The Lighthouse of Language: Enhancing LLM Agents via Critique-Guided Improvement

Ruihan Yang^{1,♥}, Fanghua Ye♠ Jian Li♠ Jiyu Yuan^{1,♥}, Yikai Zhang^{2,♥},

Zhaopeng Tu♠, Xiaolong Li♠, Deqing Yang^{1,♥†}

1,♥ School of Data Science, Fudan University ♠ Tencent Hunyuan

2,♥ College of Computer Science and Artifcial Intelligence, Fudan University

(rhyang17, yangdeqing)@fudan.edu.cn ⟨ syyuan21, ykzhang22}@m.fudan.edu.cn

Project Page: https://github.com/rhyang2021/CGI

♠{fanghua.ye.21, lijianjack, tuzhaopeng}@gmail.com

Abstract

Large language models (LLMs) have recently transformed from text-based assistants to autonomous agents capable of planning, reasoning, and iteratively improving their actions. While numerical reward signals and verifiers can effectively rank candidate actions, they often provide limited contextual guidance. In contrast, natural language feedback better aligns with the generative capabilities of LLMs, providing richer and more actionable suggestions. However, parsing and implementing this feedback effectively can be challenging for LLM-based agents. In this work, we introduce Critique-Guided Improvement (CGI), a novel two-player framework, comprising an actor model that explores an environment and a critic model that generates detailed nature language feedback. By training the critic to produce fine-grained assessments and actionable revisions, and the actor to utilize these critiques, our approach promotes more robust exploration of alternative strategies. Experiments in three interactive environments show that CGI outperforms existing baselines by a substantial margin. Notably, even a small critic model surpasses GPT-4 in feedback quality. The resulting actor achieves stateof-the-art performance, demonstrating the power of explicit guidance to enhance decision-making in LLM-based agents.

1 Introduction

Recent advances in large language models (LLMs) [1; 2; 3] have broadened their applicability from conventional chat and text-based interactions to more dynamic, agentic roles. In these roles, LLM-based agents [4; 5; 6] harness the reasoning and planning capabilities of LLMs to automate complex tasks across diverse domains, including code generation [7; 8], software engineering [9; 10], and web applications [11; 12]. A key requirement in these agentic tasks is the ability to iteratively acquire, store, and use new information to improve performance. Thus, a central challenge emerges: *How can high-quality feedback be obtained and utilized efficiently?*

A common approach for providing feedback is to rely on numerical signals, such as verifiers [13; 14] or reward models [15; 16]. In this paradigm, LLMs are typically trained as discriminators to evaluate and select the best action from a set of candidates. One widely used method is Best-of-N (BoN) [17; 18], where the agent generates N possible actions and a trained reward model picks the

^{*}Work done during an internship at Tencent Hunyuan.

[†]Corresponding authors.

highest-scoring one. However, numerical feedback has limited informativeness: it reduces guidance to scalar scores that merely filter candidate actions, overlooking context-specific suggestions, avenues for exploration, or strategies for refining suboptimal behavior.

A more direct and flexible approach involves using natural language feedback. Rather than providing only a "best" action, LLMs can generate detailed critiques and explicit, context-aware recommendations for improvement. However, transitioning to natural language feedback presents two main challenges: 1) Weak Feedback: Many techniques rely on self-refinement or self-correction [19; 20; 21], which heavily depend on the LLM's inherent capacity. This reliance can lead to degraded performance, especially when models hallucinate or encounter complex problems [20; 22]; 2) Poor Utilization: Unlike numerical feedback, which directly selects an action to execute, verbal critiques require the agent to accurately interpret and integrate the guidance into its decision-making process. LLM agents often struggle with this, sometimes demonstrating limited flexibility in adjusting or discarding previously generated plans [23; 24].

In response, we propose Critique-Guided Improvement (CGI), a two-player framework in which an actor model interacts with the environment while a critic model provides supervisory feedback. CGI features two main stages: Critique Generation and Action Refinement. In the Critique Generation stage, we train a critic model to produce precise evaluations and actionable revision instructions. In the Action Refinement stage, the actor model learns to effectively apply these critiques through iterative supervised fine-tuning, thereby improving both its reasoning and its ability to integrate external feedback.

We conduct extensive experiments in three interactive environments [11; 25; 26], demonstrating that CGI substantially enhances performance over baseline approaches. Notably, our critic model, when trained on a small backbone (e.g., Llama-3-8B), outperforms even GPT-4 as a critic by a large margin (+29.16%). Additionally, the action-refinement process further boosts results, surpassing state-of-the-art metrics by 26.74%.

In summary, our contributions are as follows: 1) We tackle the challenge of obtaining and utilizing high-quality nature language feedback in agentic tasks. 2) We present Critique-Guided Improvement (CGI), a novel two-player framework that provides more informative feedback to LLM-based agents by combining a specialized critic model with an actor model. 3) We develop a training strategy that iteratively refines the actor's actions based on natural language critiques, enabling both improved reasoning and enhanced adaptability. 4) We conduct extensive experiments across three diverse interactive environments. By leveraging CGI, the agent can continuously improve its performance in long-horizon tasks and enhance task efficiency by achieving higher scores in fewer steps.

2 Related Work

Learning from Feedback Current feedback can be categorized into numerical and natural language feedback. Numerical feedback [27; 28; 29; 30] is typically provided by training a reward model (RM) or verifiers. In this approach, the model learns to predict either the correctness of a solution [13; 29] or the preference between multiple solutions [17; 31]. Specifically, the RM generates a continuous numerical score, which is then incorporated into a classification objective. In contrast, natural language feedback offers denser rewards by evaluating the model's actions using natural language. This feedback can be generated in two main ways: by prompting an off-the-shelf LLM to act as a verifier (*i.e.*, LLM-as-judge) [32; 33], or through self-refinement and self-critique [20; 34], where the model evaluates and improves its own outputs. Our trained critic model provides natural language feedback, offering denser and more accurate rewards.

Agent Learning in Interactive Environments. Previous approaches to agent learning in interactive environments can be classified into three main categories: *1)* Prompt-based methods [21; 35; 36; 37] utilize human-written prompts to guide LLMs in summarizing experiences. These summaries, which may include causal abstractions from both successful and failed attempts [20; 38] or transferable skills [39], are integrated into the model's memory to enhance its knowledge and performance. *2)* Training-based methods [40; 12; 41] rely on techniques such as Supervised Fine-Tuning (SFT) [42; 43] or Direct Preference Optimization (DPO) [44; 45]. to train LLMs. The training data may come from expert models or be generated through exploration strategies like Monte Carlo Tree Search (MCTS). *3)* Inference-time sampling methods [14; 15; 16] employ techniques such as Best-of-N

(BoN) [17; 18] and Tree-of-Thought (ToT) [46] to identify optimal actions during inference. These methods leverage the prior knowledge in LLMs, enabling more efficient search processes. Our CGI method introduces a novel variation of inference-time sampling by using a trained critic that evaluates candidate actions and suggests revisions to enhance inference performance.

3 Preliminary

Partially Observable Markov Decision Process We define the collection of environments as \mathcal{E} . For a specific environment $e \in \mathcal{E}$, the embodied tasks for LLM agents are typically modeled as a Partially Observable Markov Decision Process (POMDP): $(\mathcal{X}, \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T})_e$. Here, \mathcal{X} denotes the set of instructions, \mathcal{S} represents the set of environment states, \mathcal{A} is the set of available actions at each state, and \mathcal{O} represents the observations available to the agent. The transition function $\mathcal{T}: \mathcal{S} \times \mathcal{A} \to \mathcal{S}$ is defined by the environment, while the reward function $\mathcal{R}: \mathcal{S} \times \mathcal{A} \to [0,1]$ specifies the reward received by the agent from the environment. For an actor model π_{θ} parameterized by θ , the actor selects the next action $a_t \sim \pi_{\theta}(a|\tau_t,e)$ at each time step t, based on the interaction history τ_t , which is defined as

$$\tau_t = \{x, a_0, o_0, \dots, a_{t-1}, o_{t-1}\}, \quad \tau_0 = \{x\}.$$

The trajectory is then represented as

$$\tau = (x, a_0, o_0, \dots, a_T, o_T) \sim \pi_{\theta}(\tau \mid x, e).$$

Iterative Supervised Fine-Tuning in Agentic Task Iterative Supervised Fine-Tuning (SFT) is a process of exploration and learning [47; 48; 49]. It iteratively utilizes the correct responses from the actor model to enhance the model's problem-solving abilities. The process involves K iterations, each consisting of two steps: exploration and learning. For the collection of environments \mathcal{E} , in the exploration step of iteration k, the model π_{θ}^{k-1} from the previous iteration is applied to each environment e, resulting in $|\mathcal{E}|$ trajectories $\{\tau^{(j)}\}_{j=1}^{|\mathcal{E}|}$. These trajectories are then filtered using the reward function $\mathcal{R}(\tau)$, retaining only the correct ones. The filtered trajectories form a new dataset $\mathcal{D}_{\text{correct}} = \{\tau^{(j)}\}_{j=1}^{N'}$, where N' is the number of trajectories retained. In the learning step of iteration k, this new dataset is used to fine-tune the actor model π_{θ} to obtain π_{θ}^k .

4 Methodology

In this section, we first provide an overview of the proposed Critique-Guided Improvement (CGI) framework (§ 4.1), explaining how the actor and critic collaborate to enhance performance. We then detail the two key stages of CGI, critique generation and action refinement, which respectively address the challenges of critic's *weak feedback* and the actor's *poor utilization*. In the **Critique Generation** stage (§ 4.2), the critic model learns to evaluate the actor's candidate actions and provide actionable feedback. In the **Action Refinement** stage (§ 4.3), the actor refines its actions by integrating critiques from the critic model into its interaction with the environment. Algorithm 1 summarizes the CGI framework (see Appendix C for definitions of all notifications).

4.1 Overview of the CGI Framework

The CGI framework follows a two-player setting, where the actor model π_{θ} generates multiple candidate actions, and the critic model π_{ϕ} provides feedback to refine them. For a specific environment $e \in \mathcal{E}$, at each time step t, given the refined trajectory history

$$\tau'_t = \{x, a'_0, o_0, \dots, a'_{t-1}, o_{t-1}\},\$$

the actor generates M candidate actions, stored in the action buffer $A_t = \{a_{t,i}\}_{i=1}^M$. The critic then evaluates each candidate and generates corresponding critiques $c_t = \{\pi_\phi(c \mid \tau_t', a_{t,i}, e)\}_{i=1}^M$. Using these critiques, the actor refines its decision and generates the final action $a_t' = \pi_\theta(a \mid \tau_t', c_t, e)$, which is then executed in the environment. This iterative process produces the refined trajectory:

$$\tau' = \{x, a_0', o_0, \dots, a_T', o_T\}.$$

In this collaborative framework, higher-quality critiques c_t and better utilization of feedback lead to improved refined actions a'_t , ultimately enhancing the final performance $\mathcal{R}(\tau')$.

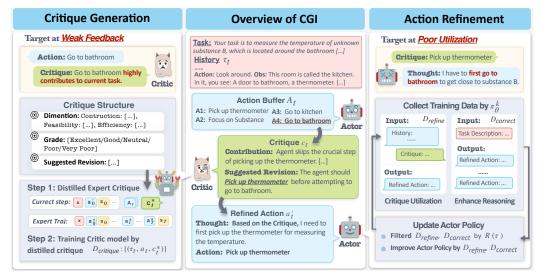


Figure 1: An overview of CGI, illustrated with a ScienceWorld example. The actor first generates candidate actions, and the critic then assesses each action and generates suggested revisions. The actor refines its actions based on the critique. The **Critique Generation** stage enhances the critic model's ability to provide effective feedback, while the **Action Refinement** stage helps the actor model better utilize the critique.

4.2 Critique Generation

The **Critique Generation** stage trains the critic model to assess the alignment between LLM-generated actions and optimal ones at each step, and convert this assessment into structured natural language critiques. We first define the critique structure, then outline the training method for critic.

Critique Structure The critique generated by the critic model π_{ϕ} consists of two components: *Discrimination*, which involves analyzing and assessing the quality of the candidate action, and *Revision*, which provides actionable suggestions for improvement.

For the discrimination component, the critic model evaluates the candidate action in three predefined dimensions:

- Contribution: Assessing how well the candidate action contributes to solving the overall task. Irrelevant actions receive a low score.
- Feasibility: Determining whether the candidate action is valid according to the agent's predefined action list.
- Efficiency: Evaluating whether the candidate action approaches the task optimally, avoiding unnecessary steps or redundancy.

For the revision component, the critic model first assigns an overall grade to each candidate action (*i.e.*, [Excellent/Good/Neutral/Poor/Very Poor]), then generates concise and actionable suggestions based on its analysis of the candidate actions (see Appendix D for detailed prompts).

Fine-tuning the Critic Model To collect high-quality step-level expert critiques, we employ an expert critic annotator (e.g., GPT-4 [1]) to generate critiques for each candidate action based on a reference expert trajectory. For a specific environment $e \in \mathcal{E}$, at time step t, the actor generates candidate actions a_t . Given the refined history trajectory τ_t' and the expert trajectory τ^{exp} , the expert critic π_{exp} assesses the alignment between candidate actions and the optimal ones, converting this assessment into structured critiques according to a predefined format. Formally, the critiques for each candidate action are represented as: $c_t \sim \pi_{\text{exp}}(c \mid a_t, \tau_t', \tau^{\text{exp}}, e)$.

Subsequently, the actor refines its next step by generating a'_t based on the critique c_t . Critiques are collected only when $\mathcal{R}(\tau') = 1$, ensuring they originate from correct trajectories. We collect step-level

Algorithm 1 Critique-Guided Improvement (CGI)

1: **Input:** Initialized policy actor model π_{θ} , critic model π_{ϕ} , environment set \mathcal{E} , full instruction set \mathcal{X} , reward function \mathcal{R} , the expert critiques $\mathcal{D}_{\text{critique}}$, the expert path $\mathcal{D}_{\text{expert}}$, the general dataset $\mathcal{D}_{\text{general}}$, max time step T, and the iteration rounds K.

2: Procedure Fine-tune the Critic Model:

Minimize the following objective to obtain the critic model π_{ϕ} :

$$\mathcal{L}_{\text{critic}}(\phi) = \mathbb{E}_{(c_t, \tau_t', a_t, e) \sim \mathcal{D}_{\text{critique}}} \left[\log \pi_{\phi} \left(c_t \mid \tau_t', a_t, e \right) \right];$$

4: Procedure Iterative Action Refinement:

```
5: \pi_{\theta}^{0} \leftarrow \pi_{\theta_{base}};
6: for iteration k=1 to K do
```

 $D_{\text{train}} \leftarrow D_{\text{expert}}$; // Initialize $\mathcal{D}_{\text{train}}$ with D_{expert} . 7:

Perform Exploration Step 8:

9:

Actor
$$\pi_{\theta}^{k-1}$$
 interacts with $e \in \mathcal{E}$ under the guidance of critic π_{ϕ} ;

9: // Collect critique-action pairs with $\mathcal{R}(\tau') = 1$.

10: Collect $D_{\text{refine}} = \bigcup_{e \in \mathcal{E}} D_{\text{refine}}^e$, where $D_{\text{refine}}^e = \{(\tau'_t, c_t, a'_t)\}_{t=1}^T$;

// Update training set with correct trajectory. 11:

12: Collect
$$D_{\text{correct}} = \bigcup_{e \in \mathcal{E}} D_{\text{correct}}^e, D_{\text{train}} \leftarrow D_{\text{train}} \cup \mathcal{D}_{\text{correct}};$$

13: **Perform** Learning Step

Minimize the following objective to obtain actor model π_{α}^{k} : 14:

Minimize the following objective to obtain actor model
$$\pi_{\theta}^{\kappa}$$
:
$$\mathcal{L}_{\text{actor}}(\theta) = \beta \left\{ \mathbb{E}_{(\tau, x, e) \sim \mathcal{D}_{\text{train}}} \left[\log \pi_{\theta}(\tau \mid x, e) \right] + \mathbb{E}_{(a'_{t}, \tau'_{t}, c_{t}, e) \sim \mathcal{D}_{\text{refine}}} \left[\log \pi_{\theta}(a'_{t} \mid \tau'_{t}, c_{t}, e) \right] \right\} + (1 - \beta) \mathbb{E}_{(x, y) \sim \mathcal{D}_{\text{general}}} \left[\log \pi_{\theta}(y \mid x) \right];$$

15: end for

expert critiques from each environment, forming the dataset $\mathcal{D}_{\text{critique}} = \cup_{e \in \mathcal{E}} \mathcal{D}^e_{\text{critique}}$. The critic model is then fine-tuned using supervised learning with the collected expert critiques. Specifically, we apply the standard language modeling loss, defined as:

$$\mathcal{L}_{\text{critic}}(\phi) = \mathbb{E}_{(c_t, \tau_t', a_t, e) \sim \mathcal{D}_{\text{critique}}} \Big[\log \pi_{\phi} \left(c_t \mid \tau_t', a_t, e \right) \Big].$$

This approach enables the critic model to generate structured, step-level critiques that provide both discrimination and revision, which are crucial for guiding the actor model toward more effective decision-making.

4.3 Action Refinement

Although the critic model is trained to provide high-quality critiques, the actor may not fully utilize the feedback. Therefore, action refinement is necessary to enhance the actor's ability to effectively leverage critiques for improving its actions. A key challenge in this process is policy misalignment. During training, the actor model learns to incorporate critiques based on its current policy. However, after training, the model's policy may evolve, making it difficult to integrate critiques effectively for newly generated candidate actions. This misalignment can lead to suboptimal action refinement.

To address this issue, we propose an iterative action refinement method based on supervised finetuning (SFT). This method consists of two main components: exploration and learning. In the exploration step at iteration k, for a specific environment $e \in \mathcal{E}$, the actor model π_{θ}^{k-1} interacts with the environment under the guidance of the critic model π_{ϕ} . At each time step \check{t} , the refined action is given by $a'_t = \pi_{\theta}(a \mid \tau'_t, c_t, e)$, forming a set of critique-action pairs $\{(\tau'_t, c_t, a'_t)\}_{t=1}^T$. Here, τ_t' represents the previously refined trajectory, and c_t denotes the current critique. To ensure the quality of the trajectories, we filter them based on the environment reward \mathcal{R} . Only trajectories for which $\mathcal{R}(\tau') = 1$, along with their corresponding critique-action pairs, are retained. We collect correct trajectories and critique-action pairs from each environment, yielding two datasets: $\mathcal{D}_{correct}$ $\cup_{e \in \mathcal{E}} \mathcal{D}_{\text{correct}}^e$ and $\mathcal{D}_{\text{refine}} = \cup_{e \in \mathcal{E}} \mathcal{D}_{\text{refine}}^e$. The dataset $\mathcal{D}_{\text{correct}}$ enhances the model's reasoning ability by providing correct trajectories, while $\mathcal{D}_{\text{refine}}$ improves its capacity to utilize critiques and generate refined actions.

In the learning step of iteration k, the newly collected datasets, combined with the expert path $\mathcal{D}_{\text{expert}}$ for each environment, are used to fine-tune the actor model π_{θ} . To avoid overfitting, we follow previous work [50] by fine-tuning the original model π_{θ} rather than the previous iteration model π_{θ}^{k-1} . Furthermore, following AgentTuning [51], we incorporate general datasets such as ShareGPT3³ to improve generalization. The training objective is defined as:

$$\mathcal{L}_{\text{actor}}(\theta) = \beta \left\{ \mathbb{E}_{(\tau, x, e) \sim \mathcal{D}_{\text{train}}} \left[\log \pi_{\theta}(\tau \mid x, e) \right] + \mathbb{E}_{(a'_t, \tau'_t, c_t, e) \sim \mathcal{D}_{\text{refine}}} \left[\log \pi_{\theta}(a'_t \mid \tau'_t, c_t, e) \right] \right\}$$

$$+ (1 - \beta) \mathbb{E}_{(x, y) \sim \mathcal{D}_{\text{general}}} \left[\log \pi_{\theta}(y \mid x) \right],$$

where $\mathcal{D}_{train} = \mathcal{D}_{expert} \cup \mathcal{D}_{correct}$. After this refinement step, a new dataset with higher-quality samples is generated for further training. Overall, this iterative action refinement process allows the actor model to progressively improve its reasoning capabilities and better integrate critiques through continuous interaction with the environment.

5 Experiment Settings

In this section, we conduct extensive experiments in three interactive environments to demonstrate the effectiveness of the critic model and our CGI framework.

5.1 Interactive and Agentic Environments

Following previous work [40; 52], we conduct experiments on three types of representative interactive environments:

- WebShop [11], which is an interactive web environment for online shopping. It contains 12K instructions and offers over one million real products from amazon.com. Agents can click buttons on the webpage or perform searches using the search engine.
- ScienceWorld [25], which is a text-based scientific environment designed to evaluate agents' scientific reasoning abilities. It includes 30 types of scientific tasks at the standard elementary science curriculum level.
- TextCraft [26], which is a text-based environment to create Minecraft items. It constructs a crafting tree based on Minecraft's recipes. Each task provides a target item and a list of crafting commands generated by the tree. Agents receive a reward of 1 when they successfully craft the target item.

Evaluation Metrics Following the setup of AgentGym [40], we evaluate our model on the test sets for these three environments (200 simulations for ScienceWorld and WebShop, 100 for TextCraft). For ScienceWorld and WebShop, we use the average final score as the evaluation metric. For TextCraft, we use the success rate as the evaluation metric. Further details can be found in Appendix E.

5.2 Training Settings

We use Llama-3-8B-Instruct [3] as the backbone model for both the actor and critic models. To collect training data, we randomly sample 500 simulations from WebShop, 350 from ScienceWorld, and 374 from TextCraft. We train the critic model by using the expert critic (*i.e.*, GPT-40) to guide the actor in interacting with the environment three times, collecting expert critiques during each interaction. For action refinement, we perform three iterations and report the results of the third iteration in Table 2. The training data for the critic model and each iteration of action refinement are provided in Table 4. Additional details can be found in Appendix F.

5.3 Baselines

To evaluate the effectiveness of our critic model, we fix the actor model as Llama-3-8B-Instruct [3]. We then compare our critic model against two types of approaches (see Appendix G for implementation details): 1) **Numerical based**: We use DGAP [14], a discriminator trained to assess the alignment between actor actions and expert actions at the step level, and Explicit RM [15], which

³https://huggingface.co/datasets/Vtuber-plan/sharegpt-cleaned

Table 1: We compared the natural language feedback provided by the trained critic model with other methods, including numerical-based (*i.e.*, DGAP, Explict RM), and verbal-based (*i.e.*, self-critique, GPT-40). Our critic offers better guidance to the actor model in all three interactive scenarios.

Model	Method	WebShop	ScienceWorld	TextCraft	Average
Llama-3-8B-Instruct	No Critique	13.49	14.48	10.00	12.65
	DGAP	30.41	19.52	21.00	23.64
	Explicit RM	14.21	18.41	15.00	15.87
Liama-3-0B-Instruct	Self-Critique	1.50	10.06	19.00	10.19
	GPT-4o	17.78	33.06	46.00	32.28
	Critic Model (Ours)	56.80	68.51	59.00	61.44
	No Critique	8.35	49.20	2.00	19.85
	DGAP	11.17	55.86	9.00	25.34
Llama-3-70B-Instruct	Explicit RM	11.18	48.41	14.00	24.53
Liama-5-/0B-Instruct	Self-Critique	3.00	28.17	40.00	23.72
	GPT-4o	16.24	43.73	56.00	38.65
	Critic Model (Ours)	52.20	72.44	73.00	65.88
Llama-3-8B-Instruct + SFT	No Critique	76.12	32.85	46.00	51.66
	DGAP	73.97	38.26	56.00	56.08
	Explicit RM	74.33	41.57	52.00	55.97
	Self-Critique	50.18	31.22	21.00	34.13
	GPT-4o	55.65	48.48	58.00	54.04
	Critic Model (Ours)	<u>74.68</u>	55.94	<u>56.00</u>	62.21

is trained to predict the Q-value, *i.e.*, expected accumulated rewards at each time step. 2) **Verbal based**: We employ a self-critique method where the actor model itself generates critiques for each candidate action at the step level. Additionally, we use GPT-40 (gpt-4o-2024-08-06) [1] as a critic, which serves as a strong general-purpose evaluator. The structure of the critiques aligns with that in Appendix D. Each approach guides the actor's inference, and we evaluate their effectiveness based on the actor's performance⁴.

For the assessment of CGI, consistent with the methodology used in AgentGym [40], we select a combination of closed-source models, including GPT-3.5-turbo [53], GPT-4o [54], Claude 3 [55], and DeepSeek-Chat [56], as well as open-source models such as Llama-3-70B-Instruct [3]. Additionally, we include agents trained on expert trajectories, such as AgentLM (13B and 70B) [57] and Agent-Flan [42]. We also compare our approach to Iterative SFT, which iteratively refines the actor model using correct trajectories collected from interactions with the environment, and Reflexion [20], a self-reinforcement method that concludes each iteration with a summary to guide decision-making in subsequent iterations. We report the results from the third iteration for both approaches.

6 Main Results

We first evaluate our critic model against numerical- and verbal-based methods to assess its effectiveness. As shown in Table 1, our 8B critic model significantly outperforms GPT-40 on both the Llama-3 (8B and 70B) and Llama-3-8B's fine-tuned variant trained on expert data. Furthermore, Table 2 demonstrates that through iterative action refinement of the actor model, our CGI notably enhances actor performance in interactive environments. It outperforms both advanced closed-source models (e.g., GPT-40) and agents trained on expert trajectories (e.g., AgentLM-70B and Agent-FLAN). Based on these results, we identify three key findings.

Finding 1: Verbal critique feedback is more effective than numerical signal. As shown in Table 1, the Critic Model consistently outperforms numerical feedback from the discriminator. For the Llama-3-8B model, it achieves an average improvement of 42.89% over the no-critique baseline, while the discriminator yields only a 5.09% gain. The advantage persists with the stronger Llama-3-70B and the fine-tuned Llama-3-8B models, where the Critic Model delivers a 46.03% and 10.55% improvement, respectively, compared to 5.49% and 4.42% from the discriminator. These results indicate that the discriminator's numerical feedback relies heavily on the model's inherent capabilities, as it essentially performs action-level filtering. For example, in the ScienceWorld, the discriminator enables Llama-3-70B to reach 55.86%, but only 18.52% for Llama-3-8B. In contrast, the Critic

 $^{^4}$ In the subsequent experiments of this paper, we set the number of candidate actions sampled at inference time to M=5.

Table 2: Results from three interactive environments. The CGI here employs a two-player setting with a fine-tuned LLama-3-8B critic model and an 8B actor model refined through three iterations. We compare its performance with various models.

Method	WebShop	ScienceWorld	TextCraft	Average				
Closed-source Models								
DeepSeek-Chat	11.00	16.80	23.00	16.93				
Claude-3-Haiku	5.50	0.83	0.00	2.11				
Claude-3-Sonnet	1.50	2.78	38.00	14.09				
GPT-3.5-Turbo	12.50	7.64	47.00	22.38				
GPT-40	25.48	46.91	64.00	45.46				
Open-source Models								
Llama-3-70B-Instruct	8.35	49.20	2.00	19.85				
AgentLM-13B	39.50	2.75	0.00	14.08				
AgentLM-70B	49.50	10.68	4.00	21.39				
Agent-FLAN	40.35	28.64	16.00	28.33				
Llama-3-8B-Instruct	13.49	14.48	10.00	12.66				
w/ Reflexion	14.08	12.55	8.00	11.54				
w/ Iterative SFT	78.21	41.42	55.00	58.21				
w/ CGI (Ours)	76.17	78.43	68.00	74.20				

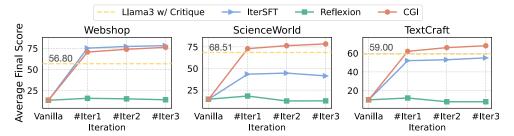


Figure 2: Performance of different iterative methods across iterations in three agentic tasks.

Model shows minimal sensitivity to the base model, achieving 68.51% with Llama-3-8B and 72.44% with Llama-3-70B. However, the self-critique approach, which also relies on verbal feedback, shows a negative effect on both Llama-3-70B, 8B and its fine-tuned variant. This suggests that while verbal feedback can convey richer information, untrained self-critique methods often produce low-quality or unstructured feedback, failing to provide effective guidance.

Finding 2: Fine-tuned models struggle to utilize critiques. While supervised fine-tuning significantly improves baseline performance (+33.11%), our experiments reveal a key limitation: fine-tuned models struggle to effectively incorporate critique feedback. In the ScienceWorld scenario, Vanilla Llama3-8B and Llama3-70B show substantial improvements with Critic guidance, achieving average scores of 68.51% and 72.44%, respectively. However, despite having a stronger base performance, fine-tuned Llama3-8B only reaches 55.94% with the same critique mechanism. This issue persists in the WebShop scenario, where critique guidance even causes a performance decline in fine-tuned Llama3-8B, dropping from 76.12% to 74.68%. These results suggest that while fine-tuning enhances base performance, it may reduce the model's responsiveness to external feedback. For example, in the ScienceWorld scenario, despite the critic model advising the actor to "look around" and assess the environment for possible actions, the actor persists in executing "Go to Bathroom", an action that seems closest to the task description. This behavior disregards the fact that the actor cannot reach the bathroom from its current position, leading to a deadlock (see Appendix H.1). Our iterative action refinement approach helps mitigate this issue; after three iterations, the model shows a significant improvement in critique utilization, with its performance surpassing that of the actor model with only supervised fine-tuning by +15.99%, achieving state-of-the-art results compared to various baselines.

Finding 3: CGI continuously enhance model performance via action refinement. As shown in Figure 2, compared to other iterative methods (*e.g.*, Reflexion and vanilla iterative SFT), CGI consistently supports model performance improvement. In contrast, Reflexion shows minimal improvement and can even cause performance degradation. This aligns with previous findings [20], where self-critique struggles to escape local minima in agentic tasks that require significant diversity and exploration. Iterative SFT achieves a notable improvement only in the first iteration, with minimal progress in the following 2-3 iterations. This issue is more pronounced in long-horizon

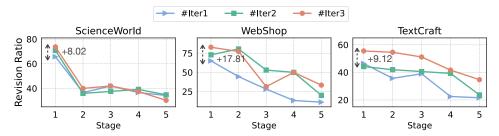


Figure 3: Revision Ration of actor model at different trajectory stages across three tasks. #Iter-*k* denotes the *k*-th iteration of action refinement. Revision Ratio represents the proportion of actions that differ after critique compared to those without critique.

tasks (e.g., ScienceWorld and TextCraft), suggesting that the model may over-sample simpler tasks while under-sampling more complex ones. Consequently, the training set for subsequent iterations becomes dominated by easier tasks, with few solutions for harder ones. As iterations progress, this bias deepens, leading to a long-tail distribution where solutions to more challenging tasks are scarce, ultimately causing the model to plateau or even degrade. In contrast, CGI leverages high-quality critiques and enhances the actor's utilization, enabling the model to consistently sample solutions to more difficult tasks, thus supporting sustained performance improvement.

7 Further Analysis

7.1 Qualitative Analysis of CGI

In this section, we conduct a qualitative analysis to examine how CGI enhances performance in agentic tasks. Specifically, we focus on two key aspects: 1) **Trajectory-level**: Identifying the stages in the trajectory where CGI provides the most significant improvement; 2) **Task-level**: Evaluating CGI's performance across tasks of varying trajectory length (see Appendix E for categorization).

CGI significantly improves early-stage performance. To determine at which stage of the trajectory CGI contributes most to performance improvements, we divide the trajectory into five stages based on its length and compute the extent to which the critique influences the actor model's behavior, measured by the Revision Ratio. As shown in Figure 3, across three tasks, actor models exhibit the highest action revision frequency in stage 1, with the revision ratio dropping sharply in later stages. This suggests that the critique primarily guides the actor during early exploration, helping reduce ineffective searches. Notably, as the actor undergoes more refinement iterations, its revision ratio in stage 1 increases (+8.02% on ScienceWorld, +17.81% on WebShop, +9.12% on TextCraft), indicating that better critique utilization accelerates effective exploration. This, in turn, enhances CGI's efficiency, enabling it to achieve higher scores in fewer steps (Figure 8).

CGI helps the model to continuously improve on long-horizon tasks. To investigate the effect of CGI on tasks of varying trajectory length, we categorize the ScienceWorld scenarios into three difficult-level groups based on the average length of the oracle agent's trajectories. Longer trajectories correspond to higher difficulty levels. As shown in Figure 4, the vanilla model (Llama3-8B) performs poorly, with performance decreasing as the task length increases. Critique-Guided Inference leads to a significant improvement, particularly on easy tasks, where performance increases by +59.64%. Although there is also a notable improvement on harder tasks (+43.95%), the gain is less pronounced compared to easier tasks. However, with iterative action refinement, the model exhibits the greatest improve-

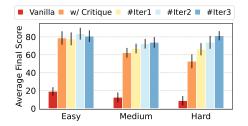


Figure 4: Performance of models across varying trajectory length. Vanilla refers to the Llama3-8B model, w/ Critique refers to the critique-guided inference with our trained Critic-Llama3, and #Iter-k represents the k-th iteration of action refinement.

ment on hard tasks, with a +28.75% increase after three iterations. These results suggest that models

without action refinement can quickly adapt to critiques for easy and medium-level tasks. However, for longer tasks, the model's ability to leverage critiques diminishes. Action refinement helps the model continuously follow critiques, thereby enhancing performance on long-horizon tasks.

7.2 Effects of Number of Candidate Actions

To further evaluate the scalability of CGI, we varied the number of candidate actions (M) sampled at each step during inference. As shown in Figure 5, we observe the following trends: CGI achieves a significant improvement with just one candidate action, reaching an average score of 56.89%. As M increases, performance improves further, rising from 56.89% to 61.72%. However, when M=7, performance plateaus, indicating saturation. Other numerical- and verbal-based methods exhibit similar trends. For GPT-4 as a critic, performance increases with M, improving from 32.28% to 38.13%. DGAP is particularly sensitive to the number of candidate actions, with performance rising from 15.02% to 25.07%. This highlights that numerical methods are highly dependent on the quality of the sampled candidate actions, as they only discriminate but not sue

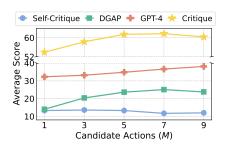


Figure 5: Performance of different methods across varying candidate actions. The average score is the mean final score across three environments.

candidate actions, as they only discriminate but not suggest meaningful revisions. In contrast, self-critique does not benefit from an increase in M. This suggests that low-quality critiques can degrade performance, regardless of the number of candidate actions.

8 Ablation Study

To analyze the contribution of the three types of data (i.e., $\mathcal{D}_{correct}$, \mathcal{D}_{refine} , and $\mathcal{D}_{general}$) used to enhance the actor model during the action refinement stage, we conduct an ablation study. Specifically, we remove each dataset from the training data used in the first iteration of action refinement and evaluate the final performance under the guidance of the critic model. As shown in Figure 6, removing any one of the three datasets results in a performance drop compared to the actor model trained with the full dataset (CGI #Iter1). The most significant decrease occurs when the critique-action pairs are removed (w/o $\mathcal{D}_{\text{refine}}$), where the average performance across the three tasks drops from 68.50% to 50.37%. This decline is more pronounced for longer tasks, with performance reductions of 5.4% on WebShop, 22% on TextCraft, and 26.94% on ScienceWorld. These results highlight that enhancing the actor's critique utilization ability is the most critical factor in improving overall performance. A less severe drop is observed when the generalization

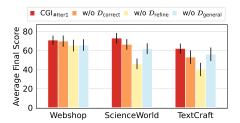


Figure 6: Ablation study of the three data types used in action refinement: $\mathcal{D}_{correct}$ (correct trajectories to improve the actor's reasoning ability), \mathcal{D}_{refine} (critique-action pairs to enhance critique utilization), and $\mathcal{D}_{general}$ (data to maintain the model's generalization capability). CGI $_{\#Iter1}$ refers to the actor model trained in the first iteration of action refinement using the full dataset.

data set is excluded (w/o $\mathcal{D}_{general}$), highlighting the importance of general instructions for model generalization, which helps the model adapt to the unseen test set.

9 Conclusion

We have introduced Critique-Guided Improvement (CGI), a two-player framework that emphasizes nature language feedback for iterative refinement of LLM-based agents. By separating the roles of an actor, which proposes actions, and a critic, which provides verbal guidance, CGI circumvents the limitations of purely numerical signals and addresses the challenges arising from self-refinement. Experimental results in three interactive and agentic environments confirm the effectiveness of CGI, with a small critic model outperforming GPT-4 in providing feedback.

Acknowledgement

We appreciate the support from the Chinese NSF General Program (No.62572129), Major Research Plan (No.92270121). We also acknowledge the use of an icon from Flaticon⁵ and thank its creators for providing this visually appealing design.

References

- [1] OpenAI. GPT-4 Technical Report, 2024.
- [2] Gemini Team. Gemini: A family of highly capable multimodal models. Technical report, Google, 12 2023.
- [3] Meta. Llama 3 model card. 2024.
- [4] Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, Jae Sung Park, Bidipta Sarkar, Rohan Taori, Yusuke Noda, Demetri Terzopoulos, Yejin Choi, Katsushi Ikeuchi, Hoi Vo, Li Fei-Fei, and Jianfeng Gao. Agent AI: Surveying the Horizons of Multimodal Interaction, 2024.
- [5] Jiangjie Chen, Xintao Wang, Rui Xu, Siyu Yuan, Yikai Zhang, Wei Shi, Jian Xie, Shuang Li, Ruihan Yang, Tinghui Zhu, Aili Chen, Nianqi Li, Lida Chen, Caiyu Hu, Siye Wu, Scott Ren, Ziquan Fu, and Yanghua Xiao. From Persona to Personalization: A Survey on Role-Playing Language Agents, 2024.
- [6] Boyu Li, Haobin Jiang, Ziluo Ding, Xinrun Xu, Haoran Li, Dongbin Zhao, and Zongqing Lu. Selu: Self-learning embodied mllms in unknown environments. arXiv preprint arXiv:2410.03303, 2024.
- [7] Siyu Yuan, Kaitao Song, Jiangjie Chen, Xu Tan, Yongliang Shen, Ren Kan, Dongsheng Li, and Deqing Yang. EASYTOOL: Enhancing LLM-based Agents with Concise Tool Instruction, 2024.
- [8] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. ToolLLM: Facilitating Large Language Models to Master 16000+ Real-world APIs, 2023.
- [9] Yue Hu, Yuzhu Cai, Yaxin Du, Xinyu Zhu, Xiangrui Liu, Zijie Yu, Yuchen Hou, Shuo Tang, and Siheng Chen. Self-Evolving Multi-Agent Collaboration Networks for Software Development. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [10] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In The Twelfth International Conference on Learning Representations, 2024.
- [11] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. WebShop: Towards Scalable Real-World Web Interaction with Grounded Language Agents, 2023.
- [12] Pranav Putta, Edmund Mills, Naman Garg, Sumeet Motwani, Chelsea Finn, Divyansh Garg, and Rafael Rafailov. Agent Q: Advanced reasoning and learning for autonomous ai agents, 2024.
- [13] Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce LLMs step-by-step without human annotations. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9426–9439, Bangkok, Thailand, August 2024. Association for Computational Linguistics.

⁵https://www.flaticon.com

- [14] Haofu Qian, Chenjia Bai, Jiatao Zhang, Fei Wu, Wei Song, and Xuelong Li. Discriminator-Guided Embodied Planning for LLM Agent. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [15] Yu Xia, Jingru Fan, Weize Chen, Siyu Yan, Xin Cong, Zhong Zhang, Yaxi Lu, Yankai Lin, Zhiyuan Liu, and Maosong Sun. AgentRM: Enhancing agent generalization with reward modeling, 2025.
- [16] Zhenfang Chen, Delin Chen, Rui Sun, Wenjun Liu, and Chuang Gan. Autonomous agents from automatic reward modeling and planning. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [17] Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback, 2022.
- [18] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, Advances in Neural Information Processing Systems, 2022.
- [19] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- [20] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language Agents with Verbal Reinforcement Learning, 2023.
- [21] Ruihan Yang, Jiangjie Chen, Yikai Zhang, Siyu Yuan, Aili Chen, Kyle Richardson, Yanghua Xiao, and Deqing Yang. SELFGOAL: Your Language Agents Already Know How to Achieve High-level Goals. In NeurIPS 2024 Workshop on Open-World Agents, 2024.
- [22] Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su, John D Co-Reyes, Avi Singh, Kate Baumli, Shariq Iqbal, Colton Bishop, Rebecca Roelofs, Lei M Zhang, Kay McKinney, Disha Shrivastava, Cosmin Paduraru, George Tucker, Doina Precup, Feryal Behbahani, and Aleksandra Faust. Training Language Models to Self-Correct via Reinforcement Learning. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [23] Alejandro Cuadron, Dacheng Li, Wenjie Ma, Xingyao Wang, Yichuan Wang, Siyuan Zhuang, Shu Liu, Luis Gaspar Schroeder, Tian Xia, Huanzhi Mao, Nicholas Thumiger, Aditya Desai, Ion Stoica, Ana Klimovic, Graham Neubig, and Joseph E. Gonzalez. The danger of overthinking: Examining the reasoning-action dilemma in agentic tasks, 2025.
- [24] Yannick Metz, Andras Geiszl, Raphaël Baur, and Mennatallah El-Assady. Reward Learning from Multiple Feedback Types. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [25] Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. ScienceWorld: Is your Agent Smarter than a 5th Grader?, 2022.
- [26] Archiki Prasad, Alexander Koller, Mareike Hartmann, Peter Clark, Ashish Sabharwal, Mohit Bansal, and Tushar Khot. ADaPT: As-Needed Decomposition and Planning with Language Models, 2024.
- [27] Binghai Wang, Rui Zheng, Lu Chen, Zhiheng Xi, Wei Shen, Yuhao Zhou, Dong Yan, Tao Gui, Qi Zhang, and Xuanjing Huang. Reward modeling requires automatic adjustment based on data quality. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4041–4064, Miami, Florida, USA, November 2024. Association for Computational Linguistics.

- [28] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training Verifiers to Solve Math Word Problems, 2021.
- [29] Fei Yu, Anningzhe Gao, and Benyou Wang. OVM, outcome-supervised value models for planning in mathematical reasoning. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, Findings of the Association for Computational Linguistics: NAACL 2024, pages 858–875, Mexico City, Mexico, June 2024. Association for Computational Linguistics.
- [30] Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Meiqi Guo, Harsh Lara, Yunxuan Li, Lei Shu, Yun Zhu, Lei Meng, Jiao Sun, and Abhinav Rastogi. Improve Mathematical Reasoning in Language Models by Automated Process Supervision, 2024.
- [31] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. WebGPT: Browser-assisted question-answering with human feedback, 2022.
- [32] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional AI: Harmlessness from AI Feedback, 2022.
- [33] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena, 2023.
- [34] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-Refine: Iterative Refinement with Self-Feedback, 2023.
- [35] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. ReAct: Synergizing Reasoning and Acting in Language Models. In *The Eleventh International Conference on Learning Representations*, 2023.
- [36] Yu Shang, Yu Li, Keyu Zhao, Likai Ma, Jiahe Liu, Fengli Xu, and Yong Li. AgentSquare: Automatic LLM Agent Search in Modular Design Space, 2025.
- [37] Minghao Chen, Yihang Li, Yanting Yang, Shiyu Yu, Binbin Lin, and Xiaofei He. AutoManual: Constructing Instruction Manuals by LLM Agents via Interactive Environmental Learning, 2024.
- [38] Bodhisattwa Prasad Majumder, Bhavana Dalvi Mishra, Peter Jansen, Oyvind Tafjord, Niket Tandon, Li Zhang, Chris Callison-Burch, and Peter Clark. CLIN: A Continually Learning Language Agent for Rapid Task Adaptation and Generalization, 2023.
- [39] Kolby Nottingham, Bodhisattwa Prasad Majumder, Bhavana Dalvi Mishra, Sameer Singh, Peter Clark, and Roy Fox. Skill Set Optimization: Reinforcing Language Model Behavior via Transferable Skills. *arXiv*, 2024.
- [40] Zhiheng Xi, Yiwen Ding, Wenxiang Chen, Boyang Hong, Honglin Guo, Junzhe Wang, Dingwen Yang, Chenyang Liao, Xin Guo, Wei He, Songyang Gao, Lu Chen, Rui Zheng, Yicheng Zou, Tao Gui, Qi Zhang, Xipeng Qiu, Xuanjing Huang, Zuxuan Wu, and Yu-Gang Jiang. AgentGym: Evolving Large Language Model-based Agents across Diverse Environments, 2024.

- [41] Siyu Yuan, Zehui Chen, Zhiheng Xi, Junjie Ye, Zhengyin Du, and Jiecao Chen. Agent-R: Training Language Model Agents to Reflect via Iterative Self-Training, 2025.
- [42] Zehui Chen, Kuikun Liu, Qiuchen Wang, Wenwei Zhang, Jiangning Liu, Dahua Lin, Kai Chen, and Feng Zhao. Agent-FLAN: Designing data and methods of effective agent tuning for large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 9354–9366, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- [43] Yifan Song, Weimin Xiong, Xiutian Zhao, Dawei Zhu, Wenhao Wu, Ke Wang, Cheng Li, Wei Peng, and Sujian Li. AgentBank: Towards generalized LLM agents via fine-tuning on 50000+ interaction trajectories. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, Findings of the Association for Computational Linguistics: EMNLP 2024, pages 2124–2141, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- [44] Wentao Shi, Mengqi Yuan, Junkang Wu, Qifan Wang, and Fuli Feng. Direct multi-turn preference optimization for language agents. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2312–2324, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- [45] Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. Trial and Error: Exploration-Based Trajectory Optimization for LLM Agents, 2024.
- [46] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R Narasimhan. Tree of Thoughts: Deliberate Problem Solving with Large Language Models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [47] Renat Aksitov, Sobhan Miryoosefi, Zonglin Li, Daliang Li, Sheila Babayan, Kavya Kopparapu, Zachary Fisher, Ruiqi Guo, Sushant Prakash, Pranesh Srinivasan, Manzil Zaheer, Felix Yu, and Sanjiv Kumar. ReST meets ReAct: Self-Improvement for Multi-Step Reasoning LLM Agent, 2023.
- [48] Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large language models can self-improve. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1051–1068, Singapore, December 2023. Association for Computational Linguistics.
- [49] Yuxiao Qu, Tianjun Zhang, Naman Garg, and Aviral Kumar. Recursive introspection: Teaching language model agents how to self-improve, 2024.
- [50] Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. STaR: Bootstrapping Reasoning With Reasoning, 2022.
- [51] Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. AgentTuning: Enabling Generalized Agent Abilities for LLMs, 2023.
- [52] Archiki Prasad, Alexander Koller, Mareike Hartmann, Peter Clark, Ashish Sabharwal, Mohit Bansal, and Tushar Khot. ADaPT: As-needed decomposition and planning with language models. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 4226–4252, Mexico City, Mexico, June 2024. Association for Computational Linguistics.
- [53] OpenAI. Chatgpt blog post. https://openai.com/blog/chatgpt, 2022. Accessed: 2024-09-06.
- [54] OpenAI. Gpt-4 technical report, 2023.
- [55] Anthropic. Introducing claude 2.1, Nov 2023. Available from Anthropic: https://www.anthropic.com/news/claude-2-1.

- [56] DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan 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, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Oiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Levi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2025.
- [57] Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. AgentTuning: Enabling generalized agent abilities for LLMs. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 3053–3077, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- [58] Diederik P Kingma. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Abstract and Section 1

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
 contributions made in the paper and important assumptions and limitations. A No or
 NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Appendix A

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper does not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Section 3, Section 4 and Appendix D, E, F

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived
 well by the reviewers: Making the paper reproducible is important, regardless of
 whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: Section 4

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Section 4, Section 5 and Appendix F

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: We set temperature to 0 for all models in experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Appendix E

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Section 5, Appendix E and F

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Appendix D

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to

generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.

- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
 necessary safeguards to allow for controlled use of the model, for example by requiring
 that users adhere to usage guidelines or restrictions to access the model or implementing
 safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
 not require this, but we encourage authors to take this into account and make a best
 faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification:Section 5, Appendix E and F

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: Section 5, Appendix E and F

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing or research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [No]

Justification: The paper does not use LLMs for writing, editing, or formatting. Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

Appendix

A Limitations

CGI incurs higher computational time compared to baseline methods but remains within a reasonable range. As shown in Table 3, CGI requires approximately four times the computational time of the baseline. However, this additional cost results in a substantial performance improvement, with CGI achieving a 43.31% gain over No Critique. This indicates that the extra computational time is effectively utilized. In contrast, other verbal-based methods, such as GPT-40 and Self-Critique, do not produce significant improvements and require more computational time than CGI.

B Broader Impacts

Positive Societal Impacts. By generating and leveraging natural language critiques, CGI enhances the performance in embodied AI tasks (*e.g.*, , ScienceWorld) and real-world scenarios (*e.g.*, , WebShop), leading to improved task performance. This contributes to enabling agents to better understand and align with human goals, ultimately supporting more helpful, interpretable AI systems. Moreover, our critic model demonstrates strong performance despite using a relatively small backbone (Llama-3-8B), which significantly reduces the computational and environmental costs compared to larger proprietary models. This also promotes broader accessibility to advanced feedback mechanisms.

Negative Societal Impacts and Risks. Natural-language critiques may inadvertently incorporate biases from the training data or expose private information through the agent's reasoning process. Moreover, as these critiques provide explicit, step-by-step explanations, they could be misused by bad actors to facilitate disinformation, social engineering, or the generation of unsafe code. The relatively small size of the critic model further reduces the cost of replication and misuse, increasing the associated risks.

C Notations

Dataset Definitions				
Meaning				
Expert critique dataset generated by expert critic (GPT-4o), used for training the critic model as described in §4.2				
Correct trajectory dataset collected during each iteration of the action refinement process				
Critique-action pairs dataset used for learning how to utilize critiques effectively				
General conversational dataset (e.g., ShareGPT) used to maintain general language modeling capabilities during fine-tuning				
Loss Function Parameters				
Meaning				
Actor model loss function with three components: learning from expert/correct trajectories, learning to utilize critiques, and maintaining general capabilities				
Weight parameter balancing agentic task learning vs. general capabilities (typically $\beta=0.8)$				
General text responses from datasets like ShareGPT for maintaining language modeling abilities				

D Instruction Prompt Examples

The system prompts for three agentic environments are presented in Listing 1. The instruction for critique generation introduced in §4.2 is presented in Listing 2.

Listing 1: Prompt details for ScienceWorld, WebShop, and TextCraft.

```
ScienceWorld Instruction:
You are an agent for the science world. Every round I will give you an
observation, you have to respond with an action based on the observation
to finish the given task.
Here are the actions you may take:
{"action": "open/close OBJ", "description": "open/close a container",} {"action": "de/activate OBJ", "description": "activate/deactivate a
device",}
{"action": "connect OBJ to OBJ", "description": "connect electrical
components", }
{"action": "disconnect OBJ", "description": "disconnect electrical
components",}
{"action": "use OBJ [on OBJ]", "description": "use a device/item",}
{"action": "look around", "description": "describe the current room",} {"action": "look at OBJ", "description": "describe an object in detail",} {"action": "look in OBJ", "description": "describe a container's contents
{"action": "read OBJ", "description": "read a note or book",}
{"action": "move OBJ to OBJ", "description": "move an object to a
container", }
{"action": "pick up OBJ", "description": "move an object to the inventory
", }
{"action": "put down OBJ", "description": "drop an inventory item",}
{"action": "pour OBJ into OBJ", "description": "pour a liquid into a
container", }
{"action": "dunk OBJ into OBJ", "description": "dunk a container into a
liquid", }
{"action": "mix OBJ", "description": "chemically mix a container",}
{"action": "go to LOC", "description": "move to a new location",} {"action": "eat OBJ", "description": "eat a food",}
{"action": "flush OBJ", "description": "flush a toilet",}
{"action": "focus on OBJ", "description": "signal intent on a task object
{"action": "wait", "description": "take no action for 10 iterations",}
{"action": "wait1", "description": "take no action for 1 iteration", }
{"action": "task", "description": "describe current task",}
{"action": "inventory", "description": "list your inventory"}
Your response should use the following format:
Thought: your thoughts.
Action: your next action
WebShop Instruction:
You are web shopping. I will give you instructions about what to do. You
have to follow the instructions. Every round I will give you an
observation and a list of available actions, you have to respond an
action based on the state and instruction.
You can use search action if search is available. You can click one of
the buttons in clickables. An action should be of the following structure
: search[keywords] click[value]. If the action is not valid, perform
nothing. Keywords in search are up to you, but the value in click must be
 a value in the list of available actions. Remember that your keywords in
 search should be carefully designed. Your response should use the
following format:
Thought: I think ...
Action: click[something] or search[something]
```

TextCraft Instruction:

You are given a few useful crafting recipes to craft items in Minecraft. Crafting commands are of the format "craft [target object] using [input ingredients]". Every round I will give you an observation, you have to respond to an action based on the state and instruction. You can "get" an object (ingredients) from the inventory or the environment, look up the game "inventory" by inventory, or "craft" (target) using any of the crafting commands. You can use ONLY these crafting commands provided, do not use your own crafting commands. However, if the crafting command uses a generic ingredient like "planks", you can use special types of the same ingredient e.g. dark oak "planks" in the command instead. Your response should use the following format:

Thought: ...

Listing 2: The instruction prompts for Critique Generation. Critique Generation Instruction: Your task is to critique the candidate's next-step action based on the agent's task goal and interaction history. {available_actions} Critique Steps Step 1: Analyze Candidate Action Examine the candidate's action based on the following criteria and assign an overall grade using this scale: Excellent, Good, Neutral, Poor, Very Poor. Critique Dimensions - Contribution: Assess whether the action contributes to completing the agent's task. This includes both direct actions (e.g., picking up the target OBJ) and indirect actions (e.g., reasonable exploration that can provide additional environmental information and facilitate future progress). - Feasibility: Assess whether the action is valid according to the agent' s predefined Allowed Action Types list. - Efficiency: Analyze whether the action optimally achieves the task without unnecessary steps or redundancy. Step 2: Provide Revision Suggestions Suggest a modification to align the candidate's action better with the task or the agent's action capabilities. Note that the suggested revision should be based on the Allowed Action and Object Types. Critique Format Please structure your critique in the following format: ## Contribution: [Analysis of Contribution]. ## Feasibility: [Analysis of feasibility]. ## Efficiency: [Analysis of efficiency]. ## Overall Grading: [Overall grade: Excellent/Good/Neutral/Poor/Very Poor ## Suggested Revision: [Brief revision suggestion, if applicable]. Inputs: The agent's task goal and interaction history: {history} Candidate next step action: {candidate_action} Now, please provide your critique:

```
Critique Generation Instruction (with expert path):
Your task is to critique the candidate next-step action based on the
agent's task goal and interaction history. The gold path for current task
is provided as a reference to guide your critique, but do not explicitly
mention it in your critique.
{available_actions}
Critique Steps
Step 1: Analyze Candidate Action
Examine the candidate's action based on the following criteria, then
assign an overall grade using this scale: Excellent, Good, Neutral, Poor,
Very Poor.
Critique Dimensions
-Contribution: Assess whether the action contributes to completing the
agent's task. This includes both direct actions (e.g., picking up the
target OBJ) and indirect actions (e.g., reasonable exploration that can
provide additional environmental information and facilitate future
progress).
-Feasibility: Assess whether the action is valid according to the agent's
predefined Allowed Action Types list.
-Efficiency: Analyze whether the action optimally achieves the task
without unnecessary steps or redundancy.
Step 2: Provide Revision Suggestions
Suggest a modification to align the candidate's action better with the
task or the agent's action capabilities. For example, if the action is
not allowed, recommend an alternative from the action list that aligns
better with the task goal.
Critique Format
Please structure your critique in the following format:
## Contribution: [Analysis of Contribution].
## Feasibility: [Analysis of feasibility].
## Efficiency: [Analysis of efficiency].
## Overall Grading: [Overall grade: Excellent/Good/Neutral/Poor/Very Poor
## Suggested Revision: [Brief revision suggestion, if applicable].
Referenced Gold Path for Current Task:
{gold_path}
The agent's task goal and interaction history:
{history}
Candidate next step action: {candidate_action}
```

Now, please provide your critique:

Table 3: Computational efficiency of different methods in three agentic environments per task.

Method	WebShop		Sciworld		Textcraft	
	Inference Time	Performance	Inference Time	Performance	Inference Time	Performance
No Critique	0.38 min	13.49	1.49 min	14.48	1.02 min	10.00
w/ DGAP	1.13 min	30.41	3.13 min	19.52	1.33 min	21.00
w/ Self-Critique	1.36 min	1.50	3.42min	10.06	2.17 min	19.00
w/ GPT-4o	2.34 min	17.78	6.28 min	33.06	3.28 min	46.00
w/ Critic Model	1.22 min	56.80	3.41 min	68.51	1.98 min	59.00

E Evaluation Settings

In the experiments presented in this paper, we used three agentic environments: WebShop, SciWorld, and TextCraft.

E.1 Details of Environments

WebShop WebShop is an interactive web environment designed for web shopping. In this environment, agents are given instructions and must purchase a product that meets specified criteria. Agents can either click a button on the webpage or use the search engine to find the product. We evaluate performance using the success rate, with a maximum of 10 rounds per task.

ScienceWorld ScienceWorld is a benchmark environment for testing agents' scientific reasoning abilities, based on a standard elementary science curriculum. It includes 30 types of tasks, such as using measurement instruments and conducting mechanics experiments. The action space is task-specific, with the environment simulator providing the effects of actions. We use reward as the evaluation metric, with a maximum of 30 rounds per task. Task lengths are classified based on the ScienceWorld environment, which provides gold-standard trajectories from 30 hand-coded oracle agents. We adopt the following categorization: Short (11.76), Medium (28.58), and Long (94.30). These lengths correspond to the average number of steps taken by the oracle agent.

TextCraft TextCraft is a text-only environment for crafting Minecraft items. It constructs a crafting tree based on Minecraft's crafting recipes, consisting of 544 nodes, each representing a target item. For each task, the agent is given a target item and a list of crafting commands generated by the tree. Tasks vary in complexity, ranging from 1 to 4 steps. The environment supports three valid actions: "craft <item> using <ingredients>", "get <item>", and "inventory". After each round, the environment checks the agent's actions and returns the execution state. Agents are rewarded with 1 point only upon successfully crafting the target item. We evaluate performance using the success rate, with a maximum of 20 rounds per task.

E.2 Computational Efficiency Analysis

We evaluated the computational efficiency of CGI by conducting experiments in three agentic environments, using Llama-3-8B as the actor model. We measured the average inference time per task. As shown in Table 3, although CGI incurred higher computation times compared to inference without critique guidance (i.e., No Critique), these times remained within an acceptable range. Furthermore, CGI demonstrated a substantial reduction in inference time compared to using GPT-4 as the critic, while achieving significant performance improvements.

F Implementation Details

In our experiments, we use the Llama-3-8B-Instruct [3] as the base model for both the actor and the critic

F.1 Training Data

Critique Generation The training data for the critic model consists of expert critiques generated by the expert critic (*i.e.*, GPT-4o) in the ScienceWorld, Webshop, and TextCraft environments, as described in § 4.2. The specific training set sizes during the SFT phases for Llama3-8B-Instruct are 14K from ScienceWorld, 10K from Webshop, and 8K from TextCraft. Examples of our SFT training data are provided below:

```
SFT Example:
{"conversations":
[{"from": "human", "value": "Your task is to critique the candidate's
next-step action based on the agent's task goal and interaction history.
Allowed Action Types:- search[KEYWORDS]: perform a web search - click[OBJ]: click an object or link on a webpage. Step 1: Analyze Candidate Action
.Examine the candidate's action based on the following criteria and
assign an overall grade using this scale: Excellent, Good, Neutral, Poor,
Very Poor. Critique Dimensions: Contribution: Assess whether the action
contributes to completing the agent's task. This includes both direct
actions and indirect actions. Feasibility: Assess whether the action is
valid according to the agent's predefined action list. Specifically,
check if the action is part of the Allowed Action Types, whether the OBJ/
KEYWORDS within the action is clearly specified, and whether it is
formatted correctly according to the system's requirements. Efficiency:
Analyze whether the action optimally achieves the task without
unnecessary steps or redundancy. Step 2: Provide Revision Suggestionn:
Suggest a modification to align the candidate's action better with the
task or the agent's action capabilities. For example, if the action is
not allowed, recommend an alternative from the action list that aligns
better with the task goal. Critique Format:Please structure your critique
 in the following format : \#\# Contribution: [Analysis of Contribution].
## Feasibility: [Analysis of feasibility]. ## Efficiency: [Analysis of
efficiency].## Overall Grading: [Overall grade: Excellent/Good/Neutral/
Poor/Very Poor]. ## Suggested Revision: [Brief revision suggestion, if
applicable]. ### InputsThe Agent's Task Goal and Interaction History:
WebShop [SEP] Instruction: [SEP] Find me machine wash, moisture wicking
men's t-shirts & tanks with polyester spandex [SEP] Search Candidate Next
Step Action: search[machine wash moisture wicking men's t-shirts with
polyester spandex long sleeve red blue XXL price: under 50]. Now, please
provide your critique:"},
{"from": "gpt", "value": "## Contribution: The candidate's action
contributes to the task by attempting to perform a search for the
specified product. It includes most of the necessary attributes such as
material, sleeve length, color, size, and price, which are crucial for
finding the desired item. ## Feasibility: The action is feasible as it
uses the allowed action type "search[KEYWORDS]". However, there are some
discrepancies in the keywords used. The candidate action uses "XXL"
instead of "xx-large" and "price: under 50" instead of "price < 50.00".
While these variations might still yield relevant results, they do not
exactly match the specified format. ## Efficiency: The action is
relatively efficient as it attempts to gather relevant search results in
one step. However, the slight variations in keyword formatting could
potentially lead to less precise search results, which might require
additional steps to refine the search. ## Overall Grading: Good. ##
Suggested Revision: Modify the search keywords to exactly match the
specified format: search[machine wash moisture wicking men's t-shirts &
tanks polyester spandex long sleeve red blue size xx-large price <
50.00]. This ensures consistency and potentially more accurate search
results."}],
"system": ""}}
```

Table 4: Training dataset size for each iteration of Action Refinement stage.

Env	Gold		#Iter 1		#Iter 2		#Iter 3	
	Path	Critique	Path	Critique	Path	Critique	Path	Critique
ScienceWorld	378	4106	133	1464	158	1803	193	2205
WebShop	427	1659	106	663	283	1637	344	2010
TextCract	441	1701	233	1423	310	2070	340	2049
OverAll	1246	6740	472	3550	751	5510	877	6264

Action Refinement As described in § 4.3, we performed action refinement three times. Each training dataset consists of two parts: Correct Path, which includes the correct trajectories obtained by the actor from interactions with the environment in the previous iteration; and Correct Critique, which includes the <critique, action> pairs that the actor correctly followed during the previous iteration. The dataset size for each iteration of Llama-3-8B-Instruct is detailed in Table 4.

F.2 Finetuning Details

We ran SFT experiments using 8 NVIDIA A100-40GB GPUs. For action refinement, to prevent policy drift, we only trained the base model (i.e., Llama-3-8B-Instruct) each time. We conduct experiments with the LlamaFactory code base⁶. The configurations of our hyper-parameters are detailed in Table 5.

Table 5: Fine-tuning hyper-parameters for Critique Generation and Action Reinement stage.

Configuration	Critique Generation	Action Refinement		
Model	Llama-3-8B-Instruct	Llama-3-8B-Instruct		
Number of epochs	3	3		
Devices	8 A100 GPU (40 GB)	8 A100 GPU (40 GB)		
Total Batch size	64 samples	64 samples		
Optimizer	Adam [58]	Adam [58]		
	$(\beta_1 = 0.9, \beta_2 = 0.98, \epsilon = 1 \times 10^{-8})$	$(\beta_1 = 0.9, \beta_2 = 0.98, \epsilon = 1 \times 10^{-8})$		
Learning rate	2×10^{-5}	2×10^{-5}		
Warmup Ratio	0.05	0.05		
Cutoff Length	4096	4096		
Training Time	5h 46m 24s	11h 11m 55s		

F.3 Effect of Iteration Count

Here we analyze the impact of the number of iterations on the performance of CGI. As shown in Figure 7, the fourth iteration leads to a performance drop in ScienceWorld and TextCraft, while the improvement in WebShop is marginal (+0.67%). The third iteration achieves the highest average performance across all tasks. These results suggest that additional iterations do not yield further gains. Therefore, we report the results from the third iteration in the main results section (§ 6).

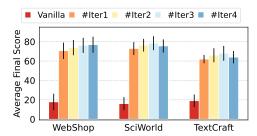


Figure 7: Performance of models with varying iteration counts across three tasks.

F.4 Effect of Training with Expert Trajectories and Updating the Critic

In this section, we analyze two key factors: 1) whether our critic heavily depends on expert trajectories; 2) whether updating the critic model through iterations can continue to improve performance. To evaluate the first factor, we replace the gold-standard expert trajectories with GPT-40-generated

⁶https://github.com/hiyouga/LLaMA-Factory

Table 6: Performance comparison across different methods and datasets

Method	WebShop	ScienceWorld	TextCraft	Average
Llama3-8B	13.49	14.48	10.00	12.66
w/ GPT-4o	17.78	33.06	46.00	32.28
w/ Critic (Trained on GPT-4o Traj)	43.21	56.93	52.00	50.71
w/ Critic (Trained on Expert Traj)	56.80	68.51	59.00	61.44
w/ Updated Critic	62.51	74.56	64.00	67.02

trajectories across three environments, following the exact pipeline described in §4.2 for training the critic model. For the second factor, we perform iterative updates on the critic model after action refinement to adapt to the updated actor's policy. Specifically, we use the critic from the first iteration to generate new critique data and retrain the critic.

As shown in Table 6, CGI remains highly effective even without gold-standard expert data. Our method achieves an average performance of 50.71, significantly outperforming both the baseline and the GPT-40 model that generated the initial data. This highlights the robustness of our framework. Additionally, the critic can be dynamically updated. The results indicate that updating the critic leads to further performance improvements (from 61.44% to 67.02%), demonstrating that the critic is not limited to being static and can evolve alongside the actor.

G Baselines

G.1 Critique Methods

In Section 6, we compare our trained Critic-Llama3 model with other verbal- and numerical-based methods. For the verbal-based approaches, such as self-critique and GPT-4, we generate critiques by prompting the model with the instructions provided in Appendix D. For the numerical-based method (i.e., DGAP), following [14], we compute the cosine similarity between the actor model (Llama-3-8B) and expert data at each step. We collected 23K data points for Scienceworld, 16K for WebShop, and 12K for TextCraft. Examples of the reward model training data are provided below:

```
DGAP Example:
{"input": "Your task is to measure the melting point of lead, which is
located around the kitchen. 10. look around. Action: go to hallway",
Score": "10"}
{"input": "Your task is to measure the melting point of lead, which is
located around the kitchen. 10. look around. Action: look at art studio",
 "Score": "0"}
{"input": "Your task is to measure the melting point of lead, which is
located around the kitchen. 10. look around. 9. go to hallway. Action: put down orange", "Score": "0"}
{"input": "Your task is to measure the melting point of lead, which is
located around the kitchen. 10. look around. Action: look at hallway",
Score": "9.03"}
{"input": "Your task is to measure the melting point of lead, which is
located around the kitchen. 10. look around. 9. look at hallway. Action:
open door to outside", "Score": "6.13"}
{"input": "Your task is to measure the melting point of lead, which is
located around the kitchen. 10. look around. 9. look at hallway. 8. open
door to outside. Action: teleport to kitchen", "Score": "8.87"}
```

G.2 Iterative Methods

In Section 6, we compare CGI with two iterative baselines: Reflexion and Iterative SFT. Following [20], we prompt the model for self-reflection (Listing 3) at the end of each iteration. The insights gained from the current round are then incorporated into the system prompt to guide the model's decision-making in the subsequent round. For Iterative SFT, we use the actor model π_{θ}^k from the previous round to interact with the environment. The newly collected correct trajectories are combined with the expert path to form the training data for the next round. Specifically, for iterations 1 to 3, we

used 1,676, 1,719, and 1,750 trajectories, respectively, mixing them with the general dataset (*i.e.*, ShareGPT) at a 1:4 ratio for training. To prevent policy drift, we update π^0_θ at each iteration. Each iteration consists of 3 epochs. The learning rate for Iterative SFT is set to 2e-5, with a 3% warm-up and a cosine scheduler. We use the AdamW optimizer with a weight decay of 0 and apply a maximum gradient clipping norm of 1 to prevent gradient explosion. The batch size is set to 2 per device, with a sequence length of 4,096 and an accumulation count of 4, adjusted based on sequence parallelism.

Listing 3: The instruction prompt for Reflexion.

```
Reflexion Instruction:
Your task is to provide a next plan of action based on the previous attempt's interaction history. The next plan should reflect the lessons learned from the failure and adjust for better results.

Previous Trial:
{history}

Next Plan:
```

H Case Study

H.1 Failure Cases of Fine-tuned Llama3

We provide negative examples where the fine-tuned Llama3-8B does not follow the critique.

Listing 4: Failure Cases of Fine-Tuned Llama3-8B in ScienceWorld

```
Task Description:
Your task is to measure the temperature of unknown substance B, which is
located around the bathroom. First, focus on the thermometer. Next, focus
on the unknown substance B. If the unknown substance B temperature is
above 100.0 degrees celsius, place it in the orange box. If the unknown
substance B temperature is below 100.0 degrees celsius, place it in the
yellow box. The boxes are located around the bathroom.
History:
Action: open door to kitchen
Observation: The door is now open.
Action: go to kitchen
Observation: You move to the kitchen.
Action: open door to hallway
Observation: The door is now open
Action: go to hallway
Observation: You move to the hallway.
Critique:
. . . . . .
### Feasibility: The candidate action "open door to bathroom" is
efficient in terms of task progression, as it directly addresses the need
to access the bathroom. However, the agent should consider checking the
current state first to ensure the action is feasible by "Look Around".
### Overall Grading: Netual
Action: go to bathroom
Observation: No known action matches that input.
```

H.2 Critic-Llama3 Outperforms GPT-40

We provide the following examples to explain why our critic outperforms GPT-40 when the actor model is fixed as Llama3-8B. Compared to GPT-40, our critic demonstrates superior spatial reasoning and a higher degree of goal-directed behavior.

Listing 5: Comparison between Critic-Llama3 and GPT-40 as Critics

```
Example 1: Spatial Reasoning
```

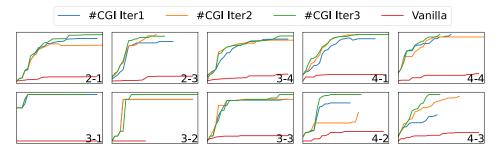


Figure 8: Visualizing trajectories of CGI, and base model (Llama-3-8B-Instructt), the x-axis represents time steps, ranging from 0 to T, while the y-axis denotes scores, which vary from 0 to 100. Each graph illustrates the trajectories corresponding to different tasks in test variation. Task identifiers are positioned at the bottom right of each figure.

```
State: Agent is in the workshop. Task requires an object in the bathroom.
Action: focus on thermometer
GPT-40 (Incorrect): Overall Grading: Excellent.
It approves the action based on the task description, ignoring the agent'
s location.
Our Critic (Correct): Overall Grading: Very Poor.
It correctly identifies the thermometer is not in the current room and
suggests navigating first.
Example 2: Goal-Directedness
Task: Find a living thing.
Action: move to the closet (an exploratory but inefficient action).
GPT-40 (Inefficient): Overall Grading: Neutral.
It allows this exploratory step.
Our Critic (Efficient): Overall Grading: Poor.
It rejects the inefficient action and suggests a more direct one, like
opening the door to explore other rooms.
```

H.3 Trajectory Visualization of ScienceWorld

We visualize the cumulative scores of the actor performing different numbers of action refinements (#Iter k) and the base model in the ScienceWorld scene⁷. As shown in Figure 8, our CGI method shows a notable improvement in efficiency, achieving higher scores in fewer steps. Furthermore, as the number of action refinements increases, efficiency also improves.

⁷Detailed information of each task can be found in https://github.com/allenai/ScienceWorld