rea-
913 soning with environment interaction. Subsequent
914 work extends this framework by incorporating skill
915 abstraction (Nottingham et al., 2024; Qiao et al.,
916 2024; Zhao et al., 2023), reflection or revision
917 mechanisms (Shinn et al., 2023; Yao et al., 2024),
918 and improved prompt engineering strategies (Sun
919 et al., 2023a). While these approaches are flexible
920 and easy to deploy, they largely rely on pre-trained
921 knowledge and often struggle in complex or par-
922 tially observable environments.

923 Training-based methods directly optimize LLMs
924 for environmental interaction, typically via super-
925 vised fine-tuning or behavioral cloning on expert
926 trajectories (Yin et al., 2024; Zeng et al., 2024a;
927 Chen et al., 2023a; Qiao et al., 2024). To enrich
928 training signals, some works incorporate negative
929 or suboptimal trajectories (Song et al., 2024; Wang
930 et al., 2024b). However, these methods usually de-
931 pend on static, pre-collected datasets, which limits
932 exploration diversity and leads to distribution shifts
933 between training and test environments.

934 Our proposed Interactive Planner-R1 addresses
935 these limitations by learning planning policies
936 through online interaction with the environment,
937 allowing LLM agents to actively explore and adapt
938 their strategies under partial observability using
939 sparse, outcome-level feedback.

940 A.2 Self-Evolution and Interactive Agents

941 Reinforcement learning has been shown to effec-
942 tively drive self-improvement in LLMs across a va-
943 riety of reasoning and problem-solving tasks (Sun
944 et al., 2023b; Lightman et al., 2024; DeepSeek-AI
945 et al., 2025). Building on this idea, prior work has
946 achieved strong performance in domains such as
947 mathematics, logical reasoning, and visual under-
948 standing (Hu et al., 2025; Lyu et al., 2025; Xie

949 et al., 2025; Liu et al., 2025). However, these set-
950 tings are predominantly static or single-turn, and
951 do not involve long-horizon interaction with an
952 external environment.

953 Beyond reinforcement learning for static rea-
954 soning, several recent studies investigate interac-
955 tive agents under alternative forms of supervision.
956 APIGen (Liu et al., 2024b) focuses on synthe-
957 sizing tool-using trajectories to train agents for
958 API invocation, emphasizing tool generation and
959 execution rather than long-horizon environment
960 planning. LAM Simulator (Hoang et al., 2025)
961 leverages language-based human feedback to guide
962 agent behavior, providing rich semantic supervi-
963 sion through simulated human interactions. These
964 approaches rely on dense or structured supervision
965 signals and primarily target tool use or instruction-
966 following scenarios.

967 In contrast, our work focuses on text-based em-
968 bodied environments and studies whether LLM
969 agents can acquire long-horizon interactive plan-
970 ning capabilities purely from sparse, outcome-level
971 environment rewards via trajectory-level reinforce-
972 ment learning, without dense supervision, human
973 feedback, or learned evaluators.

974 B Pseudo Code of Group rollout of 975 interaction with environment

976 Algorithm 1 presents the detailed procedure for our
977 group rollout strategy, which is a key component
978 of our Interactive Planner-R1 approach. The group
979 rollout mechanism enables efficient data collection
980 by running multiple parallel trajectories for each
981 task, thereby increasing sample diversity and im-
982 proving training stability.

983 The algorithm operates in three main phases.
984 First, it freezes the current policy to ensure consis-
985 tent behavior during the rollout phase and samples
986 a batch of task-environment pairs from the training
987 dataset. Second, for each task in the batch, it cre-
988 ates multiple independent environment instances
989 and initializes separate trajectory buffers. This par-
990 allel setup allows the agent to explore different
991 action sequences for the same task, capturing vari-
992 ous solution strategies and potential failure modes.
993 Third, the algorithm executes the interaction loop
994 where the frozen policy generates responses for
995 each trajectory simultaneously, extracts actions
996 from these responses, and updates the environment
997 states accordingly.

998 The key advantage of this group rollout approach

Algorithm 1 Group rollout of interaction with environment

Input: Initial policy π_{θ_0} , dataset $\mathcal{D} = \{(q, e) : q \in \mathcal{Q}, e \in \mathcal{E}\}$, where \mathcal{Q} is the set of instructions and \mathcal{E} is the set of environments with the group size set to n .

```
1: Freeze old policy  $\pi_{\theta_{\text{old}}} \leftarrow \pi_{\theta}$ 
2: Sample batch  $\mathcal{B} = \{(q_i, e_i)\}_{i=0}^K \sim \mathcal{D}$ 
3: Initialize global trajectories buffer  $\mathcal{T} \leftarrow \emptyset$ 
4: for  $(q_i, e_i) \in \mathcal{B}$  do
5:   Create  $n$  environment instances  $\{e_j\}^n$  copied from  $e_i$ 
6:   Initialize  $n$  trajectory buffers  $\{\tau_j\}^n \leftarrow \emptyset$  for environment instances  $\{e_j\}^n$ 
7:   for  $j = 1$  to  $n$  do
8:      $o_{j,0} \leftarrow e_j.\text{reset}()$ 
9:      $\tau_j \leftarrow (q_i, o_{j,0})$ 
10:  end for
11:  for  $k = 1$  to max steps do
12:    Generate responses  $\{m_{j,k} \sim \pi_{\theta_{\text{old}}}(\cdot | \tau_j)\}^n$  for each trajectory  $\tau_j$ 
13:    Extract actions  $\{a_{j,k}\}^n$  from responses  $\{m_{j,k}\}^n$ 
14:    Get observations  $\{o_{j,k}, \text{done}_{j,k} \leftarrow e_j.\text{step}(a_{j,k})\}^n$ 
15:    Update trajectories  $\{\tau_j \leftarrow \tau_j \oplus (m_{j,k}, o_{j,k})\}^n$ 
16:    if all done then
17:      break
18:    end if
19:  end for
20:  Save trajectories  $\mathcal{T} \leftarrow \mathcal{T} \cup \{\tau_j\}^n$ 
21: end for
```

Output: The grouped trajectories \mathcal{T}

999 lies in its ability to collect diverse training data effi-
1000 ciently. By running n parallel trajectories for each
1001 task, we obtain multiple examples of agent behav-
1002 ior under identical initial conditions, which helps
1003 the learning algorithm better understand the conse-
1004 quences of different action choices. The parallel ex-
1005 ecution also reduces the overall data collection time
1006 compared to sequential rollouts. The algorithm ter-
1007 minates when all environments reach their done
1008 states or the maximum number of steps is reached,
1009 ensuring that we collect complete episodes for train-
1010 ing. The collected trajectories are then used in
1011 subsequent training phases to update the policy
1012 parameters through our reward modeling and opti-
1013 mization procedures.

1014 C Illustration of Parallel Group Rollout

1015 In this section, we provides a concrete illustration
1016 of the group rollout mechanism used in Interactive
1017 Planner-R1, with the goal of clarifying its seman-
1018 tics and distinguishing it from step-level branching
1019 or tree-search-based rollout methods.

C.1 Conceptual Clarification

1020 Our group rollout mechanism performs parallel and
1021 independent simulations at the *episode level*, rather
1022 than branching from a shared intermediate environ-
1023 ment state at each timestep. For a given task, we
1024 initialize n identical environment instances at the
1025 beginning of an episode. Each instance maintains
1026 its own environment state, observation history, and
1027 reward sequence throughout the rollout. 1028

1029 A single frozen policy $\pi_{\theta_{\text{old}}}$ is used to inter-
1030 act with all environment replicas. For any trajec-
1031 tory i at timestep t , the action a_i^t is sampled from
1032 $\pi_{\theta_{\text{old}}}(\cdot | h_i^t)$, where h_i^t denotes the unique history
1033 of that trajectory. Importantly, trajectories do not
1034 share intermediate states or observations, and no
1035 information is exchanged between replicas during
1036 rollout. Any divergence among trajectories arises
1037 solely from policy stochasticity and environment
1038 dynamics.

1039 This procedure should be understood as a set
1040 of parallel episode-level simulations, rather than a
1041 step-level branching search from a common state.

1042 C.2 Example: Parallel Rollout for a Simple 1043 Task

1044 We illustrate the group rollout process using a simple
1045 example with two parallel replicas ($n = 2$).

1046 **Task.** Put the apple on the table.

1047 **Step 1.** The frozen policy $\pi_{\theta_{\text{old}}}$ generates actions
1048 for two independent environment instances initial-
1049 ized from the same task configuration.

- 1050 • Trajectory 1 history: (obs: see kitchen) \rightarrow
1051 Action: go to fridge
- 1052 • Trajectory 2 history: (obs: see kitchen) \rightarrow
1053 Action: go to counter

1054 **Step 2.** Each environment updates independently
1055 based on the executed action. The policy then gen-
1056 erates the next action conditioned on each trajec-
1057 tory’s unique observation and history.

- 1058 • Trajectory 1 history: (... , obs: at fridge) \rightarrow
1059 Action: open fridge
- 1060 • Trajectory 2 history: (... , obs: at counter, see
1061 apple) \rightarrow Action: take apple

1062 As the rollout proceeds, trajectories continue
1063 to evolve independently, potentially reaching task
1064 completion at different times or following different
1065 action sequences. This example highlights that tra-
1066 jectories within a group do not branch from shared
1067 intermediate states and that the frozen policy serves
1068 as a consistent action generator across all replicas
1069 during a rollout phase.

1070 D Implementation Details

1071 D.1 Evaluation and Benchmark Details

1072 **ALFWorld (Shridhar et al., 2021)** ALFWorld en-
1073 compasses six categories of planning tasks set pri-
1074 marily in home environments, covering not only ba-
1075 sic object manipulation (such as pick and place) but
1076 also tasks requiring complex interaction sequences.
1077 For example, the heating task requires models to
1078 first identify target objects, move them to heating
1079 devices (like microwave), execute the heating oper-
1080 ation, and finally place them in designated locations
1081 to complete the task.

1082 Following prior research (Yao et al., 2023; Song
1083 et al., 2024), we evaluate model performance under
1084 two conditions: seen and unseen scenarios. Seen
1085 scenarios consist of task instances from rooms en-
1086 countered during training, unseen scenarios com-
1087 prise task instances from entirely new rooms with

1088 container arrangements and scene organizations
1089 distinctly different from training tasks, designed
1090 to evaluate the model’s zero-shot generalization
1091 capabilities.

1092 **ScienceWorld (Wang et al., 2022)** Science-
1093 World presents a more challenging benchmark, re-
1094 quiring models to complete scientific experiments
1095 in a highly interactive environment. This environ-
1096 ment contains approximately ten interconnected
1097 rooms, over 200 object types, and 25 executable
1098 actions. Compared to ALFWorld, ScienceWorld
1099 integrates a more sophisticated physics simulation
1100 engine, including thermodynamic and electrical
1101 systems, which places more rigorous demands on a
1102 model’s planning capabilities and causal reasoning.

1103 Following prior research (Yao et al., 2023; Song
1104 et al., 2024), we assessed model performance under
1105 two conditions: seen and unseen scenarios. Seen
1106 scenarios consist of task instances from rooms en-
1107 countered during training, including the same task
1108 types, objects, containers, and room layouts, but
1109 with variations in object positions, quantities, and
1110 visual features (e.g., two blue pencils on a shelf
1111 instead of three red pencils in a drawer as seen
1112 in training data). Unseen scenarios comprise task
1113 instances from entirely new rooms with container
1114 arrangements and scene organizations distinctly
1115 different from training tasks, designed to evalu-
1116 ate the model’s zero-shot generalization capabili-
1117 ties. During evaluating ScienceWorld, we main-
1118 tained the “easy” setup consistent with previous
1119 work (Song et al., 2024), enabling all five assistive
1120 features: instant teleportation, pre-opened contain-
1121 ers, self-watering flower pots, and two additional
1122 aids that remove similar micro-action bottlenecks,
1123 thus reducing the burden of low-level operations
1124 and focusing the evaluation on high-level planning
1125 abilities. We also adopted the same balancing strat-
1126 egy as previous work (Song et al., 2024), imple-
1127 menting fair sampling at the task level.

1128 D.2 Baseline

1129 To benchmark our Interactive Planner-R1 against
1130 other training-based approaches, we compared
1131 it with three state-of-the-art methods that in-
1132 corporate human priors: SFT (Zeng et al.,
1133 2024a), which fine-tunes on expert trajectories to
1134 learn interaction capabilities; NAT (Wang et al.,
1135 2024b), which incorporates learning from re-
1136 jected trajectories; and ETO (Song et al., 2024),
1137 which learns from expert trajectories through
1138 DPO. Our evaluation also includes larger mod-

Category	Setting / Description
Backbone model	Qwen2.5-7B-Instruct
Context length	4096
Random seed	42
Rollout batch size	128
Group size	5
Learning rate	1e-6
Batch size	32
Max trajectory length	30
Clip parameter ϵ	0.2
KL penalty coefficient β	0.001
Optimizer	AdamW
Learning rate schedule	Constant

Table 4: Training, optimization, and implementation details for Interactive Planner-R1.

els: GPT-3.5 Turbo (Brown et al., 2020)(gpt-3.5-turbo-0125), GPT-4 (OpenAI, 2023)(gpt-4o-2024-08-06), DeepSeek-V3 (DeepSeek-AI et al., 2024)(DeepSeek-V3-0324), and InternLM 2.5 20B (Team et al., 2024), to demonstrate the current performance of foundation large language models (both open-source and closed-source) on interactive planning tasks. All our implemented training-based methods are exclusively fine-tuned on Qwen2.5-7B-Instruct (Yang et al., 2024a), and all models are evaluated using a consistent set of prompts. Additionally, while many methods cannot be directly compared due to differences in models and evaluation approaches, we still provide these results in Table 5 for reference.

D.3 Training Details

We implemented Interactive Planner-R1 based on the ver1 framework¹, a flexible reinforcement learning library for large language models, and extended it with an interaction module to support multi-turn environment rollouts. For clarity and reproducibility, we summarize the key training hyperparameters, optimization settings, and rollout configurations in Table 4.

Rollout and Data Collection. At each training iteration, we sample a batch of 128 task–environment pairs. For each task, we perform grouped rollouts with a group size of 5, where each group consists of independent environment replicas initialized from the same task configuration. Each trajectory follows the ReAct paradigm and is capped at a maximum of 30 interaction steps. A special done action is introduced to explicitly terminate a trajectory when the model determines

¹<https://github.com/volcengine/ver1>

that the task has been completed or is unsolvable. If done is generated prematurely without satisfying the task completion condition, the trajectory receives zero reward.

Policy Optimization. All trajectories within a group share a binary, completion-driven reward (1 if the task is completed, 0 otherwise). Advantages are computed using group-normalized returns as described in Eq. (8), which requires multiple parallel trajectories and is not well-defined for a single-trajectory setting. We optimize the policy using Interactive Policy Optimization (IPO), with clipped probability ratios ($\epsilon = 0.2$) and a KL divergence penalty with respect to a fixed reference model. Unless otherwise stated, the KL coefficient is set to $\beta = 0.001$, which we found to provide a good trade-off between training stability and exploration performance.

Model and Optimization. We use Qwen2.5-7B-Instruct (Yang et al., 2024a)² as the backbone policy model. Training is conducted with a maximum context length of 4096 tokens, which includes the instruction, observation history, and previous ReAct steps. The reference model is initialized from the same checkpoint and kept frozen throughout training. We use the AdamW optimizer with default parameters provided by ver1. The random seed is fixed to 42 for all experiments.

Training Setup and Evaluation. All experiments are conducted on 8 NVIDIA A100 GPUs with 80GB memory. Grouped rollouts are only used during training; during evaluation, we perform single-trajectory inference with a maximum of 30 steps and report pass@1 task completion rates. The sampling temperature is set to 1.0 during evaluation unless otherwise specified. Each experiment is repeated 5 times with different random seeds, and we report the mean and standard deviation.

E Prompt Design

In our research, prompt engineering accomplishes three interconnected objectives: (i) enumerating permissible action, guiding the model to execute effective commands; (ii) mitigating common hallucinations, preventing the model from falling into blind trial-and-error loops; and (iii) enforcing a concise ReAct paradigm that clearly demonstrates

²<https://huggingface.co/Qwen/Qwen2.5-7B-Instruct>

ALFWorld Templates:

Environment description

You are an intelligent agent in a household environment and your target is to perform actions to complete the task goal. At the beginning of your interactions, you will be given the detailed description of the current environment and your goal to accomplish.

For each of your turn, you will be given the observation of the last turn. You should first think about the current condition and plan for your future actions, and then output your action in this turn. Your output must strictly follow this format: **Thought:** <your thoughts> ; **Action:** <your next action>.

Available actions & required format

The available actions are:

1. go to (receptacle)

2. open (receptacle)

:

13. inventory:check your current inventory

Explanations for abstract actions

14. done:Indicate that you believe the task is complete

Where (object) refers to manipulable objects and (receptacle) refers to receptacles or locations.

Avoid hallucination

After your each turn, the environment will give you immediate feedback based on which you plan your next few steps, if the environment output: Nothing happens, that means the previous action is invalid and you should try more options.

You can only hold one object at a time. Before taking a new object, make sure you have placed down any object you are currently holding.

You should not assume or anticipate the feedback. Even if you have planned multiple steps ahead, you should only execute one action at a time.

Do not proceed with any further exploration or actions until you receive the feedback from the environment after your action.

Output template

Your response should use the following format:

Thought: [your thoughts]

Action: [your next action]

Figure 6: Prompt for ALFWorld

reasoning chains and actions, facilitating straight-forward extraction of selected operations. Following (Qiao et al., 2024), we enhanced the system instructions to explicitly highlight the allowable action set, and required output format, while specifically addressing common hallucination patterns. This approach significantly increasing the probability of obtaining initial successful trajectories necessary for guided learning. The prompt is detailed in Figure 6 and Figure 7.

F More experimental results

Figure 9a illustrates the metrics of the Interactive Planner-R1 during the training process in ALFWorld, which yields results similar to those in ScienceWorld. The reward curve exhibits a similar upward trend. Meanwhile, the response length significantly decreases from the outset, dropping from approximately 750 tokens initially to around 450 tokens by the end. The relationship between response length and action steps is also evident in Figure 5b,

ScienceWorld Templates:

System Prompt:

Environment description

You are a helpful assistant to do some scientific experiment in an environment.

You should explore the environment and find the items you need to complete the experiment.

In the environment, there are several rooms: kitchen, foundry, workshop, bathroom, outside, living room, bedroom, greenhouse, art studio, hallway. You can teleport to any room in one step.

Available actions & required format

The available actions are:

activate OBJ

close OBJ

:

wait1: wait 1 step

done: indicate that you believe the task is complete

Avoid hallucination

When arrive a new location, you should use look around to check the OBJ you can interact with. Use focus on OBJ only necessary as incorrect use will cause environment ends. Do not proceed with any further exploration or actions until you receive the feedback from the environment after your action.

Output template

Your response should use the following format:

Thought: [your thoughts]

Action: [your next action]

Figure 7: Prompt for ScienceWorld

1240 where total actions, invalid actions, and response
1241 length show a consistent downward trend. Due to
1242 the reduced difficulty compared to ScienceWorld,
1243 invalid steps have been almost entirely eliminated
1244 in the final results.

1245 In Figure 8, we present the task completion per-
1246 formance of our trained agents in ScienceWorld,
1247 alongside a comparison with ReAct and ETO. This
1248 comparison underscores the effectiveness and gen-
1249 eralization capabilities of our Interactive Planner-
1250 R1. While prompt-based and previous training-
1251 based methods often exhibit strong performance
1252 on tasks that emphasize common sense reasoning
1253 (e.g., chemistry), their success rates plummet when
1254 faced with tasks that heavily rely on environmen-
1255 tal interaction (e.g., classification). In fact, many
1256 of these tasks remain unsolvable for prompt-based
1257 approaches. In contrast, our method demonstrates

exceptional completion rates across both seen and
1258 unseen tasks, highlighting its robustness and adapt-
1259 ability in complex and varied scenarios. 1260

1261 Table 5 presents a comprehensive comparison of
1262 our Interactive Planner-R1 against existing meth-
1263 ods on both ALFWorld and ScienceWorld bench-
1264 marks. The table is organized into two main sec-
1265 tions to ensure fair evaluation: referenced results
1266 from original papers (marked with *) and our own
1267 implementations using identical experimental set-
1268 tings. The referenced results section includes a
1269 wide range of state-of-the-art methods, from tra-
1270 ditional approaches like CoLA and IPR to more
1271 recent advances such as CoPS and AgentPRM,
1272 with some methods achieving impressive perfor-
1273 mance on ALFWorld (e.g., CoPS reaching 94.0%
1274 average completion rate). However, it is impor-
1275 tant to note that these results may not be directly

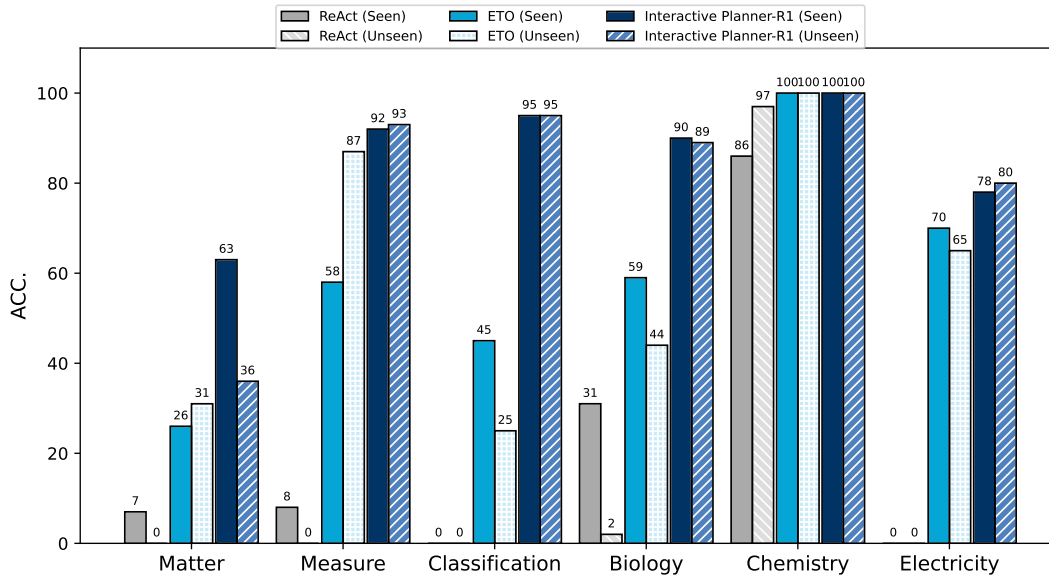
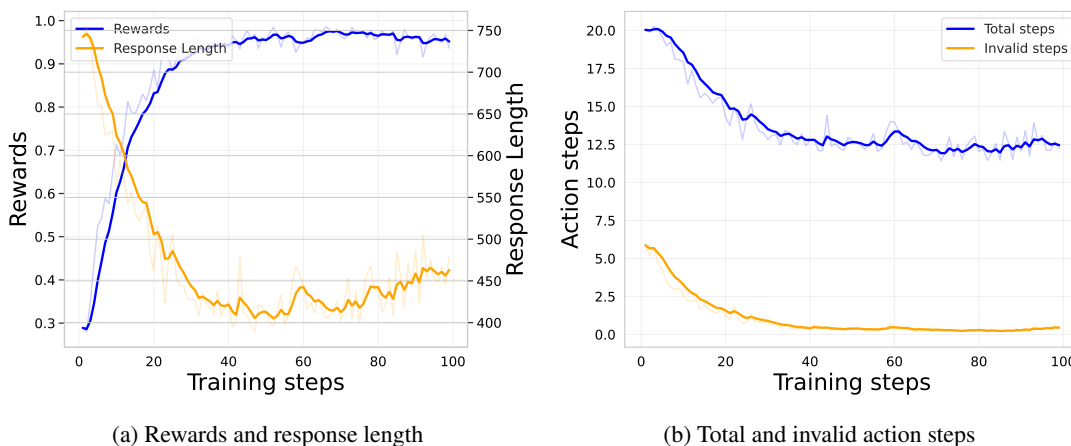


Figure 8: Completion rates of Qwen2.5-7B-Instruct on various ScienceWorld tasks, covering both seen and unseen scenarios. We compare the prompt-based method ReAct with the training-based method ETO and Interactive Planner-R1.



(a) Rewards and response length

(b) Total and invalid action steps

Figure 9: Training progress of Interactive Planner-R1 in ALFWorld: (a) shows the rewards (task completion rate) and total response length in multi-round exploration; (b) illustrates the decreasing trends in both total and invalid action steps during training.

comparable due to potential differences in experimental configurations, evaluation protocols, or implementation details. Our implementation section provides a controlled comparison using consistent experimental settings across all methods. The results demonstrate that vanilla large language models with ReAct prompting generally struggle on these tasks, with even GPT-4o achieving only 55.41% on ALFWorld and 58.18% on ScienceWorld. Training-based methods show significant improvements, with ETO achieving 82.07% on ALFWorld, though its performance on ScienceWorld remains limited at 56.20%. In contrast, our

Interactive Planner-R1 achieves remarkable performance across both benchmarks, reaching 97.55% average completion rate on ALFWorld and 79.92% on ScienceWorld. Notably, our method demonstrates strong generalization capabilities, maintaining high performance on both seen and unseen scenarios, with minimal performance degradation between the two settings. This consistent performance across different task distributions highlights the robustness and effectiveness of our approach in complex interactive environments.

G More ablation study

Method	ALFWorld			ScienceWorld		
	Seen	Unseen	Avg.	Seen	Unseen	Avg.
<i>Referenced Results (Not Directly Comparable)</i>						
KnowAgent* (Zhu et al., 2025)	66.71	62.69	64.70	58.67	49.18	53.93
CoLA* (Jia et al., 2025)	77.90	74.60	76.30	28.40	21.80	25.10
IPR* (Xiong et al., 2024)	70.30	74.70	72.50	-	-	-
STeCa* (Wang et al., 2025)	74.30	76.10	75.20	-	-	-
Agent-R* (Yuan et al., 2025)	-	-	-	-	-	70.23
DEP* (Wang et al., 2023)	-	-	76.00	-	-	-
WKM* (Qiao et al., 2024)	68.57	65.93	67.25	58.67	49.18	53.93
OREO* (Wang et al., 2024a)	79.10	80.70	79.90	-	-	-
QLASS* (Lin et al., 2025)	77.90	82.80	80.40	75.30	66.40	70.90
AgentLM* (Zeng et al., 2024b)	-	-	86.00	-	-	20.80
AgentPRM* (Choudhury, 2025)	-	-	91.00	-	-	-
AdaPlanner* (Sun et al., 2023a)	-	-	91.79	-	-	-
CoPS* (Yang et al., 2024b)	-	-	94.00	-	-	-
<i>Our Implementation</i>						
GPT-3.5-Turbo†(Brown et al., 2020)	7.14	7.46	7.30	28.03	21.63	24.83
GPT-4o†(OpenAI, 2023)	58.57	52.24	55.41	59.6	56.75	58.18
Deepseek-v3†(DeepSeek-AI et al., 2024)	36.43	31.34	33.89	29.28	27.45	28.37
InternLM 2.5 20B Chat†(Team et al., 2024)	7.86	11.94	9.90	33.66	31.15	32.41
Qwen2.5-7B-Instruct†(Yang et al., 2024a)	30.00	32.09	31.05	28.01	16.09	22.05
Qwen2.5-7B-Instruct + SFT(Zeng et al., 2024a)	80.71	79.10	79.91	68.81	55.7	62.26
Qwen2.5-7B-Instruct + NAT(Wang et al., 2024b)	63.57	66.42	65.00	58.56	49.75	54.16
Qwen2.5-7B-Instruct + ETO(Song et al., 2024)	84.29	79.85	82.07	55.08	57.32	56.20
Qwen2.5-7B-Instruct + Interactive Planner-R1(Ours)	98.57 ± 0.7	96.52 ± 1.9	97.55	83.66 ± 0.5	76.18 ± 0.9	79.92

Table 5: **Complete experimental results on ALFWorld and ScienceWorld benchmarks.** The best results are marked in **bold**. † indicates vanilla models tested with ReAct prompting, and * denotes the best results reported in the original papers.

Group Size	Seen	Unseen	Avg.
2	97.86 ± 0.4	92.68 ± 0.8	95.27
3	98.14 ± 0.4	97.31 ± 0.4	97.73
5	98.57 ± 0.7	96.52 ± 1.9	97.55
7	98.93 ± 0.5	95.9 ± 0.5	97.41

Table 6: Performance across different group sizes on ALFWorld. Group size 5 provides the optimal balance between exploration diversity and computational efficiency.

Group Size Analysis The group rollout mechanism is central to Interactive Planner-R1’s exploration strategy. We systematically examine how group size affects both exploration diversity and computational efficiency by varying the number of parallel trajectories from 2 to 7 and shown in Table 6.

As the Table 6 shown, the results reveal an interesting pattern: performance initially improves with group size (from 95.27% at size 2 to 97.73% at size 3), suggesting that increased trajectory diver-

Temperature	Seen	Unseen	Avg.
0.5	96.43 ± 2	94.78 ± 0.3	95.61
0.7	97.86 ± 2	97.02 ± 1.1	97.44
1.0	98.57 ± 0.7	96.52 ± 1.9	97.55

Table 7: Impact of temperature parameter on ALFWorld performance. Higher temperature enables better exploration while maintaining action quality.

sity enhances learning quality. However, beyond group size 5, we observe diminishing returns, particularly on unseen environments (95.9% for size 7 vs. 96.52% for size 5). This suggests that while larger groups provide more exploration diversity, they may also introduce noise that interferes with policy learning. Group size 5 emerges as the optimal configuration, balancing exploration breadth with learning stability.

Temperature Sensitivity Analysis Sampling temperature critically controls the exploration-exploitation trade-off during trajectory generation. We evaluate three temperature values to understand

1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324

1325 how randomness affects both exploration quality
 1326 and final performance.

1327 Temperature 1.0 achieves the highest overall per-
 1328 formance (97.55%), indicating that sufficient explo-
 1329 ration randomness is essential for discovering ef-
 1330 fective action sequences. Lower temperatures (0.5
 1331 and 0.7) result in more conservative exploration,
 1332 limiting the model’s ability to discover optimal
 1333 strategies. Interestingly, temperature 0.7 shows the
 1334 best generalization (only 0.84% gap between seen
 1335 and unseen), while temperature 1.0 provides the
 1336 best overall performance despite a larger general-
 1337 ization gap (2.05%).

KL Weight	Seen	Unseen	Avg.
0	99.29 ± 1	96.27 ± 2.1	97.78
0.001	98.57 ± 0.7	96.52 ± 1.9	97.55
0.1	97.15 ± 1	97.01 ± 1	97.08

Table 8: Effect of KL penalty on ALFWorld perfor-
 mance. Removing KL constraints yields optimal per-
 formance, suggesting sparse rewards provide sufficient
 regularization.

1338 **KL Divergence Constraint Effects** Kullback-
 1339 Leibler (KL) penalty is commonly used in RL to
 1340 prevent the policy from deviating too far from the
 1341 reference model. We investigate whether such con-
 1342 straints are necessary given our sparse reward struc-
 1343 ture.

1344 As shown in Table 8, removing KL penalty en-
 1345 tirely ($\beta = 0$ in Equation 7) yields the highest
 1346 performance (97.78%), outperforming both light
 1347 ($\beta = 0.001$) and moderate ($\beta = 0.1$) regular-
 1348 ization. This finding suggests that our sparse,
 1349 completion-based reward structure inherently pro-
 1350 vides sufficient regularization against reward hack-
 1351 ing and policy drift. The strong KL constraint (0.1)
 1352 actually hurts performance by overly restricting
 1353 policy updates, preventing the model from learning
 1354 effective exploration strategies. However, we ob-
 1355 served some degree of instability in unconstrained
 1356 training, manifested as significant training curve
 1357 fluctuations and inconsistent convergence rates.
 1358 Based on this finding, we adopted light KL reg-
 1359 ularization ($\beta = 0.001$) as a compromise in our
 1360 final implementation.

1361 H Error Analysis

1362 To better understand the behavioral differences be-
 1363 tween Interactive Planner-R1 and prior multi-turn

reinforcement learning baselines, we conduct a de-
 tailed error analysis against ETO under identical
 evaluation settings. The analysis is performed on
 the same set of tasks using fixed checkpoints for
 both methods, without any additional training.

For each failed trajectory, we manually assign
 one or more high-level error categories, including:
 (i) execution errors after sufficient information has
 been obtained, (ii) oscillatory or contradictory ac-
 tion sequences, (iii) goal misunderstanding after
 locating key objects, (iv) insufficient exploration
 (failing to locate required objects or locations), (v)
 inefficient execution, (vi) parsing errors, and (vii)
 unknown failure causes. The aggregated statistics
 are summarized in Table 9.

As shown in Table 9, failures of ETO are dom-
 inated by execution-stage errors. A substantial
 portion of failed trajectories execute incorrect ac-
 tion sequences despite already observing the key
 object or location, while many others exhibit re-
 peated or oscillatory behaviors. This indicates that
 ETO often struggles to convert acquired informa-
 tion into stable and coherent multi-step plans. In
 contrast, Interactive Planner-R1 significantly re-
 duces execution-stage failures, with only a small
 fraction of errors occurring after sufficient informa-
 tion has been obtained. Moreover, the proportion
 of unknown failure cases is lower, suggesting more
 interpretable and structured failure modes.

The remaining failures of Interactive Planner-
 R1 are primarily attributed to occasional goal mis-
 understanding or insufficient exploration, rather
 than breakdowns in planning or execution once
 the relevant information is available. This analysis
 suggests that the performance gains of Interactive
 Planner-R1 mainly stem from improved planning
 stability and execution quality under partial observ-
 ability, rather than solely from increased explo-
 ration.

I Case Study: Model’s Knowledge about Environment

Beyond the model’s successful generalization to
 both seen and unseen examples, which already in-
 dicates its understanding of environmental knowl-
 edge, we conducted a case study in ALFWorld to
 further compare the differences in environmental
 knowledge mastery between Interactive Planner-R1
 and the SFT model.

To demonstrate that the model has acquired en-
 vironmental knowledge through exploration, we

Error Type	ETO	Planner-R1
Executed wrong action sequence after having enough information	9 (0.33)	1 (0.20)
Repeated contradictory or oscillatory actions	9 (0.33)	2 (0.40)
Found key object but misunderstood the goal	6 (0.22)	2 (0.40)
Never found key object / location (insufficient exploration)	6 (0.22)	2 (0.40)
Insufficient exploration (generic)	3 (0.11)	–
Inefficient execution	1 (0.04)	–
Misunderstood goal	1 (0.04)	–
Parsing error	1 (0.04)	–

Table 9: Breakdown of failure modes for ETO and Interactive Planner-R1. Each failed trajectory may be assigned to multiple error categories; therefore, percentages do not sum to 1. Values are reported as number of trajectories followed by the ratio in parentheses.

Search CD with Original Task Format	
Interactive Planner R1	SFT
<p>System: You are an intelligent agent in a household environment...</p> <p>User: Your task is to: put some cd on sidetable.</p> <p>Observation: <Observation from Environment ></p> <p>Assistant: Thought: The task is to put a cd on the sidetable. I need to first find a cd, then take it to one of the sidetables. Let's start by checking the diningtable and drawers as CDs are commonly found there. Action: go to diningtable 1</p>	<p>System: You are an intelligent agent in a household environment...</p> <p>User: Your task is to: put some cd on sidetable.</p> <p>Observation: <Observation from Environment ></p> <p>Assistant: Thought: I need to locate a CD to place on the side table as per the task. I should look around the room to find one. Action: go to diningtable 1</p>

Figure 10: Comparison between Interactive Planner-R1 and SFT using the original task format, where Interactive Planner-R1 pointed out the possible position of CD and SFT exhibited inconsistency between thought and action .

Ask Model's Knowledge about CD	
Interactive Planner R1	SFT
<p>System: You are an intelligent agent in a household environment...</p> <p>User: Your task is to: answer where is cd most likely to appear.</p> <p>Observation: <Observation from Environment ></p> <p>Assistant: Thought: CDs are typically found in places like a dining table or a side table in a room. I will check the diningtable 1 and sidetable 1 first. Action: go to diningtable 1</p>	<p>System: You are an intelligent agent in a household environment...</p> <p>User: Your task is to: answer where is cd most likely to appear.</p> <p>Observation: <Observation from Environment ></p> <p>Assistant: Thought: I need to locate a CD. I will start by checking the nearest surfaces which could hold small items, such as side tables. Action: go to sidetable 1</p>

Figure 11: Comparison between Interactive Planner-R1 and SFT when asked about the location of CD, showing SFT suffered from prior knowledge bias while Interactive Planner-R1 still succeeded.

1414 identified an object, CD, that appears in numerous
1415 tasks positioned on dining tables, which contra-
1416 dicts common expectations where CDs would more
1417 likely be stored in drawers or on shelves.

1418 In Figure 10, we first maintained the input for-
1419 mat used during training and found that although
1420 both SFT and Interactive Planner-R1 models made
1421 correct actions, the SFT model's thoughts did not
1422 strongly correlate with its actions. In contrast, In-
1423 teractive Planner-R1 demonstrated awareness of
1424 CD locations in its thought process.

1425 To further verify that Interactive Planner-R1 has
1426 mastered knowledge related to the environment, we

1427 modified the input format as shown in the Figure 11,
1428 changing the task to ask the model where CDs are
1429 most likely to appear. The SFT model failed to
1430 complete this type of task, with its thoughts still
1431 reflecting the influence of prior knowledge without
1432 recognizing that CD locations in ALFWorld differ
1433 from common expectations. Interactive Planner-
1434 R1, however, still made the correct analysis in its
1435 thoughts, proving that the model has successfully
1436 captured this implicit spatial relationship during its
1437 exploration phase.