# SELU: SELF-LEARNING EMBODIED MLLMS IN UN KNOWN ENVIRONMENTS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Recently, multimodal large language models (MLLMs) have demonstrated strong visual understanding and decision-making capabilities, enabling the exploration of autonomously improving MLLMs in unknown environments. However, external feedback like human or environmental feedback is not always available. To address this challenge, existing methods primarily focus on enhancing the decision-making capabilities of MLLMs through voting and scoring mechanisms, while little effort has been paid to improving the environmental comprehension of MLLMs in unknown environments. To fully unleash the self-learning potential of MLLMs, we propose a novel actor-critic self-learning paradigm, dubbed SELU, inspired by the actor-critic paradigm in reinforcement learning. The critic employs self-asking and hindsight relabeling to extract knowledge from interaction trajectories collected by the actor, thereby augmenting its environmental comprehension. Simultaneously, the actor is improved by the self-feedback provided by the critic, enhancing its decision-making. We evaluate our method in the AI2-THOR and VirtualHome environments, and SELU achieves critic improvements of approximately 28% and 30%, and actor improvements of about 20% and 24% via self-learning.

026 027 028

029

025

004

010 011

012

013

014

015

016

017

018

019

021

## 1 INTRODUCTION

Multimodal Large Language Models (MLLMs) have demonstrated impressive perceptual and understanding capabilities across various domains, *e.g.*, web applications (Ma et al., 2024; Tao et al., 2024; Liu et al., 2024), robotics (Xiong et al., 2024; Li et al., 2024b), gaming (Li et al., 2024c; Qi et al., 2024; Xu et al., 2024), and autonomous driving (Wen et al., 2024; Zhang et al., 2024). Thanks to their powerful capabilities, many works, *e.g.*, Jarvis-1 (Wang et al., 2023b), STEVE-1 (Lifshitz et al., 2023), and Cradle (Tan et al., 2024b), directly utilize the pre-trained MLLMs to complete various decision-making tasks in different embodied environments.

However, the generalization ability of existing pre-trained MLLMs cannot meet the needs of all environments. For some uncommon environments, embodied MLLMs often exhibit hallucinations and poor visual understanding (Huang et al., 2024; Jiang et al., 2024). In more detail, they cannot 040 distinguish left from right and fail to recognize where objects are (Tan et al., 2024b). The reason is 041 that MLLMs have not been further grounded with the environments (Su et al., 2022; Sun et al., 2024). 042 Grounding can be realized by fine-tuning on the experiences from interacting with the environments. 043 Based on the evaluation methods, experience can be categorized into three types: human feedback 044 (Dai et al., 2024; Kirk et al., 2024), environmental feedback (Tan et al., 2024a; Wang et al., 2024b), and self-feedback (Pang et al., 2024; Madaan et al., 2023). The first two types require additional efforts as illustrated in Figure 1(a). Human feedback necessitates expert annotations, which can be 046 costly and influenced by preferences (McAleese et al., 2024). Environmental feedback assumes we 047 can obtain a dynamics model of the environment and build the reward model (Sontakke et al., 2023; 048 Urcelay et al., 2024). Unfortunately, many environments including the real world do not meet such requirements, so these grounding methods are not general. Therefore, we are committed to finding a *general way* to fill in the remaining gaps. 051

Assuming we are facing an unknown embodied environment and cannot obtain external feedback,
 we can only rely on the capabilities of the MLLM itself. Some work, see Figure 1(b), utilizes the evaluation (discriminative) ability of the pre-trained model itself (Wang et al., 2023a; Huang et al.,



Figure 1: Comparison of our framework with other frameworks in terms of the feedback type.

064 2023) to evaluate its own decision-making, and uses such self-feedback to enhance the model's 065 decision-making (generative) capabilities. However, it is easy to see this kind of self-learning is 066 constrained by the static evaluation ability of the model to improve decision-making capability. Ide-067 ally, we hope that self-learning can work similarly to the actor-critic paradigm (Konda & Tsitsiklis, 068 1999) in reinforcement learning, where the actor and critic can mutually enhance each other's per-069 formance. If so, the potential for actor enhancement will be greatly expanded. Unlike reinforcement learning, which can obtain external rewards for training, we do not assume that we can get any external feedback. Therefore, we aim to develop a new actor-critic self-learning paradigm for 071 embodied MLLMs in unknown environments. 072

073 In this paper, we introduce a novel SElf-Learning paradigm in Unknown environments, dubbed 074 **SELU**, as illustrated in Figure 1(c). Inspired by the actor-critic paradigm in reinforcement learning, 075 our paradigm learns to simultaneously optimize the MLLM's ability to understand the environment 076 and to make decisions. For the actor module, we fine-tune the model based on the self-feedback 077 from the critic. As the actor gets improved, it can roll-out more successful trajectories to fine-tune the critic. However, without environmental feedback, the critic may provide inaccurate feedback at the beginning of the training phase, which might mislead the overall optimization. Therefore, 079 we adopt self-asking to correct self-feedback and leverage hindsight relabeling to increase sample efficiency by turning the failure trajectory into a successful one. These high-quality and diverse 081 trajectories are deemed to enhance the critic's comprehension of the environment. Ultimately, the 082 coupling of these two components mutually promotes the improvement of each other, unleashing 083 the full self-learning potential of MLLMs. 084

- Our key contributions can be summarized as follows:
  - We propose a self-learning paradigm for embodied MLLMs, SELU, inspired by the actorcritic paradigm in reinforcement learning, which enables MLLMs to self-adapt to unknown environments.
  - We leverage self-asking and hindsight relabeling to facilitate the improvement of the critic, which greatly increase the sample efficiency and make the self-learning possible.
  - We demonstrate the effectiveness of SELU in the AI2-THOR and VirtualHome environments, achieving critic improvements of approximately 28% and 30%, and actor improvements of about 20% and 24%, respectively.
  - 2 RELATED WORK

062

063

087

090

092

093

094

096 097

098 2.1 MLLMS WITH EXTERNAL FEEDBACK

099 In recent years, MLLMs has achieved impressive results across various visual benchmarks (Mathew 100 et al., 2021; Tang et al., 2024; Lu et al., 2024), demonstrating remarkable perception and decision-101 making capabilities. However, these models still exhibit flaws and often generate unexpected out-102 puts, such as perceptual hallucinations and unreasonable decisions (Yu et al., 2024; Chen et al., 103 2024). Inspired by reinforcement learning, current approaches use external feedback to correct 104 MLLM's erroneous outputs through a cycle of interaction, feedback, and correction (Pan et al., 105 2024; Gero et al., 2023). Generally, there are two sources of external feedback: human preference feedback and environmental feedback. Utilizing manually annotated data to align MLLMs output 106 with expert preferences (Kirk et al., 2024; Zhong et al., 2024), such as DPO (Rafailov et al., 2024), 107 PRO (Song et al., 2024), etc., has proven to be effective. Environmental feedback is typically derived

from designing a reward function (Pang et al., 2024; Tan et al., 2024a) or pre-training an evaluation
model (Schick et al., 2023), which often require substantial support from expert data. For instance,
McAleese et al. (2024) first trained a critic model using expert data to score program code bugs, and
then applied these scores to train ChatGPT with PPO (Schulman et al., 2017), enhancing its debugging capabilities. Compared to existing studies, we focus on the self-learning potential of MLLMs,
as external feedback may be not professional enough in unknown environments.

114

115 2.2 SELF-IMPROVEMENT IN LLMS

Self-improvement in Large Language Models (LLMs) has gained significant attention, as re-117 searchers strive to develop models that can learn and adapt from their own outputs, interactions, and 118 internal feedback mechanisms, without relying on external human-labeled data (Yan et al., 2023; 119 Haluptzok et al., 2023). Early explorations in this area are based on unsupervised learning tech-120 niques, where models learn representations from vast datasets without explicit human guidance 121 (Winter et al., 2022; Zhao et al., 2019). Expanded to LLMs, self-improvement goes further by en-122 abling models to critique, refine, and adapt their behavior in a more autonomous manner (Tan et al., 123 2023; Choi et al., 2024). There are two common self-improvement methods: prompt engineering 124 and fine-tuning. The former is an efficient and intuitive approach for large-scale LLMs, as it allows 125 for the establishment of various chains of thought (CoTs) (Wei et al., 2022) to address the same 126 problem (Huang et al., 2023; Feng et al., 2023). For instance, Madaan et al. (2023) demonstrated that GPT-3.5 and GPT-4 can enhance the rationality of responses by simultaneously inputting a 127 question and reflecting on previous answers. The latter one is more useful for small-scale LLMs, 128 as the prompt engineering is unstable. Wang et al. (2024a) developed a fine-tuning dataset by gen-129 erating negative responses to optimize the LLMs, thereby reducing the occurrence of unreasonable 130 answers. Based on these studies, we choose to use fine-tuning to optimize small-scale MLLMs. 131 However, existing methods overlook the enhancement of MLLM's environment understanding ca-132 pabilities. Therefore, we employ an actor-critic framework to facilitate comprehensive self-learning 133 in MLLMs, optimizing both perception and decision-making abilities.

134 135

136 137

138

# 3 PRELIMINARIES

# 3.1 ACTOR-CRITIC IN REINFORCEMENT LEARNING

Actor-critic (Konda & Tsitsiklis, 1999) is a widely adopted framework in reinforcement learning.
 The agent consists of two learning modules: an actor and a critic, which are optimized iteratively.
 The actor selects and executes actions based on current observations. The critic evaluates these
 observations (and actions) by estimating their values based on reward signals received from the
 environment, thereby guiding the actor to make improved choices in future. This framework takes
 advantage of both policy-based learning and value-based learning and is popular nowadays like PPO
 (Schulman et al., 2017).

146

# 147 3.2 ACTOR-CRITIC FOR MLLMS148

With the development of MLLMs, the feedback provided to the agent is no longer constrained to 149 scalar values, like rewards; it can now include diverse modalities, such as natural language (Dong 150 et al., 2024). This enables the critic gain more specific and informative feedback on the outputs of 151 the actor. Consequently, it can provide more accurate guidance for actor improvement. A prevalent 152 approach in this domain is incorporating human feedback into the critic and building a static evalua-153 tion module that can reflect human preference (Ouyang et al., 2022; Kirk et al., 2024). Specifically, 154 human annotated data is used to train a critic model (McAleese et al., 2024) or a reward model 155 (Sontakke et al., 2023; Wang et al., 2024b) to align the MLLM with human preferences better. 156 More rigorous approaches leverage external evaluation mechanisms, such as tool-interactive learn-157 ing (Gou et al., 2024; Chen et al., 2021), or external knowledge sources like Wikipedia and the 158 Internet (Xu et al., 2023; Li et al., 2024a). However, regardless of whether preference labels or 159 external tools are used, human intervention remains inevitable. To overcome this reliance, methods like self-consistency (Wang et al., 2023a; Schick et al., 2023) employ a voting mechanism to 160 enable the model to evaluate its own behavior without relying on external information. However, 161 self-consistency lacks a learnable critic module, thus it cannot improve its grounding knowledge

166

167

169 170

171 172 173

174 175

182

187 188

189

162



Figure 2: The framework of SELU. (*lower*) The actor MLLM, represented as a robot, collects trajectories for the given instructions. (*upper*) The critic MLLM, denoted as a brain, evaluates these trajectories and determines whether they complete the tasks, guiding the update of the actor MLLM. In addition, the critic MLLM implements self-asking and hindsight relabeling to build a dataset for optimizing itself. The whole framework does not require any external feedback, such as environmental rewards or human annotations.

of the environment it interacts with. In contrast to previous work, we propose a novel actor-critic
based paradigm aimed at achieving self-learning for both the actor MLLM and critic MLLM, enabling them to iteratively improve decision-making and grounding abilities without external human
feedback or environmental rewards.

4 Method

The framework of our method is shown in Figure 2. It consists of two components: the actor MLLM and the critic MLLM. The actor MLLM follows instructions and collects trajectories in the environment. The critic MLLM evaluates the collected trajectories and acquires bootstrapped data via self-asking and hindsight relabeling to optimize itself (Section 4.1). Guided by the success detection results from the critic MLLM, the actor MLLM subsequently improves its decision-making performance in the environment (Section 4.2). By combining the two processes, we can achieve coupled improvements of the critic MLLM and the actor MLLM (Section 4.3).

197

199

4.1 CRITIC: SELF-ASKING AND HINDSIGHT RELABELING

As introduced in Section 1, enhancing the interpretation of environmental grounding information is crucial for an MLLM to improve its performance. In our framework, we achieve this objective via self-asking and hindsight relabeling to acquire bootstrapped data for optimizing the critic MLLM.

Specifically, given an instruction I, the actor MLLM collects a trajectory by following this instruction. The critic takes the last frame  $o_T$  of this trajectory as input, using it as the detection frame to determine whether the task depicted by I is completed,

206 207

$$l_d = M_c(I, p_d, o_T),\tag{1}$$

where  $M_c$  denotes the critic MLLM,  $l_d \in \{"yes", "no"\}$  is the result of the success detection, and  $p_d$  is a prompt for the detection. If the detection result is  $l_d = "yes"$ , we consider this trajectory to be a successful sample for the given instruction and store this trajectory directly into the critic fine-tuning dataset  $\mathcal{D}_{critic}$  in the format  $(I, p_d, o_T, l_d)$ , as shown by trajectory 1 in Figure 2. This trajectory includes environmental grounding information that aligns with the knowledge contained in the critic MLLM.

If the detection result is  $l_d = "no"$ , which means the critic MLLM views this trajectory as a failure for instruction *I*. We first apply **self-asking** to examine the state of task-related objects, as the decision made by the critic MLLM might not be precise due to potential hallucinations. The critic

# MLLM is used to obtain the object states,

$$l'_d = M_c(l_s, I), \quad l_s = M_c(j_I, p_s, o_T)$$
(2)

where  $j_I$  is the object name extracted from instruction I by a text-processing function,  $p_s$  is the prompt for object state analysis,  $l_s$  is the analysis result, and  $M_c$  provides a new success detection  $l'_d$  based on  $l_s$  and I. The format is corrected to  $(I, p_d, o_T, l'_d)$  and will be stored into the critic fine-tuning dataset  $\mathcal{D}_{\text{critic}}$  if  $l'_d = "yes"$ . For example, in trajectory 2 of Figure 2, the critic MLLM initially misjudged the completion of the "open cabinet" task. However, when prompted to focus on the state of the cabinet, it successfully self-corrected its judgment.

If the critic MLLM still considers the trajectory as failure, we propose to use hindsight relabeling 225 to make use of this trajectory, since it might be helpful for learning the environmental grounding of 226 other instructions. Hindsight relabeling is a method that originated in goal-conditioned reinforce-227 ment learning (Andrychowicz et al., 2017). It is based on a simple principle: if a trajectory does 228 not complete the target task, it can be viewed as having accomplished other tasks or subtasks. For 229 example, as shown in trajectory 3 in Figure 2, although it does not complete the task "open cabinet", 230 it successfully completes another task "open drawer". Therefore, we relabel this trajectory with 231 the instruction "open drawer" to help the critic MLLM recognize the completion of the relabeled instruction. We can write this process as, 232

$$I' = M_c(l_h, a_I), \quad l_h = M_c(a_I, p_h, o_T)$$
 (3)

where  $a_I$  is the verb extracted from the instruction I by another text-processing function.  $p_h$  is the prompt for hindsight relabeling, and  $l_h$  is the output, which is usually an object name or None.  $M_c$ checks whether any objects, other than the target object, in the observation  $o_T$  have completed the task associated with  $a_I$ .  $M_c$  generates a new instruction I' if  $l_h$  is not None. After that, we store the data  $(I', p_d, o_T, "yes")$  into the critic fine-tuning dataset  $\mathcal{D}_{critic}$ . Finally, if a failed trajectory proves meaningless after hindsight relabeling, it is considered as not helpful for the MLLM to understand the environment and discarded.

By applying self-asking and hindsight relabeling, we create a fine-tuning dataset  $\mathcal{D}_{critic}$  containing the last frames considered as successful by the critic itself and the last frames relabeled as successful after self-asking and hindsight relabeling.

244 245

246

233

218

#### 4.2 ACTOR: CRITIC-GUIDED IMPROVEMENT

Recent work has shown that the discriminative ability of an LLM exceeds its generative ability (Pang et al., 2024). As MLLMs are typically trained the same way as LLMs, we believe that MLLMs' evaluation abilities would also surpass their generation abilities. For instance, we can easily prompt MLLMs to extract understanding from a given image, but it is challenging to prompt them to choose an appropriate action based on perceived task-relevant information like distance or direction. Our experiment in Section 5.2 also supports this conclusion, where the critic module always performs better than the actor. Therefore, we propose using the critic MLLM to guide the improvement of the actor MLLM in the environment without external feedback.

254 255 Specifically, the actor MLLM interacts with the environment and collects online trajectories. At 256 each timestep t, the actor generates an action plan  $l_{a,t}$  by,

$$l_{a,t} = M_a(I, p_a, o_t), \tag{4}$$

257 where  $M_a$  represents the actor MLLM, I is the task instruction,  $p_a$  is the prompt for action plan and 258  $o_t$  is the current image observation. After collecting a whole trajectory, the critic MLLM determines 259 whether this trajectory completes the instruction I, as described in Section 4.1. If the answer is 260 yes, we will put the whole trajectory into the actor fine-tuning dataset  $\mathcal{D}_{actor}$  with a format of 261  $\{(I, p_a, o_t, l_{a,t})\}_{t=0}^T$ . The relabeled successful trajectories after hindsight relabeling are also added 262 into the actor fine-tuning dataset  $\mathcal{D}_{actor}$ . Since this dataset only contains task completion trajectories, the actor can quickly converge towards completing tasks in the current environment by fine-tuning 264 on the dataset. Note that the actor fine-tuning dataset  $\mathcal{D}_{actor}$  consists of trajectory data, while the 265 critic fine-tuning dataset  $\mathcal{D}_{\rm critic}$  only contains the last frames of these trajectories.

266 267

268

4.3 ACTOR-CRITIC COUPLING IMPROVEMENT

We employ Supervised Fine-Tuning (SFT) (Devlin et al., 2019; Brown et al., 2020) and Low-Rank Adaptation (LoRA) (Hu et al., 2022) to update both the actor and critic MLLMs. Initially, the actor

MLLM interacts with the environment to collect online trajectories containing grounding information, as depicted in the lower part of Figure 2. The critic module then evaluates and classifies these trajectories based on the last frame, as shown in the upper part of Figure 2. We select successful trajectories identified by the critic MLLM to create the actor fine-tuning dataset  $\mathcal{D}_{actor}$ . Subsequently, we utilize the last frame to construct the critic fine-tuning dataset  $\mathcal{D}_{critic}$ . The update of the actor and critic can be performed iteratively, and make them both improve step-by-step. A detailed pseudo-code for our algorithm is available in Appendix A.1.

277 278 279

280

281

286

287

289

291

293

295

306

307

308

310

311

312

313

314

315

319

320

5 EXPERIMENTS

5.1 EXPERIMENTAL SETUP

Environments. In order to simulate embodied MLLM interactions in unknown environments, we
select AI2-THOR (Kolve et al., 2022) and VirtualHome (Puig et al., 2018) for our experiments.
Both environments offer open-ended tasks, various interactive objects, and selectable camera perspectives, facilitating data collection for the actor and critic.

- **AI2-THOR** is an interactive simulation environment designed for embodied AI research. The primary tasks require agents to navigate and interact with household objects. It offers highly realistic 3D environments that simulate kitchens, living rooms, and other indoor settings. AI2-THOR provides locobot and robotic arms as agents and enables agents to perform complex actions like picking up, opening, toggling, and so on.
  - VirtualHome is also an embodied simulation platform designed to imitate human activities and tasks in home environments. It offers a rich set of virtual household scenes where agents can perform tasks, such as sitting on a bench in the kitchen or grasping a waterglass in the bedroom. This environment focuses on task completion through multi-step action sequences, making it ideal for testing long-term planning.

296 **Task Selection.** Considering the training costs while demonstrating the feasibility of the framework, 297 we focus on three typical task categories in the AI2-THOR environment: pick up, open, and break. 298 Locobot is selected as the agent, as this work does not consider low-level control of robotic arms. To 299 ensure task diversity and feasibility, we first prompt the MLLM to explore the environment and use 300 the explore objects to initialize the instruction list which serves as the task set. We randomly sample 2-3 objects for each type of task. Considering the training costs, we restrict the maximum step for 301 all tasks to 10. We apply a similar approach in the VirtualHome environment, selecting "female1" 302 as the agent and primarily testing in grab, open, and sit tasks. 303

MLLMs. To demonstrate the generalization capability of our framework, we conduct experiments
 using two MLLMs: LLaVA (Liu et al., 2023) and Qwen-VL (Bai et al., 2023).

- LLaVA integrates a visual encoder with a language model to effectively process and respond to multimodal inputs. It has gained prominence as one of the most popular MLLMs due to its simple architecture and lower training data requirements. These features enable LLaVA to generate responses more swiftly, and suitable for inference to investigate self-learning.
- **Qwen-VL** is the multimodal version of the large model series. Compared to other opensource MLLMs, Qwen-VL is the first model to use a 448x448 resolution image input. Due to its higher resolution, this model exhibits enhanced visual understanding capabilities. We opt for Qwen-VL with the aim of better success detection, thereby facilitating more efficient self-learning of MLLMs.

Baselines. We compare SELU with five methods to investigate the feasibility of self-learning of MLLMs in embodied environments:

- **DG** refers to the results obtained through direct generation from the initial MLLM without any fine-tuning.
- SC (Wang et al., 2023a) represents an optimization method of MLLMs through selfconsistency. Specifically, we employ multiple chains of thought (CoT) to prompt an MLLM to answer the same question, followed by majority voting. In our experiments, we utilize three different CoTs to guide the MLLM, ultimately voting for the most reasonable action.

• LMSI (Huang et al., 2023) is a self-improvement method based on SC. It generates "highconfidence" answers for unlabeled questions to build fine-tuning datasets. This approach enables the LLM to iteratively improve its performance based on the voting mechanism, and we extend this approach to MLLMs.

- Self-Refine (Madaan et al., 2023) involves multiple rounds of self-reflection, followed by self-optimization based on the reflection results. This method focuses on prompt optimization and has been validated for feasibility in large-scale LLMs, such as GPT-4. In our experiments, we reflect 3 rounds to get the final result.
  - SELU-One represents the method of using the same MLLM to simultaneously perform actor and critic tasks, and fine-tuning with a combination of actor and critic datasets. This approach aims to investigate the feasibility of utilizing a single MLLM to meet the requirements of our framework.

#### 5.2 AI2-THOR

LLaVA. We first demonstrate the effectiveness of SELU in the AI2-THOR environment. After on line interactions with the environment and fine-tuning of LLaVA, the critic exhibits an average per formance improvement of approximately 27%, while the actor achieves an improvement of around
 20% compared to the original model. Table 1 and Table 2 present the accuracy of task success
 detection and task success rate respectively.

Table 1: Accuracy of task success detection in the AI2-THOR environment.

Method	Pick up	Open	Break	Avg.
DG	80.67%	36.50%	50.50%	55.89%
SELU-One	68.67%	30.50%	25.50%	41.56%
SELU	94.33%	67.50%	87.50%	83.11%

Table 2: Task success rate in the AI2-THOR environment. SC and Self-Refine use prompt engineering to realize self-learning, whereas LMSI and SELU utilize fine-tuning.

Method	Pick Up	Open	Break	Avg.
DG	68.33%	65.00%	15.50%	49.61%
SC	65.67%	68.50%	17.50%	50.56%
Self-Refine	69.67%	70.50%	14.50%	51.56%
LMSI	75.67%	52.50%	19.50%	49.22%
SELU-One	91.33%	85.50%	27.50%	68.11%
SELU	94.67%	83.50%	30.50%	69.56%

In Table 1, it is evidenced that the unified fine-tuning of the actor and the critic (SELU-One) leads
 to a decline in success detection, even worse than the original critic. Although SELU-One achieves
 a task success rate comparable to that of SELU as shown in Table 2, the compromised critic will
 result in SELU-One incorrectly analyzing the trajectory in subsequent epochs.

Table 2 demonstrates that baselines are not suitable for the self-learning of embodied MLLMs. Both prompt-engineering and fine-tuning baselines struggle to improve the decision-making ability of the MLLM. The reason is that the embodied MLLM cannot give a correct task detection with a lack of environmental understanding. As we can see in Table 1, the initial judgment of the MLLM (DG) on tasks is only about 55%. In this case, merely optimizing the prompt to create multiple CoTs for repeated reflection does not help the MLLM gain task achievement details in an unknown environment. Therefore, neither SC nor Self-Refine can substantially enhance the task success rate. For fine-tuning baselines, relying on statistical voting to validate its own behavior even leads to worse performance. For instance, in the Open task, for LMSI the task detection accuracy of 36.5% causes the actor's performance to drop from 65% to 52.5% after fine-tuning. These results demonstrate the necessity of the critic module in SELU, and optimizing the critic is crucial for enhancing the actor's performance.

Notably, in the Open task, we can observe a low accuracy of success detection for SELU; however, it still helps the actor improve. This highlights the role of hindsight relabeling, which will be discussed in detail in Section 5.4. 

Qwen-VL. In order to prove that SELU can help different MLLMs achieve self-learning, we select Qwen-VL and test it in the AI2-THOR environment under the same setting. The result is shown in Table 3, which indicates that SELU can help Qwen-VL improve the task evaluation capability by about 24% and the decision-making performance by approximately 23%. 

Table 3: Self-learning performance of SELU on Qwen-VL in the AI2-THOR environment.

Task	С	ritic	А	ctor
TUSK	DG-Qwen-VL	SELU-Qwen-VL	DG-Qwen-VL	SELU-Qwen-VL
Pick Up	73.33%	95.67%	57.67%	95.33%
Open	51.00%	81.50%	46.50%	68.00%
Break	63.50%	83.50%	12.50%	21.50%
Avg.	62.61%	86.89%	38.89%	61.61%

#### 5.3 VIRTUALHOME

We then conduct experiments in the VirtualHome environment, which incorporates a greater variety of items and human agents, thereby enriching the experimental environments to demonstrate the effectiveness of our method. The experimental results are presented in Tables 4 and 5. In this en-vironment, SELU enhances LLaVA task evaluation capability by approximately 30% and improves decision-making performance by around 24%, and it also outperforms baselines. As the environ-ment becomes more complex, the lack of environmental understanding causes SC and Self-Refine to negatively impact the decision-making of the original embodied MLLM, as the performance of SC and Self-Refine are even lower than DG, shown in Table 5. 

Table 4: Accuracy of task success detection in the VirtualHome environment.

Method	Grab	Open	Sit	Avg.
DG	52.67%	35.33%	44.50%	44.17%
SELU-One	45.33%	15.67%	48.50%	36.50%
SELU	93.67%	83.33%	47.50%	74.83%
Table 5: Task	success rate	in the Virt	ualHome e	nvironmer
Table 5: Task	success rate	in the Virt	ualHome e	nvironmen Avg.
Table 5: Task	success rate Grab	e in the Virt Open	ualHome e Sit	nvironmer Avg.
Table 5: Task Method DG	success rate Grab 65.00%	e in the Virt Open 83.33%	ualHome e Sit 56.50%	nvironmer Avg. 68.28%
Table 5: Task Method DG SC	success rate Grab 65.00% 52.67%	e in the Virt Open 83.33% 81.67%	ualHome e Sit 56.50% 61.50%	nvironmen Avg. 68.28% 65.28%
Table 5: Task Method DG SC Self-Refine	success rate Grab 65.00% 52.67% 59.67%	e in the Virt Open 83.33% 81.67% 74.33%	ualHome e Sit 56.50% 61.50% 60.50%	Avg. 68.28% 65.28% 64.83%
Table 5: Task Method DG SC Self-Refine LMSI	success rate Grab 65.00% 52.67% 59.67% 35.67%	e in the Virt Open 83.33% 81.67% 74.33% 93.67%	ualHome e Sit 56.50% 61.50% 60.50% 52.50%	Avg. 68.28% 65.28% 64.83% 60.61%
Table 5: Task Method DG SC Self-Refine LMSI SELU-One	success rate Grab 65.00% 52.67% 59.67% 35.67% 83.67%	e in the Virt Open 83.33% 81.67% 74.33% 93.67% <b>98.67%</b>	ualHome e Sit 56.50% 61.50% 60.50% 52.50% 94.50%	Avg. 68.28% 65.28% 64.83% 60.61% 92.28%
Table 5: Task Method DG SC Self-Refine LMSI SELU-One SELU	success rate Grab 65.00% 52.67% 59.67% 35.67% 83.67% 93 33%	e in the Virt Open 83.33% 81.67% 74.33% 93.67% <b>98.67</b> % 97.67%	ualHome e Sit 56.50% 61.50% 60.50% 52.50% 94.50% 93 50%	Avg. 68.28% 65.28% 64.83% 60.61% 92.28% 94.83%

#### 432 5.4 ABLATION STUDY 433

We conduct ablation experiments on the critic module and hindsight relabeling in the AI2-THOR environemnt, and the results are shown in Table 6. SELU w/o HR means that we do not perform hindsight relabeling to reanalyze the trajectory when evaluating the task. SELU w/o critic means that we remove both self-asking and hindsight relabeling and use the data obtained from environment interaction to directly fine-tune the actor. In this case, the evaluation result of the critic is derived from the original MLLM.

Critic (Success Detection Accuracy) Actor (Task Success Rate) Task SELU w/o HR w/o critic SELU w/o HR w/o critic Pick Up 94.33% 94.67% 83.67% 80.67% 67.33% 56.33% 83.50% Open 67.50% 31.50% 36.67% 66.50% 72.50% 87.50% 30.50% Break 83.50% 50.50% 27.50% 17.50% 83.11% 66.22% 55.95% 69.56% 57.11% 48.78% Avg.

Table 6: Ablation study in the AI2-THOR environment.

451 By comparing SELU and SELU w/o critic, we can see the importance of critic clearly. Only by 452 understanding the environment can we achieve the improvement of decision-making in all tasks. By 453 comparing SELU w/o HR and SELU w/o critic, we find that self-asking can correct the critic's com-454 prehension of the environmental task, but reflection on a single task is not enough. In Open tasks, 455 we find that the lack of hindsight relabeling directly leads to the disappearance for some instruc-456 tions, which causes the declined performance of success detection and decision-making. We can 457 observe the improvement from SELU w/o HR to SELU. By incorporating hindsight relabeling, we 458 can perform a comprehensive multi-task evaluation for each trajectory, ensuring that the embodied 459 MLLM achieves self-learning on each task. Consequently, self-asking and hindsight relabeling are 460 essential components of the critic.

461 462

463

440 441

442 443

444

445

446

447

448

449

450

#### 5.5 Hyperparameter Analysis

Online Dataset Size. Since the MLLMs fine-tuning process is sensitive to the dataset size, we
 explore the amount of interaction data required to achieve effective learning for embodied tasks. We
 conduct multiple tests on this variable based on picking up tasks in the AI2-THOR environment.
 The results are presented in Figure 3(a).

The size of dataset is described by the number of trajectories for a single task. For instance, dataset-10k indicates that 10k trajectories are collected for a specific task, such as picking up an apple, during online interactions. We evaluate the performance through the actor in terms of task success rate. We can see that the performance of dataset-1k is close to that of dataset-10k. For dataset-10k, it takes about two days to collect one epoch data. To balance experimental performance with sampling efficiency, we opt to sample 1k trajectories per task in our experiments.

Based on this conclusion, we employ a data augmentation method during the actor training process.
Since MLLMs tend to prioritize text over images when making decisions, they often focus excessively on the text prompt and overlook image comprehension. To address this issue, we shuffle the action lists in the prompts of training data to provide multiple prompts for the same image. This approach not only increases the dataset size, but also strengthens the connection between the MLLM policy and the observations.

Learning Rate. The learning rate is a critical factor, which prevents us from using the same MLLM for the actor and critic. We test different learning rates for the actor and critic separately on picking up tasks in the AI2-THOR, and the results are shown in Figure 3(b). We find that, due to the varying sizes of the fine-tuning datasets, using a uniform learning rate inevitably leads to overfitting or underfitting in one of the components, thereby impacting the overall performance of SELU. Therefore, we ultimately use two MLLMs to construct the SELU framework, ensuring effective self-learning. In our experiments, the learning rate for the critic is set to 2e-6, while that for the actor is set to 2e-5.



Figure 3: Hyperparameter study of SELU on picking up tasks in the AI2-THOR environment: (a) explores the size of the interaction dataset required for embodied MLLMs, (b) illustrates why a single MLLM is not suitable for SELU from the perspective of learning rare, and (c) demonstrates that the effect of multiple training iterations.

**Training Iterations.** The goal of our framework is to achieve multi-iteration self-learning improvement; therefore, we also evaluate the results of multiple rounds of fine-tuning on picking up tasks in the AI2-THOR environment. The results are presented in Figure 3(c).

505 Our results indicate that multiple iterations of fine-tuning do not consistently improve SELU's per-506 formance in every iteration. Both the actor and critic exhibit significant performance improvements 507 during the first iteration of fine-tuning, but show fluctuations and minimal growth in the subsequent iterations. We attribute this limitation to the performance of the critic. The critic performs the suc-508 cess detection based on the last frame of the trajectory, making it difficult to compare the quality 509 of successful trajectories. Once the actor reaches a level sufficient to roughly complete the task, 510 we lack the nuanced supervisory signals to guide the actor for further improvement. Consequently, 511 while there is a notable improvement after the first iteration, subsequent enhancements are limited. 512 We also attempt to increase the number of frames for success detection. However, the performance 513 of the small-scale MLLM does not meet our needs.

514 515 516

517

496

497

498

499

500 501

# 6 CONCLUSION AND LIMITATION

In this paper, we introduce SELU, a method for MLLMs to achieve self-learning in unknown en-518 vironments. SELU facilitates interaction with the environment, analyzes interaction trajectories, 519 builds an online dataset, and performs coupled optimization of the actor and critic. We employ 520 self-asking and hindsight relabeling to enhance the critic task evaluation capabilities. Ablation 521 experiments demonstrate that relabeling significantly expands the critic task judgment range. By 522 leveraging the principle that MLLMs possess stronger perceptual abilities than decision-making 523 abilities, we improve the performance of the actor policy. We test SELU in the AI2-THOR and 524 VirtualHome environments, achieving critic improvements of approximately 28% and 30%, and 525 policy improvements of about 20% and 24%, respectively. Additionally, to validate the applicability 526 of SELU across different MLLMs, we evaluate it on Qwen-VL, resulting in a 23% performance enhancement. 527

One limitation of SELU is the lack of a detailed evaluation of the trajectories that complete the tasks.
SELU can help embodied MLLM to self-learn how to accomplish tasks in unknown environments.
The next urgent problem to address is how to complete tasks more efficiently and optimize actions
more effectively. In future work, we will explore finer-grained critic signals to perform more accurate quality assessments of trajectories, guiding embodied MLLMs to tackle more complex tasks,
like long-horizon combination tasks.

- 534
- 535
- 536

537

538

540	REFERENCES
541	REI EREITEES

566 567

568

569

570

574

575

576

577

581

582

583

584

585

586

- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob
   McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. Advances in neural information processing systems, 30, 2017.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang
  Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020.
- Xiang Chen, Chenxi Wang, Ningyu Zhang, Yida Xue, xiaoyan yang, YUE SHEN, Jinjie GU, and
   Huajun Chen. Unified hallucination detection for multimodal large language models. In *ICLR* 2024 Workshop on Reliable and Responsible Foundation Models, 2024.
- Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. FinQA: A dataset of numerical reasoning over financial data. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 3697–3711. Association for Computational Linguistics, 2021.
  - Wonje Choi, Woo Kyung Kim, Minjong Yoo, and Honguk Woo. Embodied cot distillation from LLM to off-the-shelf agents. In *Forty-first International Conference on Machine Learning*, 2024.
  - Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe RLHF: Safe reinforcement learning from human feedback. In *The Twelfth International Conference on Learning Representations*, 2024.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
   bidirectional transformers for language understanding. In North American Chapter of the Association for Computational Linguistics, 2019.
  - Yihong Dong, Kangcheng Luo, Xue Jiang, Zhi Jin, and Ge Li. PACE: Improving prompt with actor-critic editing for large language model. In *Findings of the Association for Computational Linguistics ACL 2024*, pp. 7304–7323, 2024.
- Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. Towards revealing
  the mystery behind chain of thought: A theoretical perspective. In *Advances in Neural Informa- tion Processing Systems*, volume 36, pp. 70757–70798. Curran Associates, Inc., 2023.
  - Zelalem Gero, Chandan Singh, Hao Cheng, Tristan Naumann, Michel Galley, Jianfeng Gao, and Hoifung Poon. Self-verification improves few-shot clinical information extraction. In *ICML 3rd Workshop on Interpretable Machine Learning in Healthcare (IMLH)*, 2023.
  - Zhibin Gou, Zhihong Shao, Yeyun Gong, yelong shen, Yujiu Yang, Nan Duan, and Weizhu Chen. CRITIC: Large language models can self-correct with tool-interactive critiquing. In *The Twelfth International Conference on Learning Representations*, 2024.
- Patrick Haluptzok, Matthew Bowers, and Adam Tauman Kalai. Language models can teach themselves to program better. In *The Eleventh International Conference on Learning Representations*, 2023.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.

617

626

634

635

636

594	Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large
595	language models can self-improve. In Proceedings of the 2023 Conference on Empirical Meth-
596	ods in Natural Language Processing, pp. 1051–1068. Association for Computational Linguistics,
597	2023.
598	

- Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming
  Zhang, and Nenghai Yu. Opera: Alleviating hallucination in multi-modal large language models
  via over-trust penalty and retrospection-allocation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13418–13427, 2024.
- Chaoya Jiang, Haiyang Xu, Mengfan Dong, Jiaxing Chen, Wei Ye, Ming Yan, Qinghao Ye, Ji Zhang,
   Fei Huang, and Shikun Zhang. Hallucination augmented contrastive learning for multimodal large
   language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 27036–27046, 2024.
- Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward
   Grefenstette, and Roberta Raileanu. Understanding the effects of RLHF on LLM generalisation
   and diversity. In *The Twelfth International Conference on Learning Representations*, 2024.
- <sup>611</sup>
  <sup>612</sup>
  <sup>613</sup>
  <sup>614</sup>
  <sup>614</sup>
  <sup>614</sup>
  <sup>614</sup>
  <sup>614</sup>
  <sup>614</sup>
  <sup>614</sup>
  <sup>614</sup>
  <sup>614</sup>
  <sup>615</sup>
  <sup>614</sup>
  <sup>614</sup>
  <sup>615</sup>
  <sup>614</sup>
  <sup>616</sup>
  <sup>617</sup>
  <sup>617</sup>
  <sup>618</sup>
  <sup>618</sup>
  <sup>619</sup>
  <sup>619</sup>
  <sup>619</sup>
  <sup>619</sup>
  <sup>619</sup>
  <sup>611</sup>
  <sup>611</sup>
  <sup>612</sup>
  <sup>612</sup>
  <sup>613</sup>
  <sup>614</sup>
  <sup>614</sup>
  <sup>615</sup>
  <sup>615</sup>
  <sup>616</sup>
  <sup>617</sup>
  <sup>617</sup>
  <sup>618</sup>
  <sup>618</sup>
  <sup>619</sup>
  <
- Vijay Konda and John Tsitsiklis. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.
- Miaoran Li, Baolin Peng, Michel Galley, Jianfeng Gao, and Zhu Zhang. Self-checker: Plug-and-play modules for fact-checking with large language models. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 163–181. Association for Computational Linguistics, 2024a.
- Kiaoqi Li, Mingxu Zhang, Yiran Geng, Haoran Geng, Yuxing Long, Yan Shen, Renrui Zhang, Jiaming Liu, and Hao Dong. Manipllm: Embodied multimodal large language model for object-centric robotic manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18061–18070, 2024b.
- Zaijing Li, Yuquan Xie, Rui Shao, Gongwei Chen, Dongmei Jiang, and Liqiang Nie. OptimusHybrid multimodal memory empowered agents excel in long-horizon tasks. *arXiv preprint arXiv:2408.03615*, 2024c.
- Shalev Lifshitz, Keiran Paster, Harris Chan, Jimmy Ba, and Sheila A. McIlraith. STEVE-1: A
   generative model for text-to-behavior in minecraft. In *Thirty-seventh Conference on Neural In- formation Processing Systems*, 2023.
  - Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In Advances in Neural Information Processing Systems, volume 36, pp. 34892–34916. Curran Associates, Inc., 2023.
- Junpeng Liu, Yifan Song, Bill Yuchen Lin, Wai Lam, Graham Neubig, Yuanzhi Li, and Xiang
   Yue. Visualwebbench: How far have multimodal LLMs evolved in web page understanding and
   grounding? In *First Conference on Language Modeling*, 2024.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning
   of foundation models in visual contexts. In *The Twelfth International Conference on Learning Representations*, 2024.
- Kinbei Ma, Zhuosheng Zhang, and Hai Zhao. CoCo-agent: A comprehensive cognitive MLLM agent for smartphone GUI automation. pp. 9097–9110. Association for Computational Linguistics, 2024.

648 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri 649 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad 650 Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: 651 Iterative refinement with self-feedback. In Advances in Neural Information Processing Systems, 652 volume 36, pp. 46534–46594. Curran Associates, Inc., 2023. 653 Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document 654 images. In Proceedings of the IEEE/CVF winter conference on applications of computer vision, 655 pp. 2200-2209, 2021. 656 657 Nat McAleese, Rai Michael Pokorny, Juan Felipe Ceron Uribe, Evgenia Nitishinskaya, Maja Tre-658 bacz, and Jan Leike. Llm critics help catch llm bugs. arXiv preprint arXiv:2407.00215, 2024. 659 660 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong 661 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-662 low instructions with human feedback. Advances in neural information processing systems, 35: 27730-27744, 2022. 663 664 Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. 665 Automatically correcting large language models: Surveying the landscape of diverse automated 666 correction strategies. Transactions of the Association for Computational Linguistics, 12:484–506, 667 2024. 668 669 Jing-Cheng Pang, Pengyuan Wang, Kaiyuan Li, Xiong-Hui Chen, Jiacheng Xu, Zongzhang Zhang, 670 and Yang Yu. Language model self-improvement by reinforcement learning contemplation. In 671 The Twelfth International Conference on Learning Representations, 2024. 672 Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Tor-673 ralba. Virtualhome: Simulating household activities via programs. In Proceedings of the IEEE 674 Conference on Computer Vision and Pattern Recognition (CVPR), June 2018. 675 676 Siyuan Qi, Shuo Chen, Yexin Li, Xiangyu Kong, Junqi Wang, Bangcheng Yang, Pring Wong, Yi-677 fan Zhong, Xiaoyuan Zhang, Zhaowei Zhang, Nian Liu, Yaodong Yang, and Song-Chun Zhu. 678 Civrealm: A learning and reasoning odyssey in civilization for decision-making agents. In The 679 Twelfth International Conference on Learning Representations, 2024. 680 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea 681 Finn. Direct preference optimization: Your language model is secretly a reward model. Advances 682 in Neural Information Processing Systems, 36, 2024. 683 684 Timo Schick, Jane A. Yu, Zhengbao Jiang, Fabio Petroni, Patrick Lewis, Gautier Izacard, Qingfei 685 You, Christoforos Nalmpantis, Edouard Grave, and Sebastian Riedel. PEER: A collaborative 686 language model. In The Eleventh International Conference on Learning Representations, 2023. 687 688 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy 689 optimization algorithms. arXiv preprint arXiv:1707.06347, 2017. 690 Feifan Song, Bowen Yu, Minghao Li, Haiyang Yu, Fei Huang, Yongbin Li, and Houfeng Wang. 691 Preference ranking optimization for human alignment. In Proceedings of the AAAI Conference 692 on Artificial Intelligence, volume 38, pp. 18990-18998, 2024. 693 694 Sumedh Sontakke, Jesse Zhang, Séb Arnold, Karl Pertsch, Erdem B1 yık, Dorsa Sadigh, Chelsea 695 Finn, and Laurent Itti. Roboclip: One demonstration is enough to learn robot policies. In Ad-696 vances in Neural Information Processing Systems, volume 36, pp. 55681–55693. Curran Asso-697 ciates, Inc., 2023. 698 Dan Su, Xiaoguang Li, Jindi Zhang, Lifeng Shang, Xin Jiang, Qun Liu, and Pascale Fung. Read 699 before generate! faithful long form question answering with machine reading. In Findings of the 700 Association for Computational Linguistics: ACL 2022, pp. 744–756. Association for Computa-701 tional Linguistics, 2022.

702	Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan,
703	Liangyan Gui, Yu-Xiong Wang, Yiming Yang, Kurt Keutzer, and Trevor Darrell. Aligning large
704	multimodal models with factually augmented RLHF. In Findings of the Association for Com-
705	putational Linguistics ACL 2024, pp. 13088–13110. Association for Computational Linguistics,
706	2024.

- Weihao Tan, Wentao Zhang, Shanqi Liu, Longtao Zheng, Xinrun Wang, and Bo An. True knowledge comes from practice: Aligning large language models with embodied environments via reinforcement learning. In *The Twelfth International Conference on Learning Representations*, 2024a.
- Weihao Tan, Wentao Zhang, Xinrun Xu, Haochong Xia, Ziluo Ding, Boyu Li, Bohan Zhou, Junpeng Yue, Jiechuan Jiang, Yewen Li, Ruyi An, Molei Qin, Chuqiao Zong, Longtao Zheng, Yujie Wu, Xiaoqiang Chai, Yifei Bi, Tianbao Xie, Pengjie Gu, Xiyun Li, Ceyao Zhang, Long Tian, Chaojie Wang, Xinrun Wang, Börje F. Karlsson, Bo An, Shuicheng Yan, and Zongqing Lu. Cradle: Empowering foundation agents towards general computer control, 2024b.
- Xiaoyu Tan, Shaojie Shi, Xihe Qiu, Chao Qu, Zhenting Qi, Yinghui Xu, and Yuan Qi. Self-criticism: Aligning large language models with their understanding of helpfulness, honesty, and harmlessness. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pp. 650–662, Singapore, 2023. Association for Computational Linguistics.
- Jingqun Tang, Qi Liu, Yongjie Ye, Jinghui Lu, Shu Wei, Chunhui Lin, Wanqing Li, Mohamad Fitri
   Faiz Bin Mahmood, Hao Feng, Zhen Zhao, et al. Mtvqa: Benchmarking multilingual text-centric
   visual question answering. *arXiv preprint arXiv:2405.11985*, 2024.
- Heyi Tao, Sethuraman T V, Michal Shlapentokh-Rothman, Tanmay Gupta, Heng Ji, and Derek Hoiem. WebWISE: Unlocking web interface control for LLMs via sequential exploration. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 3693–3711. Association for Computational Linguistics, 2024.
- Belen Martin Urcelay, Andreas Krause, and Giorgia Ramponi. Reinforcement learning from human text feedback: Learning a reward model from human text input. In *ICML 2024 Workshop on Models of Human Feedback for AI Alignment*, 2024.
- Tianlu Wang, Ilia Kulikov, Olga Golovneva, Ping Yu, Weizhe Yuan, Jane Dwivedi-Yu, Richard Yuanzhe Pang, Maryam Fazel-Zarandi, Jason Weston, and Xian Li. Self-taught evaluators. *arXiv preprint arXiv:2408.02666*, 2024a.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha
  Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
  models. In *The Eleventh International Conference on Learning Representations*, 2023a.
- Yufei Wang, Zhanyi Sun, Jesse Zhang, Zhou Xian, Erdem Biyik, David Held, and Zackory Erickson. RL-VLM-f: Reinforcement learning from vision language foundation model feedback. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235, pp. 51484–51501. PMLR, 21–27 Jul 2024b.
- Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. JARVIS-1: Open-world multi task agents with memory-augmented multimodal language models. In *Second Agent Learning in Open-Endedness Workshop*, 2023b.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, pp. 24824–24837. Curran Associates, Inc., 2022.
- Licheng Wen, Xuemeng Yang, Daocheng Fu, Xiaofeng Wang, Pinlong Cai, Xin Li, Tao MA, Yingxuan Li, Linran XU, Dengke Shang, Zheng Zhu, Shaoyan Sun, Yeqi BAI, Xinyu Cai, Min Dou,
  Shuanglu Hu, Botian Shi, and Yu Qiao. On the road with GPT-4v(ision): Explorations of utilizing
  visual-language model as autonomous driving agent. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*, 2024.

756 757 758 750	Robin Winter, Marco Bertolini, Tuan Le, Frank Noe, and Djork-Arné Clevert. Unsupervised learn- ing of group invariant and equivariant representations. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), <i>Advances in Neural Information Processing Systems</i> , 2022.
760 761 762	Chuyan Xiong, Chengyu Shen, Xiaoqi Li, Kaichen Zhou, Jiaming Liu, Ruiping Wang, and Hao Dong. Autonomous interactive correction MLLM for robust robotic manipulation. In 8th Annual Conference on Robot Learning, 2024.
763 764 765 766	<ul> <li>Wenda Xu, Danqing Wang, Liangming Pan, Zhenqiao Song, Markus Freitag, William Wang, and Lei Li. INSTRUCTSCORE: Towards explainable text generation evaluation with automatic feedback. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i>, pp. 5967–5994. Association for Computational Linguistics, 2023.</li> </ul>
767 768 769 770	Xinrun Xu, Yuxin Wang, Chaoyi Xu, Ziluo Ding, Jiechuan Jiang, Zhiming Ding, and Börje F. Karlsson. A survey on game playing agents and large models: Methods, applications, and challenges, 2024.
771 772 773 774	Yiming Yan, Tao Wang, Chengqi Zhao, Shujian Huang, Jiajun Chen, and Mingxuan Wang. BLEURT has universal translations: An analysis of automatic metrics by minimum risk train- ing. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics</i> ( <i>Volume 1: Long Papers</i> ), pp. 5428–5443. Association for Computational Linguistics, 2023.
775 776 777 778	Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwen He, Yifeng Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu, Hai-Tao Zheng, Maosong Sun, and Tat-Seng Chua. Rlhf-v: Towards trustworthy mllms via behavior alignment from fine-grained correctional human feedback. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 13807–13816, 2024.
779 780 781 782	Jiawei Zhang, Chejian Xu, and Bo Li. Chatscene: Knowledge-enabled safety-critical scenario gener- ation for autonomous vehicles. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision</i> <i>and Pattern Recognition (CVPR)</i> , pp. 15459–15469, 2024.
783 784 785	Han Zhao, Remi Tachet Des Combes, Kun Zhang, and Geoffrey Gordon. On learning invariant representations for domain adaptation. In <i>International conference on machine learning</i> , pp. 7523–7532. PMLR, 2019.
786 787 788 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804	Yinmin Zhong, Zili Zhang, Bingyang Wu, Shengyu Liu, Yukun Chen, Changyi Wan, Hanpeng Hu, Lei Xia, Ranchen Ming, Yibo Zhu, et al. Rlhfuse: Efficient rlhf training for large language models with inter-and intra-stage fusion. <i>arXiv preprint arXiv:2409.13221</i> , 2024.
805 806 807 808	

مارم	rithm 1 SELU
Inni	<b>Itumi I</b> SEEO <b>it</b> : critic MITM $M$ actor MITM $M$ critic fine tuning detect $\mathcal{D}$ actor fine tunin
mpt	dataset $\mathcal{D}_{critic}$ , maximum timestep T initial instruction list L success detection prompt n
	and action plan prompt $p_a$
Out	<b>put:</b> critic MLLM $M_c$ , actor MLLM $M_c$
1:	$\mathcal{D}_{\text{critic}}, \mathcal{D}_{\text{actor}} \leftarrow \{\}$
2:	function "SELU"
3:	for instruct I in L do
4:	while data collecting not done do
5:	for timestep $t = 1$ to T do
6:	get observation $o_t$ from env
7:	$l_{a,t} = M_a(I, p_a, o_t)$
8:	use $l_{a,t}$ to interact with env
9:	end for
10:	$l_d = M_c(I, p_d, o_T)$ if $l_1 = \text{```ucc''}$ then
11.	store $(I \ n \ o \sigma \ l \ )$ into $\mathcal{D}$
12. 13·	store $(I, p_d, o_T, v_d)$ into $\mathcal{D}_{\text{critic}}$ store $(I, p_d, o_t, l_{a,t})$ $t = 1$ T into $\mathcal{D}_{\text{critic}}$
14:	else
15:	get $l'_{l}$ through self-asking
16:	if $l'_{d} = "yes"$ then
17:	store $(I, p_d, o_T, l'_d)$ into $\mathcal{D}_{\text{critic}}$
18:	store $(I, p_a, o_t, l_{a,t}^{a}), t = 1,T$ into $\mathcal{D}_{actor}$
19:	else
20:	get $I'$ through hindsight relabeling
21:	if $I' \neq "None"$ then
22:	store $(I', p_d, o_T, yes)$ into $\mathcal{D}_{\text{critic}}$
23:	store $(T, p_a, o_t, l_{a,t}), t = 1,T$ into $\mathcal{D}_{actor}$
24:	ena li ord if
25:	end if
20. 27.	end while
27.28	end for
29:	optimization $M_c$ and $M_a$ by $\mathcal{D}_{critic}$ and $\mathcal{D}_{actor}$
30:	return critic MLLM $M_c$ , actor MLLM $M_a$
31:	end function

Figure 4 shows our experiment environments. Both environments restrict agents to only interact with visible items, limiting their operational range to guarantee behavior plans realistic. Therefore, the actor MLLM makes decisions based on first-person perspective input to ensure accuracy as Figure 4(a) and Figure 4(c) show. Given the limitations of the first-person view, the critic MLLM uses a third-person perspective to evaluate the trajectory, reducing hallucinations and obtaining accurate scene information as Figure 4(b) and Figure 4(d) show.

854

The positioning of the third-person camera is crucial, as it should accurately capture the agent's position and the objects it interacts with. Any occlusion or interference can impair the MLLM's understanding of the image, thereby affecting the results of critic success detection and hindsight relabeling.



936 937

938 939

940

969

919	Table 8: Hyperparameters of Qwen-VL	fine-tuning b
920		
921	Hyperparameters	Value
922	train_batch_size	2
923	eval_batch_size	1
924	gradient_accumulation_steps	8
925	learning_rate_actor	1e-5
926	learning_rate_critic	1e-6
927	warmup_ratio	0.01
929	weight_decay	0.1
930	adam_beta2	0.95
931	model_max_length	2048
932	lr_scheduler_type	"cosine"
933	bf16	True
934	lazy_preprocess	True
935		1140

# A.3 VISUALIZATION OF ACTOR AND CRITIC ON LLAVA IN AI2-THOR

We use the embodied actor MLLM to interact with the unknown environment, and collect trajectories for evaluation from the critic MLLM. An example for 'pick up the lettuce' is as follows.



Figure 5: A visualization of the actor MLLM interacting with the AI2-THOR environment. The agent is instructed to pick up the lettuce. As the lettuce is far away, the agent needs to move closer before attempting to pick it up.

We use the critic MLLM to perform success detection on each trajectory and use self-asking and
 hindsight relabeling techniques to build bootstrapped dataset. An example for 'break the mug' is as follows.



Critic-Self Asking 1

The image shows a third-person view from the {agent}'s perspective in a {AI2-THOR/VirtualHome} environment. Please check the state of the {instruction.objects} in the image. You should output the state and the reasoning. The output format should be: State:...

Reasoning:...

# Critic-Self Asking 2

The image shows a third-person view from the {agent}'s perspective in a {AI2-THOR/VirtualHome} environment. The {instruction.objects} in the observation is in {objects.state} state, please determine whether the {instruction} has been completed or not. You should output yes or no, and the reasoning. The output format should be: Result:...

Reasoning:...

# Critic-hindsight relabeling 1

The image shows a third-person view from the {agent}'s perspective in a {AI2-THOR/VirtualHome} environment. Please see the image carefully. Determine whether there is any object that is {instruction.verb  $\rightarrow$  adj.} by the {agent}? You should output the object name and the reasoning. The output format should be: Object:...

Reasoning:...

# Critic-hindsight relabeling 2

The image shows a third-person view from the {agent}'s perspective in a {AI2-THOR/VirtualHome} environment. The {relabeling.object} in the observation is {instruction.verb  $\rightarrow$  adj.}, you should give a new instruction based on it. The original instruction is {instruction}, what's the new instruction? The output format should be: New instruction:...

Reasoning:...