

# ATLAS: Agent Tuning via Learning Critical Steps

Zhixun Chen<sup>1\*</sup>, Ming Li<sup>2\*</sup>, Yuxuan Huang<sup>3</sup>, Yali Du<sup>4</sup>, Meng Fang<sup>3</sup>, Tianyi Zhou<sup>2</sup>

<sup>1</sup>University of Technology Sydney, <sup>2</sup>University of Maryland,

<sup>3</sup>University of Liverpool, <sup>4</sup>King’s College London,

zhixun.chen@student.uts.edu.au, {minglii, tianyi}@umd.edu  
{yuxuan.huang, meng.fang}@liverpool.ac.uk, yali.du@kcl.ac.uk

## Abstract

Large Language Model (LLM) agents have demonstrated remarkable generalization capabilities across multi-domain tasks. Existing agent tuning approaches typically employ supervised finetuning on entire expert trajectories. However, behavior-cloning of full trajectories can introduce expert bias and weaken generalization to states not covered by the expert data. Additionally, critical steps—such as planning, complex reasoning for intermediate subtasks, and strategic decision-making—are essential to success in agent tasks, so learning these steps is the key to improving LLM agents. For more effective and efficient agent tuning, we propose ATLAS that identifies the critical steps in expert trajectories and finetunes LLMs solely on these steps with reduced costs. By steering the training’s focus to a few critical steps, our method mitigates the risk of overfitting entire trajectories and promotes generalization across different environments and tasks. In extensive experiments, an LLM finetuned on only  $\sim 30\%$  critical steps selected by ATLAS outperforms the LLM finetuned on all steps and recent open-source LLM agents. ATLAS maintains and improves base LLM skills as generalist agents interacting with diverse environments.

## 1 Introduction

An intelligent agent can perceive its environment, process information, and take actions to achieve specific goals or objectives (Maes, 1995; Wooldridge and Jennings, 1995). Traditional AI agents, especially empowered by Monte Carlo Tree Search (MCTS) or Reinforcement learning, have shown great potential in complex and difficult tasks (Silver et al., 2017; Vinyals et al., 2019). However, these methods are task-specific and suffer from weak generalization across different task domains or environments. On the other hand, the recent rise of Large Language Models (LLMs) (Brown et al., 2020; Team, 2024a; Touvron et al.,

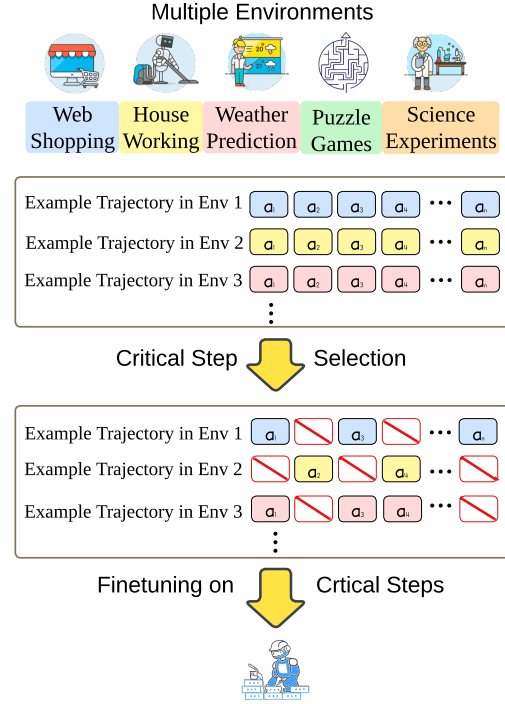


Figure 1: The proposed ATLAS identifies the critical steps in expert trajectories collected from diverse interactive environments and finetunes the agent on these steps only, where  $a_i$  represents the expert action in step- $i$ . ATLAS alleviates the potential overfitting to experts’ every-step behaviors and achieves better generalizability by training on much fewer steps (“less is better”).

2023a,b; Jiang et al., 2023) provide powerful foundation models with rich prior knowledge and reasoning capabilities, leading to remarkable generalization performance across various downstream applications (Zhao et al., 2023; Xu et al., 2024). Their unprecedented performance in different domains, including instruction following, reasoning, planning, and tool usage, makes LLMs ideal agents in multi-task settings (Wei et al., 2022; Kojima et al., 2022; Qin et al., 2023; Patil et al., 2023; Li et al., 2023, 2024a,b; Wang et al., 2024a).

Despite their effectiveness, transferring capabilities across tasks remains challenging due to

overfitting to expert trajectories and training inefficiency. Most existing agent tuning methods rely on imitation of every step of expert trajectories (Chen et al., 2023; Zeng et al., 2023; Chen et al., 2024; Zhang et al., 2024; Xi et al., 2024; Song et al., 2024). For instance, FireAct (Chen et al., 2023) generates diverse expert trajectories to fine-tune LLMs with enhanced planning abilities, while AgentTuning (Zeng et al., 2023) combines expert trajectories from various interactive environments with general instruction datasets to improve agents’ generalization capability.

However, fine-tuning on the entire expert trajectories can introduce expert bias, causing the model to overfit to specific behaviors and diminishing its ability to generalize to unseen environments (Ghosh et al., 2024). While experts in a few critical steps indeed provide key supervised information to improve the agent, the agent may already excel in other non-critical steps. Additionally, fitting the full trajectories of one task may lead to declines in others due to the potential distribution gap and negative transfer between tasks (Song et al., 2024). Moreover, imitating expert trajectories on redundant steps or sub-optimal/replaceable actions incurs excessive training costs.

To address the above challenges, we propose **ATLAS**, **A**gent **T**uning via **L**earning **C**ritical **S**teps, a novel approach that identifies the critical steps within expert trajectories and only finetunes LLMs on those selected steps, which is illustrated in Figure 1. By steering the training’s focus to a few critical steps, we not only reduce the backpropagation cost for training but also lower the overfitting risk to whole trajectories. Moreover, imitation on partial trajectories mitigates the expert bias, and encourages generalization across environments and tasks.

In ATLAS, an oracle LLM (GPT4o by default) is utilized as the selector to identify the critical steps semantically based on four criteria: key observations, plan formulation, recalling prior information, and pivotal actions. These critical steps are essential to the success of downstream tasks, but they are usually challenging to be generated by the base LLMs. Our training focuses on finetuning base models on these selected steps. Our experiments provide extensive comparisons between ATLAS and other agent-tuning strategies or LLM agents. As indicated by the results in Figure 2, ATLAS-finetuned agent maintains and improves the generalization capabilities of base LLMs in diverse environments, excelling on both

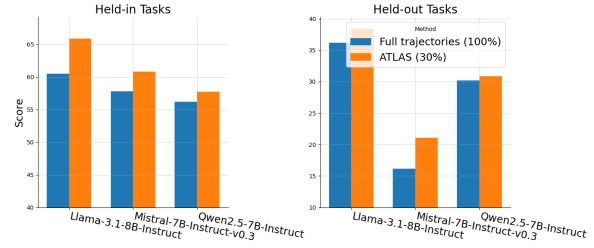


Figure 2: Three base LLMs finetuned by ATLAS vs. full trajectories (100% of the steps), evaluated on held-in and held-out agentic tasks. ATLAS consistently outperforms full-trajectory finetuning, indicating better generalizability of ATLAS by training on fewer but critical steps.

held-in and held-out tasks. Our main contributions can be summarized as:

- We introduce a novel method ATLAS reducing the tokens for agent tuning to 30%, by selecting the critical steps in expert trajectories.
- Agents finetuned on 30% critical steps outperforms the baseline agents finetuned on 100% steps, especially in the multi-task learning scenario, mitigating the expert bias and negative transfer across tasks.
- ATLAS-finetuned agent achieves a better generalization capability than baseline agents on not only held-in but also held-out tasks.

## 2 Preliminaries

As a multi-task learning framework, we define a set of  $N$  environments as  $\mathcal{E}$ . For any particular environment  $E \in \mathcal{E}$  and a task prompt space  $I$ , the agent’s dynamics is formalized as a Partially Observable Markov Decision Process (POMDP) (Sutton, 2018)  $\mathcal{M}_{E,I} \triangleq (S, A, O, T, r)$ . In this formalization,  $S$  denotes the state space,  $A$  denotes the action space,  $O$  denotes the observation space,  $T : S \times A \rightarrow S$  represents the state transition function, and  $r : S \times A \rightarrow \mathbb{R}$  defines the reward function.

Upon receiving an instruction  $i \in I$  in environment  $E$ , the agent takes action  $a_t \sim \pi_\theta(\cdot | o_t; E, i)$  according to its policy  $\pi_\theta$ , where  $\theta$  represents the policy parameters. For LLM agents,  $\pi_\theta$  is an LLM, and its output is not limited to actions and depends on the prompt. The environment then transitions to a new state  $s_{t+1} \in S$  based on the action, and the agent receives a new observation  $o_{t+1} \in O$ . This interaction continues as the agent engages with the

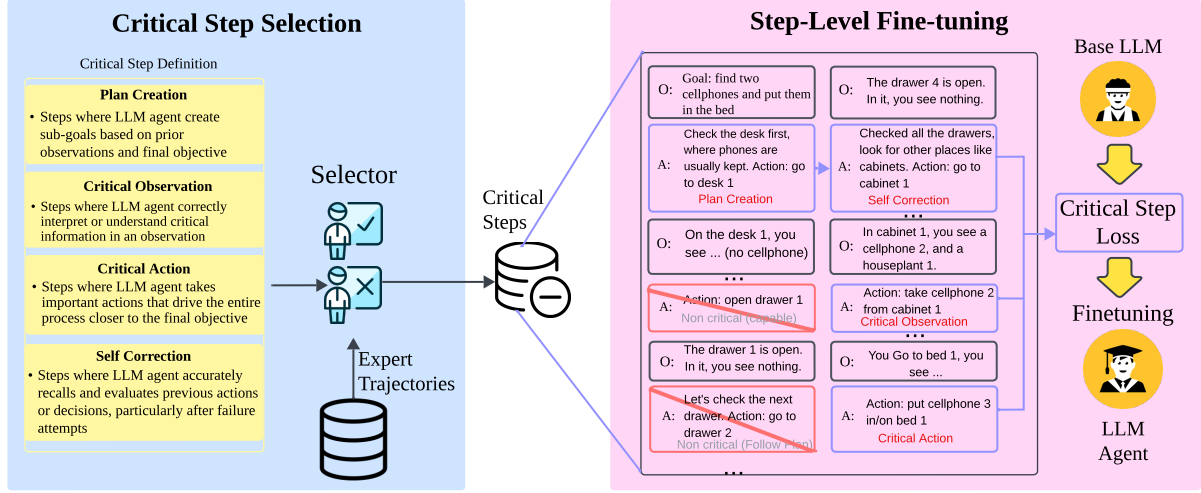


Figure 3: Overall of ATLAS. The selector identifies critical steps in expert trajectories collected in multiple environments, where “O” and “A” denote observation and action, respectively. Training loss is only computed on the critical steps. This encourages more exploration of non-critical steps, reduces the training cost, and improves the agent’s generalization performance.

environment until the task is completed or the maximum step limit is reached. We adopt the ReAct framework (Yao et al., 2022b) for the LLM agent: at each step  $t$ , the agent first generates reasoning thought  $h_t$  with a think prompt  $\text{prompt}_h$ . Define the trajectory up to step- $t$  as

$$\tau_{1:t} \triangleq (h_1, a_1, o_1, \dots, h_t, a_t, o_t), \quad (1)$$

The thought  $h_t$  is sampled by

$$h_t \sim \pi_\theta(\cdot | \tau_{1:t-1}; \text{prompt}_h, E, i). \quad (2)$$

The action is then drawn from the same model with another prompt  $\text{prompt}_a$ , i.e.,

$$a_t \sim \pi_\theta(\cdot | \tau_{1:t-1}, h_t; \text{prompt}_a, E, i). \quad (3)$$

The return of trajectory  $\tau = \tau_{1:T}$ ,  $G(\tau) = \sum_{t=1}^T \gamma^t r_t$  with a discounting factor  $\gamma$ , reflects the agent’s performance in task  $i$  within environment  $E$ . We collect a dataset  $D$  of expert trajectories for all environments in  $\mathcal{E}$  and tasks in  $I$ .

### 3 Methodology

In this section, we introduce our ATLAS, which identifies critical steps of the expert trajectories and finetunes LLM agents solely on them.

#### 3.1 Definition of Critical Steps

Critical steps are key actions or states within an agent’s trajectory that have a substantial impact on the return  $G(\tau)$ . As depicted in Figure 3, these steps are essential for ensuring the success of the

trajectory. We characterize critical steps as those requiring precise decision-making and flawless execution by the agent. In contrast, non-critical steps are more flexible, allowing for adjustments or rearrangements without compromising the overall trajectory. Typically, critical steps are challenging for the base model’s policy to generate, whereas non-critical steps are relatively easy for the base model. To identify these steps, we sample expert trajectories with dense rewards in available environments. We then empirically categorize four types of critical steps: **Plan Creation**, **Critical Observation**, **Critical Action**, and **Self Correction**.

#### 3.2 Selection of Critical Steps

Properly identifying and extracting critical steps from a trajectory is essential for ensuring effective analysis and learning. However, unlike instance-level data selection for LLMs, which has been explored (Li et al., 2024d,c), step-level selection is challenging since the task-related critical steps can differ dramatically by their semantic meaning. To address this issue, we leverage the selector to identify critical steps based on the critical step identification prompt,  $\text{prompt}_c$ . This approach ensures a more nuanced and context-aware assessment, helping to accurately identify the most influential steps for optimal performance.

The prompt initially defines the critical steps with the four categories introduced above:

- **Plan Creation:** Steps where the LLM agent formulates sub-goals by analyzing previous

observations and considering the final objective, breaking down the larger goal into manageable tasks that guide the agent’s actions toward the overall outcome.

- **Critical Observation:** Steps where the LLM agent identifies and analyzes key information from the environment, which helps the agent understand the objective or state and refine its strategy and decision-making toward more effective outcomes.
- **Critical Action:** Steps where the LLM agent takes decisive and impactful actions based on prior observations, significantly advancing the process toward the final objectives. These actions are crucial in shaping the direction of the agent’s strategy and are often pivotal moments that determine progress or failure, ensuring that the agent remains on track to achieve the desired outcome.
- **Self Correction:** Steps where the LLM agent carefully recalls and assesses its previous actions or decisions, especially after encountering failure or suboptimal outcomes. During this process, the agent reflects on what went wrong, identifies areas for improvement, and adjusts its approach to enhance future performance, which helps the agent refine its decision-making and better align with the overall objective.

After defining critical steps,  $\text{prompt}_c$  guides the selector to comprehend the given trajectory and summarize a high-level plan with sub-goals of the trajectory to enhance its understanding of the trajectory. Based on the sub-goals, the selector identifies the critical steps and categorizes them into four predefined categories. Notably, we require the selector to select at most  $m$  percentage of steps in the trajectory and return their indices as

$$C_\tau = \text{selector}(\tau; \text{prompt}_c, m) \quad (4)$$

ATLAS does not ignore all other steps in the trajectory during training. Instead, we keep the whole original trajectory, including the non-critical steps, as the input sequence but do not compute teacher-force loss on their tokens during training.

#### Critical Step Identification Prompts

1. Generate a high-level plan or strategy from the expert conversation, summarizing the key sub-goals required to complete the task.
2. Determine the most critical action steps in the expert conversation based on the plan and sub-goals, ensuring the number of selected steps does not exceed  $\{m\% \times \text{length}(\tau)\}$ .
3. Justify your selection of these critical steps, specifying the category to which each step belongs.

The full prompts are provided in Appendix D.1.

### 3.3 Agent Tuning on Critical Steps

Using the critical step dataset  $D_c$ , we optimize the model by calculating the loss exclusively on these critical steps. Training the agent on this reduced set of essential steps enhances both the efficiency and effectiveness of the fine-tuning process, ensuring the model focuses on the key actions that contribute to task success and gains knowledge rather than merely mimicking expert behavior. Thus, the objective function can be defined as:

$$\mathcal{J}(\theta) = \mathbb{E}_{\tau \sim D} \left[ \sum_{t \in C_\tau} \log \pi_\theta(h_t | \tau_{1:t-1}; \text{prompt}_h, E, i) + \log \pi_\theta(a_t | \tau_{1:t-1}, h_t; \text{prompt}_a, E, i) \right],$$

where  $\theta$  denotes the trainable parameter of the LLM agent. This expectation-based objective ensures that the model maximizes probabilities at the critical steps that most influence task success, following the same reasoning as in the full-trajectory approach but applied selectively to  $D_c$ .

## 4 Experiments

### 4.1 Environments

In our experimental framework, environments are categorized into **held-in** and **held-out** sets to evaluate the model’s performance. **Held-in environments** refer to the set of environments that are included in the training dataset. These are the environments in which the models can be fine-tuned, allowing them to learn and adapt their behavior based on the provided expert trajectories. The performance of the held-in environments measures how well the model has learned from the training data and its ability to replicate or improve upon expert behaviors within familiar settings. On the other hand, **held-out environments** consist of environments and tasks that are excluded from the training



process. These environments are used to assess the model’s generalization capabilities, determining how effectively it can apply learned knowledge to novel and unseen scenarios. The performance of held-out environments indicates the model’s ability to generalize its decision-making and problem-solving skills beyond the specific examples it was trained on, highlighting its potential for adaptability in diverse and dynamic real-world applications.

Specifically, our held-in environments include web navigation task Webshop (Yao et al., 2022a); textual digital games Textcraft (Prasad et al., 2023), Wordle and Maze (Abdulhai et al., 2023); embodied games Alfworld (Shridhar et al., 2020), Science-world (Wang et al., 2022) and Babyai (Chevalier-Boisvert et al., 2018); tool using tasks Weather, Todo and Movie (Ma et al., 2024). It is important to note that some environments provide a dense reward  $r \in [0, 1]$ , while others offer only sparse rewards  $r \in \{0, 1\}$ . To maintain simplicity and consistency, we follow the evaluation matrix definition used in previous work (Singh et al., 2024; Xi et al., 2024) across all tasks. For the held-out environments, we apply two tool usage environments Sheet and Academic, and two game environments Jericho and PDDL (Ma et al., 2024). A detailed description of these environments can be found in Appendix A.

## 4.2 Implementation Details

We apply the open-sourced dataset AgentTraj-L<sup>1</sup> (Xi et al., 2024) as our unfiltered dataset  $D_o$ , which contains the expert trajectories of all held-in tasks mentioned above. We then filter the dataset to obtain the critical step dataset  $D_c$  by prompting GPT-4o with the selection instruction  $I_c$ . For more details, please refer to Appendix D.1.

We employ Llama-3.1-8B-Instruct (Dubey et al., 2024) as the backbone model for our agent unless specified. During training, we focus on the critical step  $D_c$ , computing loss exclusively at these pivotal steps in the trajectory. We set the maximum selection ratio  $m$  for the critical step to be 30%. More information about the training setup is provided in Appendix B. The evaluation prompts are based on the React format (Yao et al., 2022b), consistent with the structure of the training data.

## 4.3 Baselines

We compare our LLM agent against three types of baselines: closed-source models, open-source models, and other fine-tuned LLM agents. Specifically, for closed-source models, we include GPT-3.5-Turbo (Ouyang et al., 2022), GPT-4-Turbo (Team, 2024a), Claude 3 (Anthropic, 2024), and DeepSeek-Chat (Liu et al., 2024). For open-source models, we evaluate Llama3-8B-Instruct, Llama3.1-8B-Instruct (Dubey et al., 2024), Gemma-2-9b-it (Gemma Team et al., 2024), Phi-3-mini-128k-instruct, Phi-3.5-mini-instruct (Team, 2024b), Mistral-7B-Instruct-v0.3 (Jiang et al., 2023) and Qwen2.5-7B-Instruct (Team, 2024c). For LLM agents fine-tuned on expert trajectories, we compare Xlam-7B-r (Zhang et al., 2024), AgentLM from AgentTuning (Zeng et al., 2023) and AgentEvol from AgentGym (Xi et al., 2024). We also fine-tune AgentTraj-L using the same base model, but with SFT on the unfiltered dataset  $D_o$  as a baseline for a fair comparison. Additionally, we introduce another baseline where the base model’s perplexity on all steps of the training dataset is used to select critical steps. Specifically, we choose the top 30% of steps with the highest perplexity and fine-tune on those steps.

## 4.4 Main Results

The results shown in Table 1 demonstrate the performance of various models on held-in and held-out tasks. For closed-source models, the reported outcomes follow those presented in AgentGym (Xi et al., 2024). On held-in tasks, our approach, fine-tuned on just 30% of the entire tokens of the expert trajectories, outperforms baseline methods across most tasks. It also achieves over a 5% average improvement compared to models fine-tuned with unfiltered data. Furthermore, in certain tasks such as Alfworld and Babyai, our model surpasses some advanced closed-source models. The perplexity-based selection performs poorly because it solely measures the token generation difficulty of one step by the base model, which is affected by both the step’s importance and expert bias. Consequently, using perplexity for selection forces the model to mimic the expert’s distribution rather than effectively learning the task knowledge to solve the problem.

The results on held-out tasks demonstrate that fine-tuning exclusively on critical steps enhances the generalization of agents to unseen tasks. Focus-

<sup>1</sup><https://huggingface.co/datasets/AgentGym/AgentTraj-L>

Model	held-in											held-out				
	Alfworld	Babyai	Maze	Movie	Sciworld	Textcraft	Todo	Weather	Webshop	Wordle	AVG	Academic	Sheet	Jericho	PDDL	AVG
Closed-Source Model																
DeepSeek-Chat	51.00	45.67	4.00	70.00	16.80	23.00	75.00	70.00	11.00	24.00	39.05	-	-	-	-	-
Claude-3-Haiku	0.00	1.93	4.00	50.00	0.83	0.00	65.00	55.00	5.50	16.00	19.83	-	-	-	-	-
Claude-3-Sonnet	13.00	79.25	4.00	50.00	2.78	38.00	80.00	65.00	1.50	36.00	36.95	-	-	-	-	-
GPT-3.5-Turbo	26.00	71.36	4.00	70.00	7.64	47.00	40.00	25.00	12.50	20.00	32.35	65.00	38.34	19.93	24.98	37.06
GPT-4-Turbo	67.50	72.83	68.00	95.00	14.38	77.00	95.00	80.00	15.50	88.00	67.32	80.00	75.83	52.44	81.16	72.36
Open-Source Base Model																
Llama-2-7B-Instruct	0.00	23.00	0.00	0.00	7.20	0.00	0.00	0.00	0.00	0.00	3.02	0.00	0.00	1.25	0.83	0.52
Llama-3-8B-Instruct	0.00	61.60	12.00	60.00	70.82	15.00	60.00	30.00	13.00	0.00	32.24	40.00	35.40	10.51	16.67	25.65
Llama-3.1-8B-Instruct	0.00	60.69	16.00	65.00	19.57	9.00	65.00	40.00	2.00	0.00	27.72	40.00	37.31	19.23	12.18	27.18
Phi-3-mini-128k-instruct	0.00	19.80	8.00	50.00	19.10	1.00	40.00	20.00	1.00	4.00	16.29	25.00	30.21	1.25	1.25	14.43
Phi-3.5-mini-instruct	0.00	21.00	12.00	60.00	18.20	3.00	60.00	30.00	2.00	6.00	21.22	30.00	33.62	0.56	0.83	16.25
Gemma2-9B-it	0.00	52.00	8.00	55.00	72.60	0.00	45.00	0.00	0.00	0.00	23.26	30.00	34.46	14.02	7.20	21.42
Mistral-7B-Instruct-v0.3	0.00	17.30	4.00	0.00	48.00	0.00	10.00	5.00	0.00	0.00	8.43	0.00	9.69	1.96	2.50	3.54
Qwen2.5-7B-Instruct	0.00	67.00	16.00	75.00	15.30	7.00	85.00	30.00	3.00	12.00	31.03	55.00	39.00	14.26	0.28	27.14
Fine-tuned Agents																
Xlam-7B-r	17.50	62.00	0.00	20.00	12.00	1.00	15.00	10.00	37.50	4.00	17.9	35.00	28.40	8.92	3.06	18.85
AgentLM-7B	71.00	0.49	12.00	5.00	1.63	4.00	15.00	0.00	36.50	4.00	14.96	10.00	13.30	15.88	2.78	10.49
AgentEvol-7B	88.00	82.70	12.00	60.00	38.00	64.00	70.00	25.00	76.50	12.00	52.82	25.00	26.20	4.05	3.10	14.59
AgentTraj-L (100% steps)	83.00	70.31	48.00	75.00	37.92	71.00	85.00	55.00	68.00	12.00	60.52	75.00	42.83	15.18	11.71	36.18
Perplexity Selection	78.50	77.59	24.00	65.00	53.54	59.00	70.00	30.00	66.50	8.00	53.21	30.00	36.48	13.28	8.76	22.13
ATLAS (30% steps)	84.50	80.98	48.00	90.00	42.02	72.00	90.00	60.00	71.50	20.00	65.91	70.00	49.39	18.21	15.84	38.36

Table 1: Performance of ATLAS-finetuned LLM agents vs. baseline LLMs and finetuned agents on held-in and held-out test tasks. “Unfilter” finetunes an LLM on the unfiltered dataset of complete expert trajectories. The best performance in each category is highlighted in **bold**. By finetuning LLMs only on the critical steps of expert trajectories, ATLAS achieves outstanding agent performance on both held-in and held-out tasks.

ing solely on these steps reduces the need to fully replicate the expert LLM’s distribution, allowing agents to learn the knowledge of how to solve problems effectively. Unlike other approaches, such as those described in (Zeng et al., 2023), which mix a large amount of general instruction data with task-specific trajectory data to enhance generalization, our method achieves comparable performance without relying on such datasets. This underscores the effectiveness of critical step training in boosting agents’ performance on held-out tasks.

To validate the effectiveness of our approach across various backbone models, we perform experiments using Mistral-7B-Instruct-v0.3 (Jiang et al., 2023) and Qwen2.5-7B-Instruct (Team, 2024c). The results, summarized in Table 2, demonstrate that our method consistently enhances performance across different models, yielding improvements on both held-in and held-out tasks. The detailed results are shown in Appendix C.

## 4.5 Ablation Study

### 4.5.1 Critical Steps vs. Non-Critical Steps

To assess the effectiveness of critical steps, we compare ATLAS with finetuning on non-critical steps **not** selected by the selector LLM in ATLAS. The number of steps for the non-critical selection dataset remains identical to our fine-tuning dataset, which is 30% of the entire training dataset.

The experimental results, presented in Table 2, indicate that fine-tuning with critical steps positively impacts the performance on both held-in and held-out tasks. In contrast, fine-tuning with

non-critical steps results in a negative impact across all tasks. These findings suggest that certain steps contribute positively to the agent’s learning by providing useful knowledge, whereas other steps introduce negative biases, leading to a decline in performance.

### 4.5.2 Different Selection Ratio

The ratio of critical steps plays a crucial role in the fine-tuning process. Therefore, we explore this hyperparameter and require the teacher model to select different proportions of steps from the expert trajectory as critical steps to identify the optimal ratio. As shown in Table 3, a 30% ratio yields the best agent performance on both held-in and held-out tasks. We also explore values of  $m$  greater than 30% but observed only a marginal increase in the number of critical steps selected by the selector. Consequently, we did not evaluate fine-tuning results for  $m$  larger than 30%. Motivated on these findings, we adopt  $m = 30\%$ . More detailed results can be found in Appendix C.

### 4.5.3 Critical Steps vs. Random Steps

To further evaluate the efficacy of using a selector LLM for selecting critical steps, we conduct comparisons to randomly selected steps of different ratios. Specifically, we randomly draw 30%, 20%, and 10% steps from the original trajectories and compare them with 30%, 20%, and 10% critical steps for agent tuning.

The results are summarized in Table 3, with only the average performance across tasks presented due

Data	held-in											held-out				
	Alfworld	Babyai	Maze	Movie	Sciworld	Textcraft	Todo	Weather	Webshop	Wordle	AVG	Academic	Sheet	Jericho	PDDL	AVG
Non-critical Steps (30%)	78.00	70.70	32.00	75.00	39.02	57.00	75.00	55.00	68.00	12.00	56.17	60.00	40.42	10.80	8.28	29.88
Complete Trajectories (100%) by AgentTraj-L	83.00	70.31	48.00	75.00	37.92	71.00	85.00	55.00	68.00	12.00	60.52	75.00	42.83	15.18	11.71	36.18
Critical Steps (30%) by ATLAS	84.50	80.98	48.00	90.00	42.01	72.00	90.00	60.00	71.50	20.00	65.91	70.00	49.39	18.21	15.84	38.36

Table 2: Critical vs. Non-critical steps & complete trajectories for agent tuning (tested on held-in and held-out tasks). The agent trained on ATLAS-selected critical steps does not only outperform the agent trained on full trajectories but also surpasses the agent trained on non-critical steps by a large margin. This indicates that the critical steps improve agent tuning while training on non-critical steps might be detrimental to the generalization.

Data	Steps	held-in	held-out
Complete	100%	60.52	36.18
Random	30%	<b>59.90</b>	<b>38.04</b>
Random	20%	59.08	34.73
Random	10%	58.68	29.04
Critical	30%	<b>65.91</b>	<b>38.36</b>
Critical	20%	60.40	36.96
Critical	10%	58.99	31.18

Table 3: ATLAS vs. random selection of steps at different ratios (including 100% as a reference).

to space constraints. Detailed results are provided in Appendix C. Our results indicate that increasing the ratio of random selection improves performance of the fine-tuned LLM agent on both held-in and held-out tasks. However, the random-selection’s performance at all the three ratios remains inferior to that achieved on ATLAS’s critical steps, particularly at higher ratios, highlighting the effectiveness of the selector.

#### 4.5.4 Value Function defined Critical Steps

To further assess the accuracy of using GPT to identify critical steps, we perform experiments to compare datasets selected based on the estimated value function. Following (Xiong et al., 2024), we estimate the value of expert actions by traversing all actions along the trajectory  $N$  times. Specifically, the agent begins from a specific action in the expert trajectory and performs rollouts from that action  $N$  times. The final rewards from these  $N$  rollouts are summed and averaged to estimate the value of the action. Subsequently, the differences in value between consecutive steps are calculated, and steps with differences exceeding a predefined threshold are identified as critical. We set  $N = 5$  for this experiment and a detailed methodology is provided in Appendix B.3. Notably, the estimated value function is designed for single-task training, so we compare the performance on the BabyAI and Weather tasks for evaluation. These environments are chosen because some others randomly generate goals upon reset, making rollouts infeasible.

Task	Method	Performance
BabyAI	ATLAS	78.93
	Value Func.	78.58
Weather	ATLAS	60.00
	Value Func.	60.00

Table 4: ATLAS vs. Value function (estimated oracle) on critical step selection for the BabyAI and Weather tasks. The value function follows (Xiong et al., 2024). It measures the importance of each step in MDP but it is computationally intractable in practice. This evaluation aims to verify the consistency between LLM-selected critical steps and their values in MDP.

Additionally, environments with long horizons are impractical due to the extensive time and computational resources required. For instance, the rollout of a dataset with 2000 trajectories, each averaging 25 steps, necessitates at least  $6.5 \times 10^5$  inferences, which is highly inefficient.

Table 4 shows that our method achieves performance comparable to the estimated value function approach. However, with limited rollouts, the estimates are often inaccurate, producing values of 0 at initial steps and 1 at later steps. As a result, only a few steps are identified as critical based on these variations. Despite this limitation, the comparable performance of the two methods suggests the potential for achieving better results with fewer critical steps, which we plan to explore in future work.

#### 4.5.5 Choices of Selector LLMs

In this section, we evaluate different models for selecting critical steps. Specifically, we use the open-source model Llama3.1-70B-Instruct<sup>2</sup> as the comparison teacher model to identify critical steps on three tasks Alfworld, BabyAI and Weather. The results in Table 5 indicate that GPT4o selected dataset can better fine-tune the LLM agents on all three tasks, showing that the selector model’s capability significantly impacts the quality of the critical step dataset. The dataset selected

<sup>2</sup><https://huggingface.co/meta-llama/Llama-3.1-70B-Instruct>

Task	Selector	Performance
Alfworld	GPT-4o	83.00
	Llama3.1-70B	78.50
BabyAI	GPT-4o	78.93
	Llama3.1-70B	67.23
Weather	GPT-4o	60.00
	Llama3.1-70B	55.00

Table 5: **Ablation study of the selector LLM** for critical step selection on Alfworld, BabyAI, and Weather tasks. Average performance over tasks is reported.

by Llama3.1-70B-Instruct fails to accurately identify critical steps, instead including many non-critical steps, which in turn reduces the agents’ performance after fine-tuning.

#### 4.6 Critical Step Verification

We conduct an additional ablation study to evaluate the improvement in the base model’s performance (LLaMA-3.1-8B-Instruct) when execution begins immediately after each identified critical step. In trajectories containing multiple critical steps, the model is reinitialized after each one in an iterative manner, ensuring that all critical steps are utilized throughout the trajectory.

Experimental results on a representative subset of 100 tasks from the BabyAI and Maze environments reveal notable performance improvements when the model is initialized from critical steps, compared to start from the beginning of the task. This indicates that the base model is generally capable of generating all non-critical steps autonomously and that training efforts should prioritize learning from critical steps specifically.

Tasks	From Scratch	From Critical Step
BabyAI	54.8	76.4
Maze	18.0	44.0

Table 6: Performance comparison of the model initialized from scratch vs. from critical steps on BabyAI and Maze tasks.

## 5 Related Work

### 5.1 LLM Agents

LLMs have demonstrated strong capabilities as general-purpose agents, effectively handling complex tasks across diverse environments without task-specific training. Several recent works explore different paradigms to enhance the reasoning, planning, and learning abilities of LLM-based agents. (Huang et al., 2022a) demonstrate that

LLMs can act as zero-shot planners, enabling agents to perform complex tasks directly from natural language instructions. Building on this, Inner Monologue (Huang et al., 2022b) proposes using internal dialogues for decision-making and planning, further eliminating the need for external training. Similarly, AUTOACT (Qiao et al., 2024) shows that agents can learn to solve question-answering tasks through self-generated strategies and planning. ReST meets ReAct (Aksitov et al., 2023) introduces a self-improvement framework in which agents enhance their multi-step reasoning through iterative feedback. TextGrad (Yuksekgonul et al., 2024) leverages LLMs for automatic differentiation via textual reasoning, bypassing conventional training pipelines.

In the realm of interactive and social intelligence, SOTOPIA (Wang et al., 2024b) emphasizes real-time human-agent communication for acquiring socially intelligent behavior. MMAC-Copilot extends this by employing collaborative chains to improve LLM agents’ performance in multimodal interaction scenarios. (Fang et al., 2024) demonstrate the effectiveness of integrating LLMs with external symbolic modules for zero-shot symbolic reasoning in text-based games. To improve agent performance in dynamic environments, MC-DML (Shi et al., 2025) combines LLMs with Monte Carlo planning for text-based decision-making. Lastly, HASARD (Tomilin et al., 2025) presents a vision-based safe reinforcement learning benchmark that challenges embodied agents with complex egocentric tasks, advancing the evaluation of safety in RL.

### 5.2 Tuning Agents

Fine-tuning large language models (LLMs) to enhance their performance in agent-related tasks has emerged as a key area of research, with various works proposing innovative frameworks, datasets, and methods. FireAct (Chen et al., 2023) fine-tunes models using various agent trajectories in the ReAct format, achieving substantial performance boosts. AgentTuning (Zeng et al., 2023) demonstrates the value of instruction tuning with datasets such as AgentInstruct, achieving superior task performance with its AgentLM models. AgentBank (Song et al., 2024) improves agent generalization and performance by fine-tuning LLMs with more than 50,000 diverse interaction trajectories. AgentOhana (Zhang et al., 2024) standardizes and aggregates agent trajectory datasets, facilitating efficient



pipelines and providing state-of-the-art results with its xLAM-v0.1 model. AGENTGYM (Xi et al., 2024) enables agent self-evolution through the interactive AGENTEVOL framework, which reduces human supervision through environmental feedback. The IPR framework (Xiong et al., 2024) refines LLM agents with step-level feedback, demonstrating strong results on complex benchmarks. Above all, these works highlight the critical role of fine-tuning in advancing LLM agent capabilities, leveraging diverse data, iterative feedback mechanisms, and innovative frameworks to improve generalization, robustness, and task efficiency.

## 6 Conclusion

In this work, we propose a novel method, ATLAS, for fine-tuning LLM agents by focusing on critical steps in expert trajectories. Using an oracle LLM as the selector, we first construct a dataset of critical steps, which highlights the most impactful actions in the trajectories. We then optimize the LLMs by computing and minimizing the loss exclusively on these critical steps. Experimental results demonstrate that by leveraging only 30% of the steps, our method achieves superior performance in training environments compared to finetuning on 100% of the trajectories. Moreover, the improved generalization ability of our finetuned agents is also reflected in enhanced performance on held-out tasks. These findings underscore the importance of critical step training in maintaining and enhancing the strengths of LLM agents while significantly reducing training costs. This approach offers a cost-effective solution for fine-tuning LLMs across diverse tasks and environments.

## 7 Limitation

While our proposed method, ATLAS, demonstrates significant improvements in efficiency and generalization for fine-tuning LLM agents on multi-environments, some aspects could be further improved. Firstly, the current approach to selecting critical steps depends heavily on powerful closed-source models. It is crucial to explore methods that enable precise critical step selection while minimizing time and computational costs. Additionally, the existing selection process primarily focuses on semantic aspects; combining it with other metrics could enhance the accuracy of identifying critical steps and further reduce the tokens needed for fine-tuning.

## References

- Marwa Abdulhai, Isadora White, Charlie Snell, Charles Sun, Joey Hong, Yuexiang Zhai, Kelvin Xu, and Sergey Levine. 2023. Lmrl gym: Benchmarks for multi-turn reinforcement learning with language models. *arXiv preprint arXiv:2311.18232*.
- 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. 2023. *Rest meets react: Self-improvement for multi-step reasoning llm agent*. *Preprint*, arXiv:2312.10003.
- Anthropic. 2024. *The claude 3 model family: Opus, sonnet, haiku*.
- 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. 2020. *Language models are few-shot learners*. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao. 2023. *Fireact: Toward language agent fine-tuning*. *arXiv preprint arXiv:2310.05915*.
- Zehui Chen, Kuikun Liu, Qiuchen Wang, Wenwei Zhang, Jiangning Liu, Dahua Lin, Kai Chen, and Feng Zhao. 2024. *Agent-flan: Designing data and methods of effective agent tuning for large language models*. *arXiv preprint arXiv:2403.12881*.
- Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Saharia, Thien Huu Nguyen, and Yoshua Bengio. 2018. *Babylai: A platform to study the sample efficiency of grounded language learning*. *arXiv preprint arXiv:1810.08272*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. *The llama 3 herd of models*. *arXiv preprint arXiv:2407.21783*.
- Meng Fang, Shilong Deng, Yudi Zhang, Zijing Shi, Ling Chen, Mykola Pechenizkiy, and Jun Wang. 2024. *Large language models are neurosymbolic reasoners*. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pages 17985–17993.
- Thomas Mesnard Gemma Team, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Laurent Sifre, Morgane Rivi re, Mihir Sanjay Kale, Juliette Love,

- Pouya Tafti, Léonard Hussenot, and et al. 2024. [Gemma](#).
- Sreyan Ghosh, Chandra Kiran Reddy Evuru, Sonal Kumar, Deepali Aneja, Zeyu Jin, Ramani Duraiswami, Dinesh Manocha, et al. 2024. A closer look at the limitations of instruction tuning. *arXiv preprint arXiv:2402.05119*.
- Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. 2022a. [Language models as zero-shot planners: Extracting actionable knowledge for embodied agents](#). *Preprint*, arXiv:2201.07207.
- Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, Pierre Sermanet, Noah Brown, Tomas Jackson, Linda Luu, Sergey Levine, Karol Hausman, and Brian Ichter. 2022b. [Inner monologue: Embodied reasoning through planning with language models](#). *Preprint*, arXiv:2207.05608.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Diederik P Kingma. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.
- Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, Jiuxiang Gu, and Tianyi Zhou. 2024a. [Selective reflection-tuning: Student-selected data recycling for LLM instruction-tuning](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 16189–16211, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, and Tianyi Zhou. 2023. [Reflection-tuning: Recycling data for better instruction-tuning](#). In *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*.
- Ming Li, Pei Chen, Chenguang Wang, Hongyu Zhao, Yijun Liang, Yupeng Hou, Fuxiao Liu, and Tianyi Zhou. 2024b. Mosaic-it: Free compositional data augmentation improves instruction tuning. *arXiv preprint arXiv:2405.13326*.
- Ming Li, Yong Zhang, Shwai He, Zhitao Li, Hongyu Zhao, Jianzong Wang, Ning Cheng, and Tianyi Zhou. 2024c. [Superfiltering: Weak-to-strong data filtering for fast instruction-tuning](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14255–14273, Bangkok, Thailand. Association for Computational Linguistics.
- Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. 2024d. [From quantity to quality: Boosting LLM performance with self-guided data selection for instruction tuning](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 7595–7628, Mexico City, Mexico. Association for Computational Linguistics.
- Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, et al. 2024. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *arXiv preprint arXiv:2405.04434*.
- Chang Ma, Junlei Zhang, Zhihao Zhu, Cheng Yang, Yujiu Yang, Yaohui Jin, Zhenzhong Lan, Lingpeng Kong, and Junxian He. 2024. Agentboard: An analytical evaluation board of multi-turn llm agents. *arXiv preprint arXiv:2401.13178*.
- Pattie Maes. 1995. Agents that reduce work and information overload. In *Readings in human-computer interaction*, pages 811–821. Elsevier.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.
- Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. 2023. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*.
- Archiki Prasad, Alexander Koller, Mareike Hartmann, Peter Clark, Ashish Sabharwal, Mohit Bansal, and Tushar Khot. 2023. Adapt: As-needed decomposition and planning with language models. *arXiv preprint arXiv:2311.05772*.
- Shuofei Qiao, Ningyu Zhang, Runnan Fang, Yujie Luo, Wangchunshu Zhou, Yuchen Eleanor Jiang, Chengfei Lv, and Huajun Chen. 2024. [Autoact: Automatic agent learning from scratch for qa via self-planning](#). *Preprint*, arXiv:2401.05268.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*.

- Zijing Shi, Meng Fang, and Ling Chen. 2025. Monte carlo planning with large language model for text-based game agents. *arXiv preprint arXiv:2504.16855*.
- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. 2020. Alfworld: Aligning text and embodied environments for interactive learning. *arXiv preprint arXiv:2010.03768*.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. 2017. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359.
- Avi Singh, John D. Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Xavier Garcia, Peter J. Liu, James Harrison, Jaehoon Lee, Kelvin Xu, Aaron Parisi, Abhishek Kumar, Alex Alemi, Alex Rizkowsky, Azade Nova, Ben Adlam, Bernd Bohnet, Gamaleldin Elsayed, Hanie Sedghi, Igor Mordatch, Isabelle Simpson, Izzeddin Gur, Jasper Snoek, Jeffrey Pennington, Jiri Hron, Kathleen Keanealy, Kevin Swersky, Kshiteej Mahajan, Laura Culp, Lechao Xiao, Maxwell L. Bileschi, Noah Constant, Roman Novak, Rosanne Liu, Tris Warkentin, Yundi Qian, Yamini Bansal, Ethan Dyer, Behnam Neyshabur, Jascha Sohl-Dickstein, and Noah Fiedel. 2024. Beyond human data: Scaling self-training for problem-solving with language models. *Preprint*, arXiv:2312.06585.
- Yifan Song, Weimin Xiong, Xiutian Zhao, Dawei Zhu, Wenhao Wu, Ke Wang, Cheng Li, Wei Peng, and Sujian Li. 2024. Agentbank: Towards generalized llm agents via fine-tuning on 50000+ interaction trajectories. *arXiv preprint arXiv:2410.07706*.
- Richard S Sutton. 2018. Reinforcement learning: An introduction. *A Bradford Book*.
- OpenAI Team. 2024a. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.
- Phi-3 Team. 2024b. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2407.10671*.
- Qwen2 Team. 2024c. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.
- Tristan Tomilin, Meng Fang, and Mykola Pechenizkiy. 2025. Hasard: A benchmark for vision-based safe reinforcement learning in embodied agents. *arXiv preprint arXiv:2503.08241*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models. *Preprint*, arXiv:2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.
- Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. 2019. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *nature*, 575(7782):350–354.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. 2024a. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345.
- Ruiyi Wang, Haofei Yu, Wenxin Zhang, Zhengyang Qi, Maarten Sap, Graham Neubig, Yonatan Bisk, and Hao Zhu. 2024b. Sotopia- $\pi$ : Interactive learning of socially intelligent language agents. *Preprint*, arXiv:2403.08715.
- Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. 2022. Scienceworld: Is your agent smarter than a 5th grader? *arXiv preprint arXiv:2203.07540*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Michael Wooldridge and Nicholas R Jennings. 1995. Intelligent agents: Theory and practice. *The knowledge engineering review*, 10(2):115–152.
- Zhiheng Xi, Yiwen Ding, Wenxiang Chen, Boyang Hong, Honglin Guo, Junzhe Wang, Dingwen Yang, Chenyang Liao, Xin Guo, Wei He, et al. 2024. Agentgym: Evolving large language model-based agents across diverse environments. *arXiv preprint arXiv:2406.04151*.



Weimin Xiong, Yifan Song, Xiutian Zhao, Wenhao Wu, Xun Wang, Ke Wang, Cheng Li, Wei Peng, and Sujian Li. 2024. Watch every step! Ilm agent learning via iterative step-level process refinement. *arXiv preprint arXiv:2406.11176*.

Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold Cheng, Jinyang Li, Can Xu, Dacheng Tao, and Tianyi Zhou. 2024. [A survey on knowledge distillation of large language models](#). *ArXiv*, abs/2402.13116.

Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022a. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022b. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations*.

Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Zhi Huang, Carlos Guestrin, and James Zou. 2024. [Textgrad: Automatic "differentiation" via text](#). *Preprint*, arXiv:2406.07496.

Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. 2023. Agenttuning: Enabling generalized agent abilities for llms. *arXiv preprint arXiv:2310.12823*.

Jianguo Zhang, Tian Lan, Rithesh Murthy, Zhiwei Liu, Weiran Yao, Juntao Tan, Thai Hoang, Liangwei Yang, Yihao Feng, Zuxin Liu, et al. 2024. Agentohana: Design unified data and training pipeline for effective agent learning. *arXiv preprint arXiv:2402.15506*.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. [A survey of large language models](#). *Preprint*, arXiv:2303.18223.

## A Environment Details

**ALFWorld** (Shridhar et al., 2020): ALFWorld is a simulated household environment developed upon the TextWorld framework, designed to require agents to perform tasks that involve spatial navigation and the application of common-sense knowledge. Within this environment, the action space includes interactions such as object manipulation (e.g., lifting, repositioning), environmental inspection, and the operational use of furniture. Agent actions are executed according to predefined logical rules, with the system generating contextually appropriate feedback. The principal evaluation criterion is the success rate, with each task sequence capped at 30 iterations to balance efficiency and feasibility. For the AGENTTRAJ-L dataset, a total of 2,420 trajectories were aggregated, consisting of 1,920 instances generated by state-of-the-art computational models and 500 human-annotated examples to ensure diversity and grounding in real-world reasoning.

**BabyAI** (Chevalier-Boisvert et al., 2018): BabyAI is a grid-world simulation platform featuring 40 instruction-following tasks that necessitate agent-object interaction. Agents are constrained to a 7x7 observational grid, permitting object manipulation only within immediate proximity. The original framework employs visual observations and primitive actions (e.g., "move forward," "turn left"), while an adapted version substitutes visual inputs with textual instructions and integrates high-level action commands (e.g., "pick up green key 1," "go through blue locked door 2"), thereby broadening the action space. A step-discounted reward is assigned upon successful task completion, with unsuccessful attempts yielding zero reward. For AGENTTRAJ-L, 810 trajectories were generated using state-of-the-art models. Performance is evaluated via task-specific rewards, capped at 20 rounds per task.

**MAZE** (Abdulhai et al., 2023): Maze is a textual environment where agents possess positional awareness, including their current location, target destination, and adjacent obstructions. Agents select movement in one of four directional axes (up, down, left, right) per step, with each action incurring a penalty of -1 until goal attainment. AGENTTRAJ-L comprises 215 trajectories. Performance is evaluated via success rate, with task sequences capped at 15 iterations.

**Movie** (Ma et al., 2024): Movie is a structured



environment enabling LLM agents to leverage a specialized tool for accessing film-related metadata (e.g., cinematic details, personnel, production entities). The framework offers 16 discrete actions to execute task-specific objectives. It integrates The Movie Database (TMDB) API, incorporating its dataset and operational functions to facilitate agent interactions. Agents obtain a binary reward (1 for correct output alignment with reference solutions, 0 otherwise). AGENTTRAJ-L comprises 215 annotated trajectories, with 20 questions reserved for evaluation and the remainder allocated to training. Success rate serves as the primary performance metric, with task sequences limited to 12 iterations.

**SciWorld** (Wang et al., 2022): SciWorld functions as an empirical assessment framework for evaluating scientific reasoning capabilities in an interactive textual environment, structured to reflect elementary-level curricular standards. The platform encompasses 30 diverse task categories, operationalizing activities such as tool-based measurement and mechanical experimentation. It employs a domain-specific action space, with a simulator generating contextual feedback on action outcomes. The AGENTTRAJ-L dataset includes 2,120 trajectories synthesized for evaluation. Task performance is measured via reward accumulation, with each trial restricted to a maximum of 30 iterations.

**TextCraft** (Prasad et al., 2023): TextCraft is a text-based environment modeled after WordCraft, designed for simulating item crafting through a hierarchical framework of 544 nodes aligned with Minecraft’s recipe hierarchy. Each task specifies a target item and a sequence of compositional crafting actions (e.g., “craft <item> using <ingredients>,” “get <item>,” “inventory”), with complexity ranging from 1 to 4 procedural steps. The environment provides real-time feedback on action validity and execution states, enabling agents to directly acquire non-craftable items. A reward of 1 is granted exclusively upon successful synthesis of the target item. For evaluation, 100 tasks are partitioned from a broader training set, with AGENTTRAJ-L comprising 374 trajectories (299 generated by state-of-the-art models, 75 human-annotated). Success rate is the primary metric, capped at 20 rounds per task for empirical validation.

**TODOList** (Ma et al., 2024): Todolist is a structured environment enabling LLM agents to manage personal agenda data through a task-oriented interface with 11 discrete operational commands. The

tool integrates the TodoList API to operationalize task management functionalities. Agents receive binary reward allocation (1 for congruence with reference outputs, 0 otherwise). AGENTTRAJ-L comprises 135 trajectories, with 20 task instances reserved for evaluation and the remainder allocated to training. Performance is assessed via success rate, constrained to a maximum of 15 iterations per task.

**Weather** (Ma et al., 2024): Weather is a structured environment enabling LLM agents to retrieve meteorological parameters (e.g., temperature, precipitation, air quality) across spatiotemporal contexts via a specialized tool. The framework offers 18 discrete operational commands, integrating the Open-Meteo API through Python-based implementation to enable data querying functionalities. Agents receive a binary reward (1 for output alignment with reference solutions, 0 otherwise). AGENTTRAJ-L includes 311 annotated trajectories, with 20 queries reserved for evaluation and the remainder allocated to training. Performance is assessed via success rate, constrained to a maximum of 10 iterations per task.

**WebShop** (Yao et al., 2022a): Webshop is a simulated e-commerce platform where agents execute product procurement tasks adhering to predefined criteria through interface interactions (button-based navigation) or text-based search functionality. The environment integrates 12,000 structured instructions and leverages over one million real-world Amazon product listings, with 6,910 instructions selected for task execution. AGENTTRAJ-L includes 3,930 annotated trajectories. Performance is quantified via success rate, with task sequences limited to 10 rounds to balance efficiency and practical applicability.

**Wordle** (Abdulhai et al., 2023): Wordle is a lexical reasoning assessment framework where agents deduce a target word from a constrained five-letter vocabulary. The environment operationalizes character-level feedback after each guess, indicating positional accuracy and presence of characters. Agents accumulate a step penalty of -1 until successful identification or attempt exhaustion. AGENTTRAJ-L contains 955 trajectories. Performance is quantified via success rate, with trials capped at 8 rounds for empirical validation.

**Academia** (Ma et al., 2024): Academia is a structured environment enabling LLM agents to access computer science research resources (e.g., publications, author metadata) via seven discrete

operational commands. The tool integrates the Citation Network Dataset to implement core functionalities. Agents receive binary rewards (1 for output alignment with reference solutions, 0 otherwise). Performance is assessed via success rate across 20 evaluation tasks, with trials capped at 12 rounds per task.

**Sheet** (Ma et al., 2024): Sheet is a structured environment enabling LLM agents to manipulate spreadsheet data via 20 discrete operational commands, leveraging the Google Sheets API. Reward is determined by structural and content congruence between the agent-modified spreadsheet and a reference template, quantified on a 0–1 scale. Evaluation employs 20 predefined tasks, with performance measured via reward accumulation and trials capped at 15 rounds per task.

**Jericho** (Ma et al., 2024): Jericho is a text-based game framework designed to evaluate agents’ capacity for interactive exploration and dynamic world modeling within fictional narratives. These tasks demand agents to infer contextual rules (e.g., magical systems) through iterative interaction rather than pre-existing commonsense knowledge. To address operational feasibility given LLM agents’ context window constraints, original game objectives (often requiring 50–300 steps) are restructured into modular subtasks achievable within 15 steps. For example, the zork1 dungeon exploration is redefined as “locate the secret passage entrance in the living room,” reducing the sequence to 8 steps. This adaptation preserves core exploratory challenges while aligning task complexity with computational tractability.

**PDDL** (Ma et al., 2024): The Planning Domain Definition Language (PDDL) serves as a framework for evaluating strategic reasoning in symbolic planning tasks, with four benchmark domains—Gripper, Barman, Blocksworld, and Tyreworld—designed to test multi-step action sequencing under efficiency constraints. Agents must navigate domain-specific objectives (e.g., transporting objects, mixology, block stacking, tire installation) by optimizing action sequences to minimize redundant operations. For instance, in the Barman domain, agents strategically allocate containers to reduce cleaning cycles during cocktail preparation. To enable natural language interaction, symbolic PDDL predicates (e.g., `ontable(shaker1)`) and actions (e.g., `clean-shaker`) are translated into textual observations (e.g., “Shaker1 is on the table”) and instructions (e.g., “Clean shaker1

with hand1 while hand2 is empty”). Each domain features 10–20 curated multi-round problems, with progress measured via a normalized matching score that quantifies alignment between the current state and the goal state. Full task completion (100% progress) requires satisfying all goal conditions, such as hierarchical block placement in Blocksworld. This adaptation bridges symbolic planning with language-agent interoperability while preserving strategic complexity.

## B More Implementation Details

### B.1 Training Configuration

We use the Adam optimizer (Kingma, 2014), with a learning rate of  $2e-5$  and a cosine scheduler for the agent fine-tuning. The models are trained with 3 epoches and a warmup rate 0.03. The batch size is 128 and the max length of 8192. All experiments are conducted on 4 NVIDIA A100 80G GPUs and fine-tunes a Llama3.1-8B-Instruct on ATLAS takes approximately 8 hours. We use PyTorch FSDP (Paszke et al., 2019) for efficient training.

### B.2 Perplexity Selection

The perplexity-based selection method identifies critical steps by measuring the generation difficulty of each step within a trajectory. We compute perplexity as the exponentiated average negative log-likelihood of a step’s tokens under the model’s distribution. Formally, the perplexity for step  $x$  is defined as:

$$\text{PPL}(x) = \exp \left( -\frac{1}{|x|} \sum_{t=1}^{|x|} \log p(x_t \mid x_{1:t-1}) \right) \quad (5)$$

where  $|x|$  is the token length of the step, and  $p(x_t \mid x_{1:t-1})$  is the model’s probability for token  $x_t$  given its context. Higher perplexity values indicate steps that are more challenging for the model to generate.

Instead of using a fixed threshold, we dynamically select the top 30% of steps with the highest perplexity within each trajectory.

### B.3 Value Estimation

We follow the rollout setup of (Xiong et al., 2024), which defines the rollout estimate reward function as  $r_s(s_t, a_t)$  as the estimated outcome reward from the exploration starting step  $t$ . The agent with policy  $\pi_s$  is employed to rollout trajectory  $\tau_{t:n}$  from

step  $t$ , based on the historical trajectory  $\tau_{t-1}$ . The environment then gives a final reward  $r_o(i, \tau_n, E)$  for the trajectory. The estimated reward can be calculated as:

$$r_s(s_t, a_t) = \mathbb{E}_{\tau_n \sim \pi_s(\tau_{t:n} | \tau_{t-1})} [r_o(i, \tau_n, E)] \quad (6)$$

Due to the complexity of computing this expectation value directly, we utilize the Monte Carlo sampling technique for estimation. By drawing  $N$  trajectories from step  $t$  using  $\pi_s$ , we create a collection of trajectories:

$$\{\tau^{(i)} | i = 1, \dots, N\} = \text{MC}_{\pi_s}(\tau_{t-1}; N) \quad (7)$$

The estimated reward is then calculated as:

$$r_s(s_t, a_t) = \begin{cases} \frac{1}{N} \sum_{i=1}^N r_o(i, \tau_n, E), & \text{for } t < n, \\ r_o(i, \tau_n, E), & \text{for } t = n. \end{cases} \quad (8)$$

Then the value function  $V^\pi(s_t)$  is estimated as the expectation over all possible trajectories starting from state  $s_t$  under policy  $\pi$ :

$$V^\pi(s_t) = \mathbb{E}_{\tau_{t:n} \sim \pi(\tau_{t:n} | s_t)} \left[ \sum_{k=t}^n \gamma^{k-t} r_s(s_k, a_k) \right], \quad (9)$$

where  $\gamma$  is the discount factor which we set to be 0.99 here. Finally, we compute the difference between consecutive steps. If this difference exceeds the threshold of 0.1, the step is identified as a critical step. At this point, we designate the expert action corresponding to the same timestep  $t$  as trainable and calculate the loss associated with that action.

## C Additional Results

Detailed results of other backbone models, random selection and different critical selection ratio  $m$  are shown in Table 7.

## D Prompts Details

### D.1 Prompts of Critical Step Selection

#### Critical Selection Prompts

A critical step is defined as a key action or decision that, if performed correctly, significantly increases the likelihood of successfully completing the task. It represents a turning point in the process that influences the outcome of subsequent actions. More specifically, critical steps include:

- **Plan Creation:** Steps where the LLM agent formulates sub-goals by analyzing previous observations and considering the final objective, breaking down the larger goal into manageable tasks that guide the agent's actions towards the overall outcome.
- **Critical Observation:** Steps where the LLM agent identifies and analyzes key information from the environment, which help agent understand the objective or state and refine its strategy and decision-making towards more effective outcomes.
- **Critical Action:** Steps where the LLM agent takes decisive and impactful actions based on prior observations, significantly advancing the process toward the final objectives. These actions are crucial in shaping the direction of the agent's strategy and are often pivotal moments that determine progress or failure, ensuring that the agent remains on track to achieve the desired outcome.
- **Self Correction:** Steps where the LLM agent carefully recalls and assesses its previous actions or decisions, especially after encountering failure or sub-optimal outcomes. During this process, the agent reflects on what went wrong, identifies areas for improvement, and adjusts its approach to enhance future performance, which helps the agent refine its decision-making and better align with the overall objective.

**Task Description:** Your task is:

1. Induce a high-level plan or strategy based on the expert conversation, summarizing the key steps needed to successfully complete the task.
2. Based on this high-level plan, identify the most critical action steps in the expert conversation. A maximum of  $\{ \}$  steps may be chosen from the conversation.
3. Provide a detailed explanation for choosing these critical steps, specifying which category (e.g., key observation, planning, recall, pivotal action) they belong to and why mastering these steps ensures the success of the task.

**Answer Format:**

1. **The high-level plan is:** [Summarize the strategy and key steps for task completion]
2. **The critical steps are:** conversation[...]
3. **Reason:** [Explain why these steps are critical, including which category they fall into (key observation, planning, recall, pivotal action) and how they enable the player to avoid mistakes in subsequent steps]

Model	held-in											held-out				
	Alfworld	Babyai	Maze	Movie	Sciworld	Textcraft	Todo	Weather	Webshop	Wordle	AVG	Sheet	Academic	Jericho	PDDL	AVG
Mistral-7B-Instruct-v0.3																
AgentTraj-L	79.50	76.43	44.00	75.00	47.32	76.00	75.00	25.00	72.00	8.00	57.83	35.00	22.39	6.55	2.50	16.61
Ours	85.00	77.73	40.00	85.00	56.59	69.00	85.00	30.00	72.00	8.00	60.83	45.00	30.44	6.32	2.53	21.07
Qwen2.5-7B-Instruct																
AgentTraj-L	78.50	73.48	32.00	80.00	39.51	55.00	90.00	40.00	65.50	8.00	56.20	55.00	41.56	14.52	9.77	30.21
Ours	76.50	80.90	40.00	80.00	41.32	64.00	80.00	40.00	62.50	12.00	57.72	55.00	39.96	16.76	11.72	30.86
Random Selection																
40%	83.50	73.08	60.00	85.00	37.51	69.00	75.00	50.00	69.50	16.00	61.86	55.00	40.71	18.97	15.96	32.66
30%	83.50	71.70	40.00	80.00	38.32	67.00	80.00	50.00	68.50	20.00	59.90	70.00	48.63	14.80	18.71	38.04
20%	81.00	72.30	36.00	80.00	39.98	65.00	80.00	60.00	64.50	12.00	59.08	65.00	43.35	13.93	16.64	34.73
10%	78.50	72.78	32.00	80.00	41.02	66.00	80.00	55.00	65.50	16.00	58.68	55.00	36.95	9.80	14.39	29.04
Critical Selection Ratio																
30%	84.50	80.98	48.00	90.00	42.01	72.00	90.00	60.00	71.50	20.00	65.91	70.00	49.39	18.21	15.84	38.36
20%	82.50	73.77	48.00	80.00	34.68	72.00	80.00	50.00	71.00	12.00	60.40	70.00	43.90	25.01	8.94	36.96
10%	78.00	83.76	44.00	75.00	36.13	64.00	80.00	50.00	67.00	12.00	58.99	60.00	35.63	13.01	16.08	31.18

Table 7: Detailed performance of other backbone models, random selection and different critical selection ratio  $m$ . The best performance in each section is highlighted in **bold**.

## D.2 Prompts for Task Inference

We adopt the setup of AgentGym (Xi et al., 2024) and utilize the same prompts for task inference. For further details, please refer to Appendix G: Prompt Details of AgentGym.