

# WALL-E: WORLD ALIGNMENT BY RULE LEARNING IMPROVES WORLD MODEL-BASED LLM AGENTS

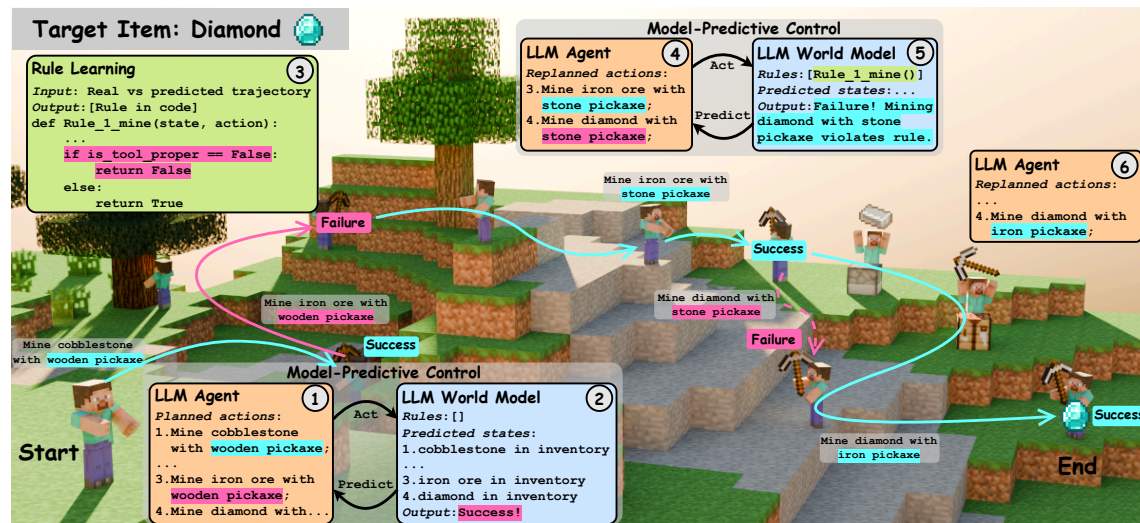


Figure 1: **Illustration of WALL-E mining a diamond in Minecraft.** Step 1-2: the agent makes a plan via MPC with the initial unaligned world model, resulting in a failed action for mining iron ore. Step 3: by comparing real trajectories with the world model predictions, WALL-E learns a critical rule that if the tool is not proper to the material being mined, the action will fail. Step 4-5: the learned rule helps the world model make accurate predictions for transitions that were predicted mistakenly in MPC. Step 6: the agent accordingly modifies its plan and replaces *stone pickaxe* with an *iron pickaxe* toward completing the task.

## ABSTRACT

Can large language models (LLMs) directly serve as powerful world models for model-based agents? While the gaps between the prior knowledge of LLMs and the specified environment’s dynamics do exist, our study reveals that the gaps can be bridged by aligning an LLM with its deployed environment and such “world alignment” can be efficiently achieved by rule learning on LLMs. Given the rich prior knowledge of LLMs, only a few additional rules suffice to align LLM predictions with the specified environment dynamics. To this end, we propose a neurosymbolic approach to learn these rules gradient-free through LLMs, by inducing, updating, and pruning rules based on comparisons of agent-explored trajectories and world model predictions. The resulting world model is composed of the LLM and the learned rules. Our embodied LLM agent “WALL-E” is built upon model-predictive control (MPC). By optimizing look-ahead actions based on the precise world model, MPC significantly improves exploration and learning efficiency. Compared to existing LLM agents, WALL-E’s reasoning only requires a few principal rules rather than verbose buffered trajectories being included in the LLM input. On open-world challenges in Minecraft and ALFWorld, WALL-E achieves higher success rates than existing methods, with lower costs on replanning time and the number of tokens used for reasoning. In Minecraft, WALL-E exceeds baselines by **15-30%** in success rate while costing **8-20** fewer replanning rounds and only **60-80%** of tokens. In ALFWorld, its success rate surges to a new record high of **95%** only after **6** iterations. Code is available here.

## 1 INTRODUCTION

While large language models (LLMs) have been successfully applied to complex reasoning, generation, and planning tasks, they are not sufficiently reliable to be deployed as an agent in specific open-world environments, e.g., games, VR/AR systems, medical care, education, autonomous driving, etc (OpenAI, 2023; Wei et al., 2022; Liu et al., 2024). A primary reason for the failures is the gap between the commonsense reasoning with prior knowledge of pretrained LLMs and the specified, hard-coded environment’s dynamics, which leads to incorrect predictions of the future states, hallucinations, or violation of basic laws in LLM agents’ decision-making process (Mu et al., 2023b; Yang et al., 2024; Das et al., 2018; Wu et al., 2024). Although the alignment of LLMs with human preferences has been widely studied as a major objective of LLM post-training, “world alignment” with an environment’s dynamics has not been adequately investigated in building LLM agents (Hao et al., 2023; Rafailov et al., 2024; Ge et al., 2024). Moreover, many existing LLM agents are model-free and their actions are directly executed in real environments without being verified or optimized in advance within a world model or simulator (Mu et al., 2023b; Yao et al., 2023; Shinn et al., 2024; Zitkovich et al., 2023; Wu et al., 2023; Micheli & Fleuret, 2021; Brohan et al., 2022). This leads to safety risks and suboptimality of generated trajectories.

In this paper, we show that aligning an LLM with environment dynamics is both necessary and crucial to make it a promising world model, which enables us to build more powerful embodied agents. In particular, we introduce a neurosymbolic world model that composites a pretrained LLM with a set of newly learned rules from the interaction trajectories with the environment. This specific form of world model combines the strengths of both in modeling the environment dynamics, i.e., (1) the rich prior knowledge, probabilistic, and deductive reasoning capability of LLMs (Hu & Shu, 2023); and (2) the hard constraints and rigorous guarantees enforced by rules (Li et al., 2024a). While creating a rule-only world model for a complex environment is challenging due to the massive amount of rules and uncertainty (Xiao et al., 2021), in our method, only a few complementary rules suffice to align a pretrained LLM to specific environment dynamics. This is achieved by simply including these rules in the LLM’s prompt without tedious training or inference. In contrast, existing LLM agents usually require expensive finetuning of LLMs via RL/imitation learning on trajectory data, or memory-heavy inference with a long input context of buffered trajectories (Mu et al., 2023b; Gao et al., 2023a; Yang et al., 2024; Shinn et al., 2024).

To this end, we propose “World Alignment by ruLe LEarning (WALL-E)”, which builds the neurosymbolic world model by learning complementary rules with LLMs’ inductive reasoning and code generation capability. Specifically, in each iteration, WALL-E interacts with the environment to collect a real trajectory and compare it with the world model predictions. The comparison results are then analyzed by an LLM, which extracts new rules or modifies existing ones to improve the consistency between the predicted and real trajectories. To keep the rule set minimal necessarily, at the end of each iteration, we prune the rules by solving a maximum coverage problem, which aims to select a subset of rules with the maximal coverage of the transitions failed being predicted by the LLM in the world model (without applying any rules). Hence, the selected rules are complementary to the LLM predictions. The above rule learning procedure repeats for multiple iterations until the LLM+rules performs as an accurate world model.

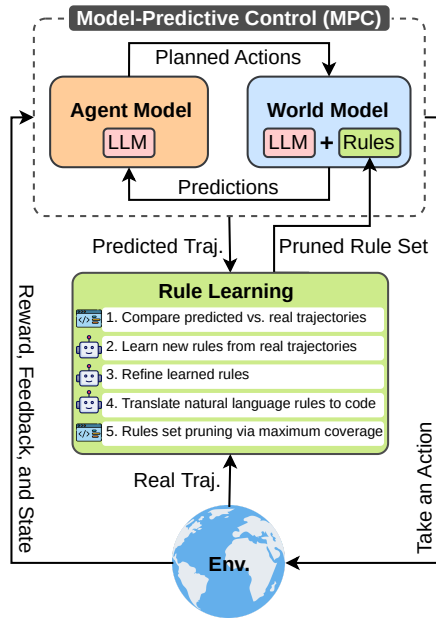


Figure 2: **Overview of WALL-E.** The agent’s action per step is controlled by MPC, where the agent model plans actions in a look-ahead window based on the LLM+rules world model’s predictions.

The precise world model achieved by WALL-E enables us to create better model-based LLM agents for challenging open-world tasks. However, model-based reinforcement learning (RL) of LLM agents in complex environments is still hindered by the expensive exploration and finetuning of LLMs. In this paper, we revisit the classical idea of model-predictive control (MPC) (Qin & Badgwell, 2003; Hafner et al., 2019; 2020; 2023), compared to RL, which does not require training a policy network but needs to optimize actions for a look-ahead time window in every step. To reduce the optimization cost per step, we instead apply the LLM agent as an optimizer searching for the optimal look-ahead actions by interacting with the WALL-E’s world model. With an aligned world model and an efficient LLM-based optimizer, MPC leads to a more promising and efficient framework of LLM agents in open-world environments.

We evaluate WALL-E on challenging tasks in open-world environments such as Minecraft and ALFWorld where the agents can explore freely and target complicated tasks. Our main contributions are threefold.

- We investigate the underexplored “world alignment” challenge for LLM agents.
- We propose a novel class of neurosymbolic world models based on rule learning on LLMs.
- We develop LLM agents based on model-predictive control (MPC) with the neurosymbolic world model.

## 2 RELATED WORK

Recent studies have integrated LLMs with rule learning to improve reasoning and generalization capabilities across various tasks, including numerical reasoning, knowledge graph exploration, and adherence to predefined rules (Yang et al., 2023a; Zhu et al., 2023c; Mu et al., 2023a; Yang et al., 2023b; Luo et al., 2023). However, prior work has not focused on aligning LLM-based world models with dynamic environments. Our research addresses this gap by applying rule learning to enhance model-based agent performance in such contexts. Several works have also used LLMs to construct world models for task planning by translating natural language into representations or combining LLMs with task-specific modules (Wong et al., 2023; Guan et al., 2023; Tang et al., 2024). Unlike these approaches, we directly employ LLMs as world models, leveraging their inherent knowledge for greater flexibility and efficiency. While some works use LLMs as world models, typically relying on fine-tuning or human defined prompts for alignment with environment dynamics (Xiang et al., 2024; Xie et al., 2024; Zhao et al., 2024; Hao et al., 2023; Liu et al., 2023). Our method advances this by automatically learning rules through exploration, reducing human intervention and improving performance. For a more comprehensive discussion of related work, please refer to Appendix A.

## 3 METHOD

### 3.1 MODEL-PREDICTIVE CONTROL (MPC) OF WORLD MODEL-BASED LLM AGENTS

We consider a scenario where a LLM, denoted as  $f$ , is deployed in a dynamic environment for agent interaction over discrete time steps. At each time step  $t$ , the agent observes the current state  $s_t$ , selects an action  $a_t$ , and transitions to the next state  $s_{t+1}$ . This transition is represented as  $\delta_t = (s_t, a_t, s_{t+1})$ . A trajectory  $\tau = (\delta_0, \delta_1, \dots, \delta_{T-1})$  comprises a sequence of such transitions, capturing the agent’s behavior from the initial to the terminal state within an episode.

The LLM-based world model  $f_{\text{wm}}$  predicts the subsequent state  $\hat{s}_{t+1}$  based on the current state and action:

$$\hat{s}_{t+1} = f_{\text{wm}}(s_t, a_t), \quad (1)$$

Model Predictive Control (MPC) is a widely recognized framework for model-based control. In this context, we integrate MPC with the LLM-based world model  $f_{\text{wm}}$  to enhance agent planning and decision-making, the whole framework is illustrated in Figure 2. The objective is to determine an optimal sequence of actions

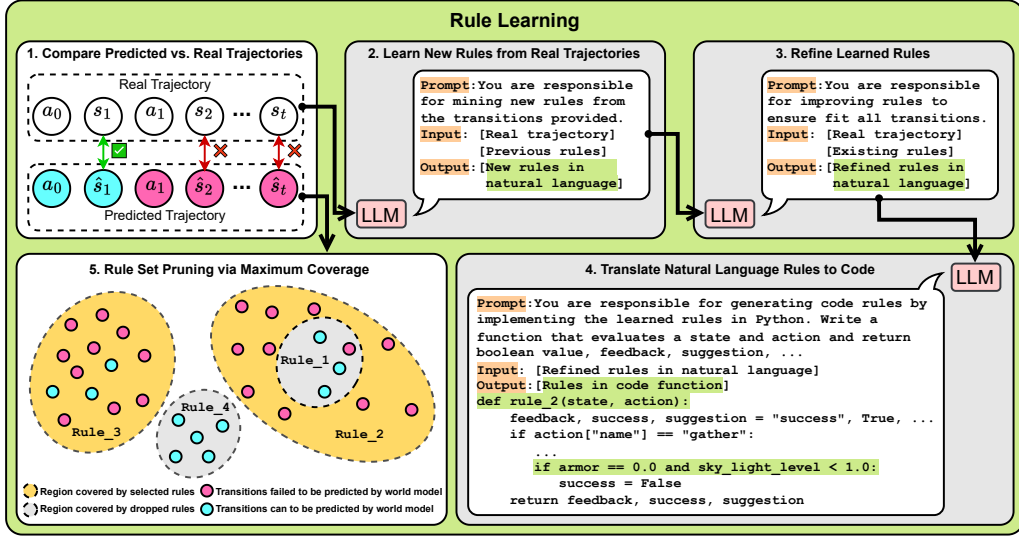


Figure 3: **Rule Learning details.** The rule learning module iteratively refines the rules by comparing the world model predicted trajectories with the agent’s actual trajectories in the environment. The rule learning takes five steps: (1) comparing predicted and actual trajectories; (2) learning new rules from real trajectories; (3) refining learned rules; (4) translating natural language rules to code; and (5) rule set pruning via solving a maximum coverage problem. (2)-(4) are handled by LLMs, while (1) and (5) are executed by programs.

$a_{t:t+H}$  over a finite horizon  $H$  that maximizes the expected cumulative reward. At each time step  $t$ , the optimization problem is formulated as:

$$a_{t:t+H}^* = \arg \max_{a_{t:t+H}} \mathbb{E} \left[ \sum_{i=0}^H \gamma^i \mathcal{F}(\hat{s}_{t+i+1}) \right], \quad (2)$$

where  $\gamma$  is the discount factor, and  $\mathcal{F}(\hat{s}_{t+i+1})$  denotes the reward function.

However, if the LLM-based world model is misaligned with the actual environment dynamics, the predicted state  $\hat{s}_{t+1}$  may not match the true state  $s_{t+1}$ . This misalignment leads to incorrect reward evaluations, resulting in inaccurate cumulative reward estimates. Consequently, the derived action sequence  $a_{t:t+H}^*$  may be suboptimal or erroneous, leading to ineffective control decisions by the agent. Therefore, addressing the misalignment between the LLM world model and the environment’s true dynamics is crucial for ensuring optimal performance within the MPC framework.

### 3.2 WORLD ALIGNMENT BY RULE LEARNING (WALL-E)

In complex environments, direct state prediction is challenging due to complexity and randomness. To address this, our world model uses a two-stage approach: first, it assesses `action_result` (e.g., success or failure), then generates the subsequent `state_info` (provides state details) based on the action success:

$$\hat{s}_{t+1} = (\text{action\_result}_{t+1}, \text{state\_info}_{t+1}) = f_{\text{wm}}(s_t, a_t), \quad (3)$$

To address potential misalignment between the  $f_{\text{wm}}$  and the real environment, we introduce a rule learning framework, illustrated in Figure 3 and detailed in the following sections. The learned rules align the  $f_{\text{wm}}$  with the environment, enhancing state prediction accuracy and improving agent performance within the MPC framework.

**Comparing Predicted and Real Trajectories.** To find misalignments between the LLM world model and the real environment, we compare action outcomes in predicted and actual next state, focusing on the binary action\_result rather than detailed state\_info. This focus provides a reliable basis for identifying discrepancies. Let the predicted trajectories be  $\tau^{\text{predicted}} = \{\delta = (s_t, a_t, \hat{s}_{t+1})\}_{t=0}^T$ . Then, we may divide  $\tau^{\text{predicted}}$  into correct and incorrect transition set, and correct the wrong  $\hat{s}_{t+1}$  (see Step 1 of rule learning in Figure 3):

$$\begin{aligned} \mathcal{D}^{\text{correct}} &= \{\delta_t^{\text{correct}} = (s_t, a_t, \hat{s}_{t+1}) \mid \hat{s}_{t+1} = s_{t+1}\}, \\ \mathcal{D}^{\text{incorrect}} &= \{\delta_t^{\text{incorrect}} = (s_t, a_t, \hat{s}_{t+1}) \mid \hat{s}_{t+1} \neq s_{t+1}\}, \end{aligned} \quad (4)$$

where  $s_{t+1}$  is the true state given by environment. Then  $\tau^{\text{predicted}} = \mathcal{D}^{\text{correct}} \cup \mathcal{D}^{\text{incorrect}}$ . By analyzing  $\mathcal{D}^{\text{incorrect}}$ , we pinpoint where the model’s predictions diverge from reality, highlighting areas needing correction through additional rules.

**Learning New Rules from Real Trajectories.** Before address these misalignments, we prompt the LLM  $f_{\text{gen}}$  to generate new natural language rules from real trajectories  $\tau^{\text{real}}$  (see Appendix B.1 for detailed prompt). The LLM is given the task setup and state-action structures to infer new natural language rules  $R_{\text{new}}^{\text{NL}}$  that explain the observed dynamics, ensuring they are distinct from previous rules  $R_{\text{previous}}^{\text{NL}}$ :

$$R_{\text{new}}^{\text{NL}} = f_{\text{gen}}(\tau^{\text{real}}, R_{\text{previous}}^{\text{NL}}), \quad (5)$$

**Refining Learned Rules.** Then, we prompt the LLM to update existing rules based on the real trajectories  $\tau^{\text{real}}$  (see Appendix B.2 for detailed prompt). Early-stage rules could be inaccurate due to data drift caused by the limited data, so the LLM identifies conflicting rules and modifies or discards them as needed. The set of all existing rules up to the current point is  $R_{\text{existing}}^{\text{NL}} = R_{\text{previous}}^{\text{NL}} \cup R_{\text{new}}^{\text{NL}}$ , where the LLM  $f_{\text{refine}}$  refines these rules with the real trajectories:

$$R^{\text{NL}} = f_{\text{refine}}(\tau^{\text{real}}, R_{\text{existing}}^{\text{NL}}). \quad (6)$$

**Translating Natural Language Rules to Code.** The next step is translating the refined natural language rules  $R^{\text{NL}}$  into executable code. We prompt the LLM  $f_{\text{code\_gen}}$  to produce the code-based rule set  $R^{\text{code}}$  (see Appendix B.3 for detailed prompt):

$$R^{\text{code}} = f_{\text{code\_gen}}(R^{\text{NL}}), \quad (7)$$

**Rule Set Pruning via Maximum Coverage.** In the final step, to address the inherent uncertainty and variability in the LLM-driven rule-learning process, we programmatically verify and refine the rule set to reduce dependence on the LLM. The code-based rules  $R^{\text{code}}$  are executed and validated against the labeled predicted transitions  $\tau^{\text{predicted}}$ . Any rule that fails to predict a transition correctly is discarded, ensuring that only accurate and effective rules are retained.

We further optimize the rule set by selecting rules that maximize coverage of the incorrectly predicted transitions  $\delta_t^{\text{incorrect}}$ , where the LLM world model’s failed. This approach focuses our efforts on correcting the most significant misalignments between the LLM and the environment. We formulate this optimization as a maximum set cover problem.  $\mathcal{D}^{\text{incorrect}} = \{\delta_1^{\text{incorrect}}, \delta_2^{\text{incorrect}}, \dots, \delta_n^{\text{incorrect}}\}$  is the set of incorrectly predicted transitions, and  $R^{\text{code}} = \{R_1^{\text{code}}, R_2^{\text{code}}, \dots, R_m^{\text{code}}\}$  is the set of code-based rules. Our goal is to select a minimal subset of rules that maximizes coverage of  $\mathcal{D}^{\text{incorrect}}$ :

$$\max_{\mathbf{x} \in \{0,1\}^m, \mathbf{y} \in \{0,1\}^n} \left\{ \sum_{j=1}^n y_j - \lambda \sum_{i=1}^m x_i \mid y_j \leq \sum_{i=1}^m x_i a_{ij}, \forall j = 1, \dots, n \right\}, \quad (8)$$

where  $x_i$  indicates whether rule  $R_i^{\text{code}}$  is selected ( $x_i = 1$ ) or not ( $x_i = 0$ ),  $y_j$  indicates whether transition  $\delta_j^{\text{incorrect}}$  is covered ( $y_j = 1$ ) or not ( $y_j = 0$ ), and  $a_{ij} = 1$  if transition  $\delta_j^{\text{incorrect}}$  is covered by rule  $R_i^{\text{code}}$ ,

235  $a_{ij} = 0$  otherwise. The constraint ensures that a transition  $\delta_i^{\text{incorrect}}$  is considered covered if it is included in  
 236 at least one selected rule. The parameter  $\lambda > 0$  balances minimizing the number of rules and maximizing  
 237 transition coverage; we set  $\lambda$  to be very small to prioritize coverage maximum. We solve this optimization  
 238 problem using a greedy algorithm (see Appendix F).

239 Through this process, we eliminate **rules covering only correct transitions**, as they do not address mis-  
 240 alignments, and **redundant rules** fully covered by more comprehensive ones (see Step 5 of rule learning in  
 241 Figure 3). This pruning process results in a pruned rule set that is both efficient and effective in correcting  
 242 the LLM’s misalignments. Additionally, any code-based rules removed from  $R^{\text{code}}$  are also excluded from  
 243 the set of natural language rules  $R^{\text{NL}}$ .

### 245 3.3 INFERENCE ON LLM AGENTS WITH LEARNED RULES

246 After completing the rule learning process, we obtain rules in two distinct forms: natural language rules  
 247  $R^{\text{NL}}$  and code-based rules  $R^{\text{code}}$ . Both types of rules enhance the LLM world model’s ability to predict the  
 248 next state  $\hat{s}_{t+1}$  within the planning framework: For **natural language rules**, these can be embedded directly  
 249 into the LLM’s input prompt to guide the model’s predictions, e.g.,  $\hat{s}_{t+1} = f_{\text{wm}}(s_t, a_t, R^{\text{NL}})$ . For **code-**  
 250 **based rules**, these are applied programmatically after the LLM generates its initial prediction, e.g.,  $\hat{s}_{t+1} =$   
 251  $\text{ApplyRules}(f_{\text{wm}}(s_t, a_t), R^{\text{code}})$ . Here, the function `ApplyRules` serves as a verification layer, overriding the  
 252 LLM’s prediction if an active rule contradicts the generated outcome. For further details on rule activation,  
 253 refer to Appendix G.

254 By integrating learned rules, the aligned LLM world model enhances the agent’s planning process signifi-  
 255 cantly. This alignment allows the agent to more effectively obtain optimal action sequences  $a_{t:t+H}$  through  
 256 two key improvements: First, the alignment leads to more **accurate reward evaluations**  $\mathcal{F}(\hat{s}_{t+1})$ , increas-  
 257 ing the likelihood of selecting optimal action sequences  $a_{t:t+H}$  within the MPC framework. Second, the  
 258 aligned world model, equipped with learned rules, provides **high-quality feedback** that helps the agent  
 259 refine  $a_{t:t+H}$  effectively. Along with predicting action results and state information, it offers auxiliary infor-  
 260 mation when an action is predicted to fail, including:

- 262 • Feedback: A textual explanation of the failure based on violated rules.
- 263 • Suggestion: Recommendations for corrective actions or improvements based on the current state,  
 264 action taken, and violated rules.

265 This information is crucial when an action fails, guiding the agent in revising its strategy by exploring  
 266 alternatives or adjusting its approach(see Appendix D.2 for examples).

267 In conclusion, integrating learned rules improves the LLM world model’s prediction accuracy and provides  
 268 actionable feedback, enabling more efficient and adaptive planning.

## 271 4 EXPERIMENTS

272 We evaluate the environment modeling and task-solving capabilities of WALL-E on open-world environ-  
 273 nments using the Minecraft (Fan et al., 2022) and ALFWorld (Shridhar et al., 2020b) benchmarks. Compared  
 274 to state-of-the-art (SOTA) LLM/VLM agents, WALL-E achieves higher success rates with lower costs in  
 275 terms of replanning time and token usage for reasoning. Notably, in Minecraft, WALL-E surpasses base-  
 276 lines by 15–30% in success rate while costing 8–20 fewer replanning rounds and only 60–80% of tokens.  
 277 In ALFWorld, it achieves a record of 95% success rate only after 6 iterations, significantly exceeding other  
 278 planning-based methods such as RAFA (Liu et al., 2023). Moreover, integrated with our proposed rule  
 279 learning method, WALL-E achieves a 15% higher success rate than methods relying on a long input context  
 280 of buffered trajectories. These highlights demonstrate WALL-E’s superior efficiency and effectiveness in  
 281 complex and open-world environments.

Table 1: Comparison of WALL-E and baselines on Minecraft tasks for success rate (%) and replanning rounds. \*-reported in previous work. VLMs = vision-language models, LLMs = large language models. The best score for each task is highlighted in **bold**. **WALL-E substantially exceeds other SOTA LLM/VLM agents and is the only method that performs better than human players in the Minedojo benchmark.**

Method		Success Rate (%) $\uparrow$ (Replanning Rounds $\downarrow$ )						
		Avg.	Wooden	Stone	Iron	Golden	Diamond	Redstone
VLMs	GPT-4V* (Li et al., 2024b)	10(-)	41(-)	21(-)	0(-)	0(-)	0(-)	0(-)
	Jarvis-1* (Wang et al., 2023b)	42(-)	94(-)	89(-)	36(-)	7(-)	9(-)	16(-)
	Optimus-1* (Li et al., 2024b)	<b>47(-)</b>	<b>99(-)</b>	<b>92(-)</b>	<b>47(-)</b>	<b>9(-)</b>	<b>12(-)</b>	<b>25(-)</b>
LLMs	GPT-3.5* (Li et al., 2024b)	10(-)	40(-)	20(-)	0(-)	0(-)	0(-)	0(-)
	DEPS (Wang et al., 2023a)	37(35.36)	83(10.67)	41(33.26)	33(35.27)	22(45.29)	24(42.46)	17(45.22)
	GITM (Zhu et al., 2023b)	54(25.49)	96(3.42)	<b>92(6.01)</b>	57(23.93)	29(37.17)	30(39.80)	22(42.63)
	WALL-E w/o WM	61(23.13)	94(5.04)	89(9.58)	<b>67(18.56)</b>	33(39.67)	41(32.73)	43(33.21)
	<b>WALL-E (ours)</b>	<b>69(15.77)</b>	<b>98(1.64)</b>	91( <b>4.58</b> )	63(19.38)	<b>69(15.61)</b>	<b>46(27.08)</b>	<b>48(26.33)</b>
<b>Human Performance*</b> (Li et al., 2024b)		59(-)	100(-)	100(-)	86(-)	17(-)	17(-)	33(-)

Table 2: Comparison of WALL-E and baselines on Minecraft tasks for average token usage and API costs (in USD). The number of tokens is calculated as the sum of prompt tokens and generation tokens. The average API cost is derived by separately calculating the costs of prompt and generation tokens and then summing both. The lowest cost for each task is highlighted in **bold**.

Method	Inference Tokens $\downarrow$ (Cost in USD $\downarrow$ )						
	Avg.	Wooden	Stone	Iron	Golden	Diamond	Redstone
DEPS	93560.95(0.65)	28223.33(0.20)	87999.46(0.61)	93313.38(0.65)	119827.88(0.84)	112346.49(0.79)	119655.16(0.84)
GITM	74638.54(0.51)	<b>10027.71(0.07)</b>	<b>17566.79(0.12)</b>	70071.99(0.48)	108816.53(0.74)	116526.40(0.80)	124821.83(0.85)
WALL-E w/o WM	72390.16(0.52)	15759.72(0.11)	29976.28(0.21)	58074.70(0.41)	124147.71(0.89)	102447.94(0.73)	103934.58(0.74)
<b>WALL-E (ours)</b>	<b>60348.71(0.41)</b>	23179.52(0.15)	36595.33(0.24)	<b>57106.20(0.39)</b>	<b>84776.25(0.58)</b>	<b>59261.31(0.40)</b>	<b>101173.64(0.68)</b>

#### 4.1 EXPERIMENTAL SETUP

**Benchmarks.** **Minecraft** is a popular open-world environment. We employ the standard evaluation pipeline provided by MineDojo’s TechTree tasks (Fan et al., 2022). These tasks can be categorized into six levels of increasing difficulty: *Wood*, *Stone*, *Iron*, *Gold*, *Diamond*, and *Redstone* (see Appendix E.1 for details). **ALFWorld** is a virtual environment designed as a text-based simulation where agents perform tasks by interacting with a simulated household environment (Shridhar et al., 2020b). This benchmark includes six distinct task types, each requiring the agent to accomplish a high-level objective, such as placing a cooled lettuce on a countertop (see Appendix E.2 for details).

**Metrics.** (1) **Success rate** (higher is better): the percentage of tasks the agent completes successfully. (2) **Replanning rounds** (lower is better): the number of times the agent revisits the same task to revise its plan for recovering from the failed task planning. (3) **Token cost** (lower is better): the number of tokens consumed by LLM agent/world

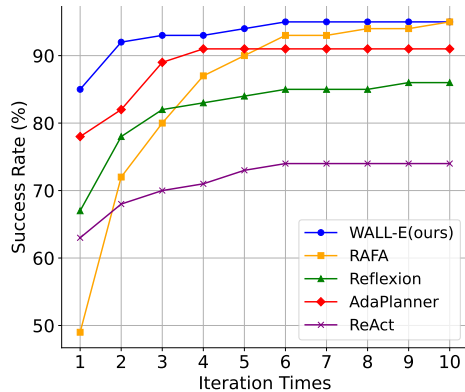


Figure 4: Comparison of WALL-E and baselines on 134 testing tasks from the ALFWorld benchmark.

by LLM agent/world

329 models during task completion. For Minecraft, we select four tasks from each level to serve as the testing  
330 set and the remaining tasks to construct the training set. All these three metrics are employed in our experi-  
331 ment. The task will be marked incomplete if the agent either dies in the environment (such as being killed by  
332 hostile mobs or falling into lava) or reaches one of the following maximal budgets: 10-minute time limit and  
333 maximum replanning rounds. In these cases, the replanning rounds and token cost will be set to the maximal  
334 value. For ALFWorld, we train WALL-E on the designated training set and evaluate its performance on a set  
335 of 134 predefined testing tasks. The averaged success rate over several trials is used as the evaluation metric  
336 to measure the performance of all baselines.

## 337 338 4.2 MAIN RESULTS

339  
340 We conduct a detailed comparison of WALL-E and existing baseline methods in Tables 1, 2 and 3, to  
341 demonstrate its superior performance in terms of success rate, planning efficiency, and token cost consumed  
342 by LLMs across diverse tasks.

343 **WALL-E demonstrates superior planning and task-solving abilities.** Tables 1 and 3 show that our  
344 method achieves the highest success rates across different environments. Specifically, in the Minecraft  
345 environment, WALL-E outperforms other baselines by an impressive margin of 15–30%. Figure 4 shows  
346 that WALL-E achieves the highest success rate after only 6 iterations, significantly surpassing other SOTA  
347 planning-based baselines such as RAFA (Hao et al., 2023) and AdaPlanner (Sun et al., 2024).

348 **Aligned world model leads to higher sample efficiency.** While the integration of the LLM world model  
349 leads to additional token costs compared to model-free methods, WALL-E demonstrates remarkably high  
350 sample efficiency, which is sufficient to offset the additional consumption caused by the world modeling.  
351 Specifically, our method requires 8–20 fewer replanning rounds than other baselines (see Table 1), resulting  
352 in overall token usage that is only 60–80% of that observed in other methods (see Table 2). It is worth noting  
353 that the advantage of WALL-E becomes more apparent in harder environments. In turn, model-free methods  
354 can only achieve comparatively high sample efficiency on those easy tasks such as *Wood* and *Stone*.

355 **WALL-E is a general and environment-agnostic method.** Unlike methods tailored to specific environ-  
356 ments, e.g., GITM (Zhu et al., 2023b) for open-world exploration in Minecraft and BUTLER (Micheli &  
357 Fleuret, 2021) for long-horizon planning in ALFWorld, WALL-E can excel at both, underscoring its gen-  
358 eralizability and effectiveness in enhancing agent’s capabilities of exploration, planning, and reflection in  
359 general, complex scenarios.

## 360 361 4.3 EFFECTIVENESS OF RULE LEARNING

362  
363 In order to demonstrate the effectiveness of our proposed rule learning method, we conduct a comparative  
364 study against GITM (Zhu et al., 2023b) - a method employing buffered trajectories as in-context examples to  
365 align LLM agents with the environment dynamics. By jointly examining the rule learning process (Figure 5)  
366 and the agent’s training progress (Figure 6), we observe an interesting phenomenon that WALL-E’s success  
367 rate hits the upper bound after 4 iterations, while the rule learning process also finds a compact set of rules  
368 for the LLM world model and keeps this set fixed after 4 iterations, reflecting that WALL-E’s improvement  
369 mainly benefits from the learning of new rules.

370 **Rule learning achieves efficient “world alignment”.** To verify whether the learned rules enable a more ac-  
371 curate world model, we first collect a dataset of transitions that cannot be predicted by the LLM world model  
372 correctly and evaluate each rule on this dataset by calculating the cover rate - the probability that the LLM’s  
373 failed predictions are addressed by the rules obtained during the rule learning process. According to Figure  
374 5, it is evident that the rules learned by our proposed framework consistently improve cover rates across  
375 different types of actions in the Minedojo benchmark. In specific, actions such as gather and fight reach



Table 3: Comparison of WALL-E and baselines on 134 testing tasks from the ALFWorld benchmark. \*-reported in previous work. VLMs = vision-language models, LLMs = large language models. The success rate (%) is the percentage of tasks completed successfully. The best score for each task is highlighted in **bold**.

Method	Success Rate (%) $\uparrow$						
	Avg.	Pick	Clean	Heat	Cool	Examine	Picktwo
<b>VLMs</b>							
MiniGPT-4* (Zhu et al., 2023a)	16	4	0	19	17	67	6
BLIP-2* (Li et al., 2023)	4	0	6	4	11	6	0
LLaMA-Adapter* (Gao et al., 2023b)	13	17	10	27	22	0	0
InstructBLIP* (Dai et al., 2023)	22	50	26	23	6	17	0
EMMA* (Yang et al., 2024)	<b>82</b>	<b>71</b>	<b>94</b>	<b>85</b>	<b>83</b>	<b>88</b>	<b>67</b>
<b>LLMs</b>							
BUTLER* (Micheli & Fleuret, 2021)	26	31	41	60	27	12	29
GPT-BUTLER* (Micheli & Fleuret, 2021)	69	62	81	85	78	50	47
DEPS (Wang et al., 2023a)	76	93	50	80	<b>100</b>	<b>100</b>	0
AutoGen* (Wu et al., 2023)	77	92	74	78	86	83	41
ReAct (Yao et al., 2023)	74	79	54	96	85	83	51
AdaPlanner (Sun et al., 2024)	91	<b>100</b>	<b>100</b>	89	<b>100</b>	97	47
Reflexion (Shinn et al., 2024)	86	92	94	70	81	90	88
RAFA (Liu et al., 2023)	<b>95</b>	<b>100</b>	97	91	95	<b>100</b>	82
<b>WALL-E (ours)</b>	<b>95</b>	<b>100</b>	97	<b>100</b>	86	85	<b>100</b>
<b>Human Performance*</b> (Shridhar et al., 2020a)	91	-	-	-	-	-	-

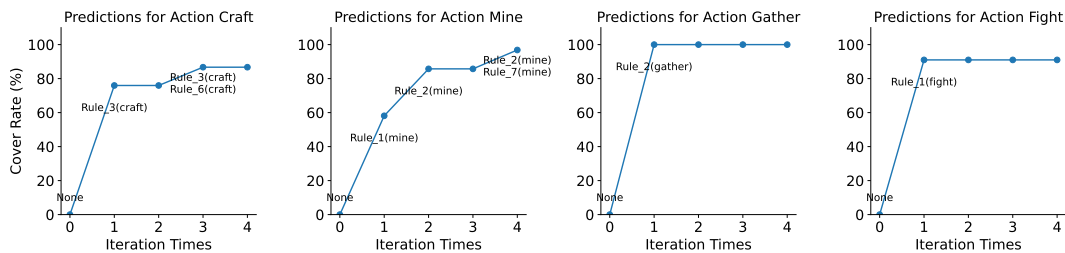


Figure 5: Cover rate of LLM failed predictions across different actions over iteration times during training. The cover rate represents the probability that the LLM’s failed predictions are addressed by the rules obtained during the rule learning process. The predictions and rules are categorized by action type: *craft*, *mine*, *gather* and *fight*. The learnt rules at each iteration are displayed in black under each node, labeled with their respective rule IDs.

100% and 91% coverage after the first iteration, while craft and mine actions demonstrate improvements over multiple iterations, with final coverage rates of 87% and 96%, respectively.

#### 4.4 ABLATION STUDY

We conduct a comprehensive ablation study to evaluate the importance of various components in WALL-E. Specifically, we separately remove the learned rules and the world model and check their effects on WALL-E’s final performance. **According to the results in Table 4, we give the following conclusions.** (1) Regardless of whether the learned rules are applied within the agent or the world model, adding them significantly enhances the total performance. The success rate increases by 20% to 30% approximately. This observation underscores the crucial role that rules play in improving the effectiveness of WALL-E. (2) When the learned rules are utilized within the world model, they contribute to nearly a 30% improvement in success rate, whereas using rules within the agent result in about a 20% improvement. This disparity may be primarily due to the fact that the learned rules are highly related to the state information (See Appendix

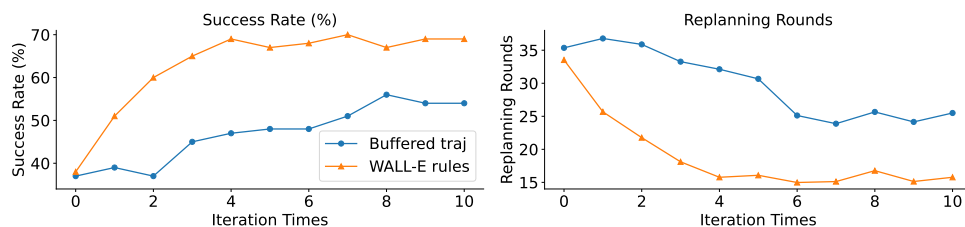


Figure 6: Learning curve comparison between rule learning (e.g., WALL-E) and buffered trajectory (e.g., GITM) over 10 iterations on Minecraft tasks during training. The left plot shows the average success rate (%) across all tasks, where a higher value indicates more tasks successfully completed. The right plot illustrates the average number of replanning rounds, with fewer rounds indicating higher efficiency in task completion.

Table 4: Ablation study of WALL-E with different configurations on Minecraft tasks, in the format of “success rate (replanning rounds)”. The success rate (%) refers to the percentage of tasks completed successfully (higher the better). Replanning rounds (lower the better) measure the inference efficiency and represent the number of revisions needed for the agent to complete a task. The row highlighted in grey represents the configuration and performance of WALL-E.

WALL-E		Success Rate (%) ↑ (Replanning Rounds ↓)						
Agent	World Model	Avg.	Wooden	Stone	Iron	Golden	Diamond	Redstone
LLM	-	37(35.36)	83(10.67)	41(33.26)	33(35.27)	22(45.29)	24(42.46)	17(45.22)
LLM	LLM	38(33.53)	86(10.35)	44(30.79)	35(34.08)	19(43.99)	26(39.51)	19(42.46)
LLM+rules	-	61(23.13)	94(5.04)	89(9.58)	67(18.56)	33(39.67)	41(32.73)	43(33.21)
LLM	LLM+rules	69(15.77)	98(1.64)	91(4.58)	63(19.38)	69(15.61)	46(27.08)	48(26.33)
LLM+rules	LLM+rules	67(16.59)	95(2.88)	93(3.75)	58(21.42)	62(19.34)	53(23.75)	43(28.41)

D for more details). (3) MPC using a world model without applying any rules cannot significantly improve WALL-E’s performance in terms of the success rate and the number of replanning times. This finding suggests that the alignment between the world model and the environment dynamics by rule learning is crucial to our appealing results.

## 5 CONCLUSION

We have shown that LLMs can effectively serve as world models for agents when aligned with environment dynamics through rule learning. Our neurosymbolic approach bridges the gap between LLMs’ prior knowledge and specific environments without gradient updates. By integrating a rule-enhanced LLM-based world model with MPC, our agent WALL-E demonstrates superior planning and task-solving abilities. Experiments indicate that WALL-E outperforms baselines in Minecraft and ALFWorld, achieving higher success rates with fewer replanning rounds and reduced token usage. Specifically, WALL-E attains a 15–30% higher success rate in Minecraft, requires 8–20 fewer replanning rounds, and uses only 60–80% of the tokens compared to baselines. In ALFWorld, it rapidly reaches a 95% success rate from the 6th iteration onward. The rule learning converges swiftly by the 4th iteration, outperforming buffered trajectory methods in both efficiency and effectiveness. These results suggest that minimal additional rules suffice to align LLM predictions with environment dynamics, enhancing model-based agents in complex environments.

## REFERENCES

- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. [arXiv preprint arXiv:2212.06817](#), 2022.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. InstructBLIP: Towards general-purpose vision-language models with instruction tuning. In [NeurIPS](#), 2023.
- Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied question answering. In [CVPR](#), 2018.
- Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. [NeurIPS](#), 2022.
- Jensen Gao, Bidipta Sarkar, Fei Xia, Ted Xiao, Jiajun Wu, Brian Ichter, Anirudha Majumdar, and Dorsa Sadigh. Physically grounded vision-language models for robotic manipulation. [arXiv preprint arXiv:2309.02561](#), 2023a.
- Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, et al. Llama-adapter v2: Parameter-efficient visual instruction model. [arXiv preprint arXiv:2304.15010](#), 2023b.
- Zhiqi Ge, Hongzhe Huang, Mingze Zhou, Juncheng Li, Guoming Wang, Siliang Tang, and Yueting Zhuang. Worldgpt: Empowering llm as multimodal world model. [arXiv preprint arXiv:2404.18202](#), 2024.
- Lin Guan, Karthik Valmeekam, Sarath Sreedharan, and Subbarao Kambhampati. Leveraging pre-trained large language models to construct and utilize world models for model-based task planning. [NeurIPS](#), 2023.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. [arXiv preprint arXiv:1912.01603](#), 2019.
- Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. [arXiv preprint arXiv:2010.02193](#), 2020.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. [arXiv preprint arXiv:2301.04104](#), 2023.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. Reasoning with language model is planning with world model. [arXiv preprint arXiv:2305.14992](#), 2023.
- Zhiting Hu and Tianmin Shu. Language models, agent models, and world models: The law for machine reasoning and planning. [arXiv preprint arXiv:2312.05230](#), 2023.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In [ICML](#), 2023.
- Ming Li, Han Chen, Chenguang Wang, Dang Nguyen, Dianqi Li, and Tianyi Zhou. Ruler: Improving llm controllability by rule-based data recycling. [arXiv preprint arXiv:2406.15938](#), 2024a.
- Zaijing Li, Yuquan Xie, Rui Shao, Gongwei Chen, Dongmei Jiang, and Liqiang Nie. Optimus-1: Hybrid multimodal memory empowered agents excel in long-horizon tasks. [arXiv preprint arXiv:2408.03615](#), 2024b.

- 517 Yang Liu, Weixing Chen, Yongjie Bai, Jingzhou Luo, Xinshuai Song, Kaixuan Jiang, Zhida Li, Ganlong  
518 Zhao, Junyi Lin, Guanbin Li, et al. Aligning cyber space with physical world: A comprehensive survey  
519 on embodied ai. [arXiv preprint arXiv:2407.06886](#), 2024.
- 520  
521 Zhihan Liu, Hao Hu, Shenao Zhang, Hongyi Guo, Shuqi Ke, Boyi Liu, and Zhaoran Wang. Reason for  
522 future, act for now: A principled framework for autonomous llm agents with provable sample efficiency.  
523 [arXiv preprint arXiv:2309.17382](#), 2023.
- 524 Linhao Luo, Jiabin Ju, Bo Xiong, Yuan-Fang Li, Gholamreza Haffari, and Shirui Pan. Chatrule: Mining log-  
525 ical rules with large language models for knowledge graph reasoning. [arXiv preprint arXiv:2309.01538](#),  
526 2023.
- 527  
528 Vincent Micheli and François Fleuret. Language models are few-shot butlers. [arXiv preprint](#)  
529 [arXiv:2104.07972](#), 2021.
- 530 Norman Mu, Sarah Chen, Zifan Wang, Sizhe Chen, David Karamardian, Lulwa Aljeraisy, Dan Hendrycks,  
531 and David Wagner. Can llms follow simple rules? [arXiv preprint arXiv:2311.04235](#), 2023a.
- 532  
533 Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang, Mingyu Ding, Jun Jin, Bin Wang, Jifeng Dai,  
534 Yu Qiao, and Ping Luo. Embodiedgpt: Vision-language pre-training via embodied chain of thought.  
535 [arXiv preprint arXiv:2305.15021](#), 2023b.
- 536  
537 OpenAI. GPT-4 technical report. [arXiv preprint arXiv.2303.08774](#), 2023.
- 538 S Joe Qin and Thomas A Badgwell. A survey of industrial model predictive control technology. [Control](#)  
539 [engineering practice](#), 11(7):733–764, 2003.
- 540  
541 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn.  
542 Direct preference optimization: Your language model is secretly a reward model. [NeurIPS](#), 2024.
- 543 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Lan-  
544 guage agents with verbal reinforcement learning. [NeurIPS](#), 2024.
- 545  
546 Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke  
547 Zettlemoyer, and Dieter Fox. ALFRED: A benchmark for interpreting grounded instructions for everyday  
548 tasks. In [CVPR](#), 2020a.
- 549 Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew  
550 Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. [arXiv preprint](#)  
551 [arXiv:2010.03768](#), 2020b.
- 552  
553 Haotian Sun, Yuchen Zhuang, Lingkai Kong, Bo Dai, and Chao Zhang. Adaplanner: Adaptive planning  
554 from feedback with language models. [NeurIPS](#), 2024.
- 555  
556 Hao Tang, Darren Key, and Kevin Ellis. Worldcoder, a model-based llm agent: Building world models by  
557 writing code and interacting with the environment. [arXiv preprint arXiv:2402.12275](#), 2024.
- 558  
559 Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, Yitao Liang, and Team CraftJarvis.  
560 Describe, explain, plan and select: interactive planning with large language models enables open-world  
561 multi-task agents. In [NeurIPS](#), 2023a.
- 562  
563 Zihao Wang, Shaofei Cai, Anji Liu, Yonggang Jin, Jinbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng He,  
Zilong Zheng, Yaodong Yang, et al. Jarvis-1: Open-world multi-task agents with memory-augmented  
multimodal language models. [arXiv preprint arXiv:2311.05997](#), 2023b.

- 564 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,  
565 and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In NeurIPS,  
566 2022.
- 567 Lionel Wong, Jiayuan Mao, Pratyusha Sharma, Zachary S Siegel, Jiahai Feng, Noa Korneev, Joshua B  
568 Tenenbaum, and Jacob Andreas. Learning adaptive planning representations with natural language guid-  
569 ance. arXiv preprint arXiv:2312.08566, 2023.
- 570 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang,  
571 Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation  
572 framework. arXiv preprint arXiv:2308.08155, 2023.
- 573 Xiyang Wu, Tianrui Guan, Dianqi Li, Shuaiyi Huang, Xiaoyu Liu, Xijun Wang, Ruiqi Xian, Abhinav Shri-  
574 vastava, Furong Huang, Jordan Lee Boyd-Graber, et al. Autohallusion: Automatic generation of halluci-  
575 nation benchmarks for vision-language models. arXiv preprint arXiv:2406.10900, 2024.
- 576 Jiannan Xiang, Tianhua Tao, Yi Gu, Tianmin Shu, Zirui Wang, Zichao Yang, and Zhiting Hu. Language  
577 models meet world models: Embodied experiences enhance language models. NeurIPS, 2024.
- 578 Wei Xiao, Noushin Mehdipour, Anne Collin, Amitai Y Bin-Nun, Emilio Frazzoli, Radboud Duintjer  
579 Tebbens, and Calin Belta. Rule-based optimal control for autonomous driving. In Proceedings of the  
580 ACM/IEEE 12th International Conference on Cyber-Physical Systems, pp. 143–154, 2021.
- 581 Kaige Xie, Ian Yang, John Gunerli, and Mark Riedl. Making large language models into world models with  
582 precondition and effect knowledge. arXiv preprint arXiv:2409.12278, 2024.
- 583 Wenkai Yang, Yankai Lin, Jie Zhou, and Jirong Wen. Enabling large language models to learn from rules.  
584 arXiv preprint arXiv:2311.08883, 2023a.
- 585 Yijun Yang, Tianyi Zhou, Kanxue Li, Dapeng Tao, Lusong Li, Li Shen, Xiaodong He, Jing Jiang, and Yuhui  
586 Shi. Embodied multi-modal agent trained by an llm from a parallel textworld. In CVPR, 2024.
- 587 Zeyuan Yang, Peng Li, and Yang Liu. Failures pave the way: Enhancing large language models through  
588 tuning-free rule accumulation. arXiv preprint arXiv:2310.15746, 2023b.
- 589 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React:  
590 Synergizing reasoning and acting in language models. In ICLR, 2023.
- 591 Zirui Zhao, Wee Sun Lee, and David Hsu. Large language models as commonsense knowledge for large-  
592 scale task planning. NeurIPS, 2024.
- 593 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-  
594 language understanding with advanced large language models. arXiv preprint arXiv:2304.10592, 2023a.
- 595 Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu,  
596 Xiaogang Wang, et al. Ghost in the minecraft: Generally capable agents for open-world environments via  
597 large language models with text-based knowledge and memory. arXiv preprint arXiv:2305.17144, 2023b.
- 598 Zhaocheng Zhu, Yuan Xue, Xinyun Chen, Denny Zhou, Jian Tang, D. Schuurmans, and Hanjun Dai. Large  
599 language models can learn rules. arXiv preprint arXiv:2310.07064, 2023c.
- 600 Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan  
601 Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic  
602 control. In CoRL, 2023.
- 603  
604  
605  
606  
607  
608  
609  
610

## A DETAILED RELATED WORK

**LLMs with Rule Learning.** Recent studies have explored integrating LLMs with rule learning to enhance reasoning and model behavior. For instance, Yang et al. (2023a) introduced rule distillation, enabling LLMs to learn from predefined rules, which improved generalization with limited training data. Similarly, Zhu et al. (2023c) proposed the Hypotheses-to-Theories (HtT) framework, which enhanced numerical and relational reasoning by generating and validating rules from training data. In the same vein, Mu et al. (2023a) developed the RuLES framework to evaluate LLM adherence to developer-specified rules, addressing challenges like rule evasion through adversarial inputs. Furthermore, Yang et al. (2023b) presented the Tuning-free Rule Accumulation (TRAN) framework, allowing LLMs to accumulate rules from incorrect cases to avoid repeating mistakes without additional tuning. Lastly, in knowledge graph reasoning, Luo et al. (2023) introduced ChatRule, a framework that mines logical rules over knowledge graphs using LLMs.

These studies show the potential of combining LLMs with rule learning to improve reasoning and generalization. However, none have integrated rule learning with LLM-based world models, which is the focus of our work. We explore how rule learning can align LLM world models with specific environment dynamics, thereby improving the performance of model-based agents in dynamic environments.

**Using LLMs to Build World Models.** Many studies have leveraged LLMs to construct world models for planning. For example, Wong et al. (2023) proposed translating natural language instructions into adaptable planning representations via LLMs, enabling flexible and context-aware world modeling. Similarly, Guan et al. (2023) showed that combining pre-trained LLMs with task-specific planning modules improves task success rates by providing a more detailed understanding of the environment. Another approach, World-Coder Tang et al. (2024), exemplified an LLM agent that constructs world models by generating and executing code to simulate various states and actions, refining its understanding iteratively.

These studies demonstrate the utility of LLMs in building world models to improve planning and reasoning in complex environments. However, unlike these works, our approach directly employs the LLM as the world model, utilizing its inherent knowledge and reasoning abilities without an explicit model-building phase. This direct use of LLMs enhances adaptability and computational efficiency.

**Using LLMs as World Models.** Several studies have explored using LLMs directly as world models by leveraging their implicit knowledge. Some methods rely on fine-tuning to align the LLM world model with the environment. For example, Xiang et al. (2024) fine-tuned LLMs with embodied experiences in a simulated world to enhance reasoning and planning abilities in embodied environments. Similarly, Xie et al. (2024) transformed LLMs into world models by incorporating knowledge of action preconditions and effects, fine-tuning the models to reason about actions and predict their outcomes accurately.

Other approaches align LLMs as world models through prompting. For instance, Zhao et al. (2024) introduced the LLM-MCTS algorithm, prompting LLMs to serve as both the policy and world model for large-scale task planning, integrating commonsense priors with guided search. In another approach, Hao et al. (2023) introduced Reasoning via Planning (RAP), where LLMs are prompted to act as reasoning agents and world models by generating reasoning trees to explore solutions. Finally, (Liu et al., 2023) used a Bayesian adaptive Markov Decision Process to guide LLMs in planning future trajectories, prompting them to predict future states.

While these approaches demonstrate the potential of using LLMs as world models, they often require extensive fine-tuning or rely heavily on human-crafted prompts, making them labor-intensive and inflexible. Our work overcomes these limitations by automatically extracting rules from exploration experiences, reducing human effort and enhancing adaptability across different environments.

## B DETAILED PROMPT

### B.1 LEARN NEW RULES FROM REAL TRAJECTORIES

#### Prompt for Learning New Rules from Real Trajectories

```

658 You are responsible for mining new rules from the given transitions, ensuring
659 that these rules differ from the ones already provided.
660 Focus on generating general and universal rules that are not tied to any
661 specific item or tool.
662 Your goal is to generalize across different objects, creating flexible rules
663 that can be applied broadly to diverse contexts and situations.
664
665 I will give you an array of transitions:
666 [
667   {
668     'state_0': {
669       "state feature 1": {"feature name": value, ...},
670       ...
671     },
672     'action': {
673       "name": "action name",
674       "action feature 1": {"feature name": value, ...},
675       ...
676     },
677     'action_result': {
678       "feedback": "the environment feedback",
679       "success": "Whether the action is executed successfully, give 'True' or
680         'False' only",
681       "suggestion": "If the 'action' fails, 'suggestion' would be given based
682         on 'state 0' and 'action'"
683     }
684   },
685   {
686     'state_0': {
687       "state feature 1": {"feature name": value, ...},
688       ...
689     },
690     'action': {
691       "name": "action name",
692       "action feature 1": {"feature name": value, ...},
693       ...
694     },
695     'action_result': {
696       "feedback": "the environment feedback",
697       "success": "Whether the action is executed successfully, give 'True' or
698         'False' only",
699       "suggestion": "If the 'action' fails, 'suggestion' would be given based
700         on 'state 0' and 'action'"
701     }
702   },
703   ...
704 ]
and an array of rules:
[

```

```

705     "Rule 1: For action ..., if..., the action will fail; Checking Method:
706     ...",
707     "Rule 2: For action ..., if..., the action will fail; Checking Method:
708     ...",
709     ...
710 ]
711 You should only respond in the format as described below:
712 RESPONSE FORMAT:
713 {
714     "new_rules":[
715         "Rule ...: For action ...,...; Checking Method: ...",
716         "Rule ...: For action ...,...; Checking Method: ...",
717         ...
718     ]
719 }
720 Instructions:
721 - Ensure the response can be parsed by Python `json.loads`, e.g.: no trailing
722   commas, **no single quotes**, etc.
723 - Please use you knowledge in <ENV>, do inductive reasoning. You need to dig up
724   as many rules as possible that satisfy all transitions.
725 - Extract and utilize only the features that influence the outcome of the
726   action.
727 - Please generate general and universal rules; the rules should not reference
728   any specific item or tool! You need to generalize across various items or
729   tools.
730 - Generate only the rules under what conditions the action will fail.
731 - While generating a rule, you also need to state how to check if a transition
732   satisfies this rule. Please be specific as to which and how 'features' need
733   to be checked

```

## B.2 REFINE LEARNED RULES

### Prompt for Refining Learned Rules

```

737 You are responsible for improving the existing rules by verifying that they
738 hold true for all transitions.
739 This involves identifying any conflicting rules, diagnosing potential issues,
740 and making necessary modifications.
741 Ensure that the refined rules are consistent and correctly align with the
742 transitions provided, avoiding any contradictions or overlaps.
743 I will give you an array of transitions:
744 [
745     {
746         'state_0': {
747             "state feature 1": {"feature name": value, ...},
748             ...
749         },
750         'action': {
751             "name": "action name",
752             "action feature 1": {"feature name": value, ...},
753             ...

```



```

752     },
753     'action_result': {
754       "feedback": "the environment feedback",
755       "success": "Whether the action is executed successfully, give 'True' or
756         'False' only",
757       "suggestion": "If the 'action' fails, 'suggestion' would be given based
758         on 'state 0' and 'action'"
759     },
760   {
761     'state_0': {
762       "state feature 1": {"feature name": value, ...},
763       ...
764     },
765     'action': {
766       "name": "action name",
767       "action feature 1": {"feature name": value, ...},
768       ...
769     },
770     'action_result': {
771       "feedback": "the environment feedback",
772       "success": "Whether the action is executed successfully, give 'True' or
773         'False' only",
774       "suggestion": "If the 'action' fails, 'suggestion' would be given based
775         on 'state 0' and 'action'"
776     }
777   }
778 ]
779 and an array of rules:
780 [
781   "Rule 1: For action ..., if..., the action will fail; Checking Method:
782   ...",
783   "Rule 2: For action ..., if..., the action will fail; Checking Method:
784   ...",
785   ...
786 ]
787
788 You should only respond in the format as described below:
789 RESPONSE FORMAT:
790 {
791   "verified_rules":[
792     "Rule ...: For action ...,...; Checking Method: ...",
793     "Rule ...: For action ...,...; Checking Method: ...",
794     ...
795   ],
796   "conflicting_rules":[
797     "Rule ...: For action ...,...; Checking Method: ...",
798     "Rule ...: For action ...,...; Checking Method: ...",
799     ...
800   ],
801   "improved_rules":[
802     "Rule ...: For action ...,...; Checking Method: ...",
803     "Rule ...: For action ...,...; Checking Method: ...",
804     ...
805   ]
806 }

```

```

799     ],
800     "final_rules":[
801         "Rule ...: For action ...,...; Checking Method: ...",
802         "Rule ...: For action ...,...; Checking Method: ...",
803         ...
804     ]
805 }
806
807 where
808 verified_rules: list rules that satisfy all the provided transitions.
809 conflicting_rules: list rules that contradict any of the transitions. Modify
810     these rules if they can be modified correctly and put them in '
811     improved_rules'.
812 improved_rules: show modified 'conflicting_rules'.
813 final_rules: combine all the rules from 'verified_rules', 'new_rules'.
814
815 Instructions:
816 - Ensure the response can be parsed by Python 'json.loads', e.g.: no trailing
817     commas, **no single quotes**, etc.
818 - Please use you knowledge in <ENV>, do inductive reasoning. You need to dig up
819     as many rules as possible that satisfy all transitions.
820 - Extract and utilize only the features that influence the outcome of the
821     action.
822 - Please generate general and universal rules; the rules should not reference
823     any specific item or tool! You need to generalize across various items or
824     tools.
825 - Generate only the rules under what conditions the action will fail.
826 - While generating a rule, you also need to state how to check if a transition
827     satisfies this rule. Please be specific as to which and how 'features' need
828     to be checked

```

### B.3 TRANSLATE NATURAL LANGUAGE RULES TO CODE

#### Prompt for Translating Natural Language Rules to Code

```

832 You are responsible for generating code rules by implementing the learned rules
833     in Python.
834 Your task is to write a function that takes the current state and an action as
835     inputs, evaluates these conditions, and returns a Boolean value based on
836     the specified rule.
837 This function should effectively mirror the logic of the rules, enabling
838     precise predictions for various state-action pairs.
839
840 The function should be defined as follows:
841
842 ```python
843 def expected_rule_code(state, action):
844     # Your code here
845     return feedback, success, suggestion
846 where
847 feedback: a string, give the action feedback based on success or not.

```

```

846 success: a bool, whether the action is executed successfully, give 'True' or '
847 False'. If the action type is not the action type in the rule, count as
848 success (e.g., success = True).
849 suggestion: a string, if the 'action' fails, 'suggestion' would be given based
850 on 'rule', 'state' and 'action'.
851
852 Here is several examples of the input format:
853 <Input Format>
854
855 The function should return a Boolean (True or False) based on an internal rule
856 which you must implement.
857
858 Ensure that the function handles the input and outputs the expected result
859 based on <ENV>'s mechanics and the provided state and action.
860
861 If the rule involves the need to use your knowledge to make a judgement about
862 an item or action then write the function, LLM_request("question"+"response
863 format").
864 LLM_request would send the "question" to gpt4, and return the gpt4's response.
865 you just need to write the "question" in the LLM_request.
866 LLM_request("question"+"response format") has already been predefined, you can
867 just use it dirtectly. Do not need to define it again in your response. But
868 you need to define the "question" and "response format" carefully.
869
870 example: i want to know if the item can be destroyed
871 the LLM function: LLM_request(f"if the {item} can be destroyed in the <ENV>?" +
872 "only reply True or False")
873
874 You should only respond in the format as described below, and do not give
875 example usage or anything else:
876 RESPONSE FORMAT:
877 def expected_rule_code(state, action):
878     # Your code here
879

```

where “input format” please refer to Appendix C.

## C ENVIRONMENTS’ STATE SPACE AND ACTION SPACE

The format of state and action information is crucial for understanding the rules we have extracted. In this section, we provide an description of the state and action space used in different environments.

### C.1 MINECRAFT

**State Space.** We collect state information directly from the observation space provided by MineDojo (Fan et al., 2022), which includes: (1) equipment status, (2) inventory details, (3) life statistics, and (4) location statistics. The specific structure is illustrated in the following example.

#### Examples for Minecraft’s State Space

```

890 state = {
891     "equipment": {
892         "dirt": 60.0,

```

```

893     "diamond boots": 1.0,
894     "diamond leggings": 1.0,
895     "diamond chestplate": 1.0,
896     "diamond helmet": 1.0,
897     "air": 0.0
898   },
899   "inventory": {
900     "dirt": 60.0,
901     "crafting table": 1.0,
902     "planks": 2.0,
903     "stick": 4.0,
904     "air": 0.0,
905     "log": 1.0
906   },
907   "life_stats": {
908     "life": 20.0,
909     "oxygen": 300.0,
910     "armor": 20.0,
911     "food": 20.0,
912     "saturation": 5.0,
913     "is_sleeping": False
914   },
915   "location_stats": {
916     "biome": "plains",
917     "rainfall": 0.4,
918     "temperature": 0.8,
919     "is_raining": False,
920     "sky_light_level": 0.2,
921     "sun_brightness": 0.0
922   }
923 }

```

**Action Space.** The action space is defined based on the action API provided by MineDojo (Fan et al., 2022), with additional modifications inspired by the action space used in GITM (Zhu et al., 2023b). The detailed action definitions are presented below.

### Minecraft’s Action Space

```

924 craft(obj, materials, platform): craft the object with the materials and
925 platform; used to craft new object that is not in the inventory or is not
926 enough.
927 - obj: a dict, whose key is the name of the object and value is the object
928 quantity, like {"crafting table": 1} and {"stone pickaxe": 1}.
929 - materials: a dict, whose keys are the names of the materials and values are
930 the quantities, like {"planks": 4} and {"cobblestone": 3, "stick": 2}.
931 - platform: a string, the platform used for crafting the current 'object', like
932 "furnace" and "crafting table". Set to null if without any platform.
933
934 mine(obj, tool, y_level): dig down to the y-level and mine the specified object
935 with the tool. This action will go underground and continuously mine the
936 object until the desired quantity is obtained.
937 - obj: a dict, whose key is the name of the object and value is the object
938 quantity, like {"stone": 5} and {"iron ore": 1}.
939 - tool (string): the tool used for mining, like "wooden pickaxe". Set to null
  if without any tool.

```

```

940 - y_level: a number, the y-level to dig down to. Different ores have different
941     probabilities of distribution in different levels.
942
943 fight(obj, target, tool): find, track, and fight the target until you collect
944     the desired number (goal_num) of object by using the chosen tool.
945 - obj: a dict, whose key is the name of the object and value is the object
946     quantity, like {"leather": 5} and {"porkchop": 3}.
947 - target: a string, The name of the entity you want to fight (e.g., "skeleton",
948     "sheep").
949 - tool: a string, the tool or weapon you will use in the fight, like "iron
950     sword" or "wooden sword". Set to null if without any tool.
951
952 equip(obj): equip the object from the inventory.
953 - obj: a string, the object to equip, like "wooden pickaxe".
954
955 apply(obj, target, tool): automates the process of using a tool on target until
956     you collect a specific number of object.
957 - obj: a dict, whose key is the name of the object and value is the object
958     quantity, like {"wool": 5}.
959 - target: a string, the name of the target you want to interact with (e.g., "
960     water", "sheep").
961 - tool: a string, the specific tool you will use for the action. (e.g., "bucket
962     ", "shears")
963
964 gather(obj, tool): collect resources (obj) directly from the environment. This
965     includes picking up flowers, seeds from grass, and wood from trees.
966 - obj: a dict, whose key is the name of the object and value is the object
967     quantity, like {"log": 10}.
968 - tool: a string, the tool you will use in the gathering. Set to null if
969     without any tool.
970
971 change_time(target_time): adjust to the specified time of day; this function
972     enables you to wait until a predefined time, such as morning, night, or
973     midnight, depending on the specified target_time.
974 - target_time: a string, specifying the desired time to change to. Valid
975     options include "morning", "night", and "midnight", each corresponding to
976     distinct values in 'sky_light_level' and 'sun_brightness' in "state
977     features" like:
978 -- "morning": 'sky_light_level': array([1.]), 'sun_brightness': array([1.])
979 -- "night": 'sky_light_level': array([0.25]), 'sun_brightness': array([0.36])
980 -- "midnight": 'sky_light_level': array([0.2]), 'sun_brightness': array([0.])

```

## 976 C.2 ALFWORLD

978 **State Space.** In the original ALFWorld setup, state information is represented as natural language dialogue  
979 history. To facilitate the rule learning process, we developed scripts to transform this dialogue history into a  
980 structured JSON format, as shown in the following example.

### 981 Examples for ALFWorld's State Space

```

983 state = {
984     "reachable_locations": [
985         "cabinet 5",
986         "cabinet 4",

```

```

987     "cabinet 3",
988     "cabinet 2",
989     "cabinet 1",
990     "coffeemachine 1",
991     "countertop 2",
992     "countertop 1",
993     "diningtable 1",
994     "drawer 2",
995     "drawer 1",
996     "fridge 1",
997     "garbagecan 1",
998     "microwave 1",
999     "shelf 3",
1000    "shelf 2",
1001    "shelf 1",
1002    "sinkbasin 1",
1003    "stoveburner 4",
1004    "stoveburner 3",
1005    "stoveburner 2",
1006    "stoveburner 1",
1007    "toaster 1"
1008 ],
1009 "items_in_locations": {
1010   "fridge 1": [
1011     "lettuce 2",
1012     "mug 2",
1013     "potato 3"
1014   ],
1015   "microwave 1": []
1016 },
1017 "item_in_hand": {
1018   "item_name": "cup 1",
1019   "status": "normal"
1020 },
1021 "current_position": {
1022   "location_name": "microwave 1",
1023   "status": "open"
1024 }
1025 }

```

**Action Space.** We utilize the action space provided by the ALFWorld directly, as demonstrated below.

### Action Space for Minecraft

```

1025 go to [location/object]: Move to a specified location or object.
1026 open [object]: Open a specified object like a cabinet or drawer.
1027 close [object]: Close an opened object.
1028 take [object] from [location]: Pick up an item from a specified location.
1029 put [object] in/on [location]: Place an item in or on a specified location.
1030 clean [object] with [location/tool]: Clean an object using a specific location
1031 or tool, like cleaning lettuce at the sink basin.
1032 heat [object] with [tool]: Use an appliance, such as a microwave, to heat an
1033 item.
1034 cool [object] with [tool]: Use a cooling tool or appliance, such as a fridge,
1035 to cool an item.

```

1034 use [tool]: Activate or use a tool, such as a desk lamp.  
 1035  
 1036

## 1037 D LEARNED RULES 1038

1039 There are two points to note about the numbering of the rules:  
 1040

- 1041 • The reason for duplicates is that the numbering is based on actions, and different actions have their  
 1042 own separate sequences. For example: Rules for Craft: [Rule 1, Rule 2, Rule 3, Rule 4, Rule 5...];  
 1043 Rules for Mine: [Rule 1, Rule 2, Rule 3, Rule 4, Rule 5...].
- 1044 • The reason the sequence may appear unordered is that some rules have been pruned (Section 3.2  
 1045 Rule Set Pruning via Maximum Coverage). For instance, Rules for Craft where [Rule 1, Rule 2,  
 1046 Rule 4, Rule 5] has been removed, Rules for Mine where [Rule 1, Rule 3, Rule 4, Rule 5, Rule 6]  
 1047 has been removed, and the final rule set is Rules for Craft: [Rule 3, Rule 6] and Rules for Mine:  
 1048 [Rule 2, Rule 7].  
 1049

### 1050 D.1 NATURAL LANGUAGE RULES 1051

#### 1052 Natural Language Rules for Minecraft

1053 "Rule 3: For action 'craft', if the specified platform is incorrect or not  
 1054 specified when required, the action will fail; Checking Method: Check if  
 1055 the 'platform' specified in the 'action' matches the required platform for  
 1056 the 'obj' being crafted.",  
 1057 "Rule 6: For action 'craft', if the player does not have enough materials to  
 1058 craft the specified object, the action will fail; Checking Method: Check if  
 1059 the 'materials' specified in the 'action' are present in the 'inventory'  
 1060 with the required quantities. If not, the action will fail.",  
 1061 "Rule 2: For action 'mine', if the 'tool' is not appropriate for the object  
 1062 being mined, the action will fail; Checking Method: Check if 'action.args.  
 1063 tool' is not suitable for 'action.args.obj'.",  
 1064 "Rule 7: For action 'mine', if the 'tool' is not in the inventory, the action  
 1065 will fail; Checking Method: Check if 'action.args.tool' is not present in '  
 1066 state\_0.inventory'.",  
 1067 "Rule 2: For action 'gather', if the 'sky\_light\_level' in 'location\_stats' is  
 1068 less than 1.0, the action will fail; Checking Method: Check if '  
 1069 sky\_light\_level' in 'location\_stats' is less than 1.0.",  
 1070 "Rule 1: For action 'fight', if the 'tool' is not present in the 'inventory' or  
 1071 'equipment', the action will fail; Checking Method: Check if the 'tool'  
 1072 specified in the action is present in either 'inventory' or 'equipment'."

#### 1073 Natural Language Rules for ALFWorld

1074 Rule 1: For action 'clean', if the object to be cleaned is not in hand, the  
 1075 action will fail; Checking Method: Check if 'item\_in\_hand.item\_name' in '  
 1076 initial\_state' matches 'action.args.obj'.  
 1077 Rule 3: For action clean, if the tool is not reachable, the action will fail;  
 1078 Checking Method: Check if the tool specified in the action is in the list  
 1079 of reachable locations in the initial state.  
 1080 Rule 5: For action 'clean', if the current position is not at the tool location  
 , the action will fail; Checking Method: Check if 'current\_position.  
 location\_name' in 'initial\_state' matches 'action.args.tool'.

```

1081 Rule 2: For action 'take', if the agent is already holding an item, the action
1082 will fail; Checking Method: Check if 'item_in_hand.item_name' in '
1083 initial_state' is not None.
1084 Rule 4: For action 'take', if the agent is not at the location of the item, the
1085 action will fail; Checking Method: Check if the 'current_position.
1086 location_name' in 'initial_state' is not the same as the 'source' in the '
1087 action'.
1088 Rule 3: For action 'open', if the current position is not the target, the
1089 action will fail; Checking Method: Check if the 'current_position' is the
1090 target.
1091 Rule 2: For action 'put', if the item to be put is not in hand, the action will
1092 fail; Checking Method: Check if 'item_in_hand.item_name' is not equal to '
1093 action.args.obj'.
1094 Rule 1: For action 'use', if the object to be used is not at the current
1095 position, the action will fail; Checking Method: Check if the object
1096 specified in the action is listed under the 'items_in_locations' of the '
1097 current_position' in the 'initial_state'.
1098 Rule 1: For action 'heat', if the tool (microwave) is not at the current
1099 position, the action will fail; Checking Method: Check if 'current_position
1100 .location_name' is equal to the tool in the action arguments.
1101 Rule 5: For action 'heat', if the item in hand is not the item to be heated,
1102 the action will fail; Checking Method: Check if 'item_in_hand' in '
1103 initial_state' is equal to 'action.args.obj'.
1104 Rule 1: For action 'go to', if the target location is the same as the current
1105 location, the action will fail; Checking Method: Check if 'current_position
1106 .location_name' is equal to 'action.args.target'.

```

## 1105 D.2 CODE-BASED RULES

1107 When a rule requires the LLM’s domain knowledge to make judgments, we instruct the LLM to use the function `LLM_request('question', 'response format')` directly within the generated code. The LLM should generate the "question" and "response format" according to the function to be implemented. The predefined `LLM_request` function sends the message to the LLM and returns its response, enabling the code to dynamically leverage the LLM’s knowledge.

1112 Additionally, the feedback and suggestions returned by each code-based rule are automatically generated by prompting the LLM with the corresponding rule. The detailed prompts used to generate these code-based rules can be found in Appendix B.3. These feedback and suggestions play a crucial role in helping the agent refine and improve its planning process (Section 3.3).

### 1117 Code-based Rules for Minecraft

```

1118 def Rule_3_craft(state, action):
1119     if action['name'] == 'craft':
1120         obj = list(action['args']['obj'].keys())[0]
1121         platform = action['args']['platform']
1122
1123         # Ask the LLM if the specified platform is required for the object
1124         being crafted
1125         question = f"Is a specific platform required to craft {obj} in
1126         Minecraft? If yes, what is the platform?"
1127         response_format = "only reply with the platform name (e.g., 'crafting
         table', 'furnace') or 'None' if no specific platform is required"

```



```

1128     required_platform = LLM_request(question + response_format)
1129
1130     if required_platform != 'None' and platform != required_platform.lower
1131         ():
1132         feedback = f"Crafting {obj} requires a {required_platform}, but {
1133             platform} was provided."
1134         success = False
1135         suggestion = f"Use a {required_platform} to craft {obj}."
1136         return feedback, success, suggestion
1137     else:
1138     feedback = f"Crafting {obj} was successful."
1139     success = True
1140     suggestion = ""
1141     return feedback, success, suggestion
1142
1143 else:
1144     feedback = "Action type is not 'craft', so it is considered successful
1145     ."
1146     success = True
1147     suggestion = ""
1148     return feedback, success, suggestion
1149
1150 def Rule_6_craft(state, action):
1151     feedback = ""
1152     success = True
1153     suggestion = ""
1154
1155     if action["name"] == "craft":
1156         materials_needed = action["args"]["materials"]
1157         inventory = state["inventory"]
1158
1159         for material, quantity in materials_needed.items():
1160             if inventory.get(material, 0) < quantity:
1161                 feedback = f"Failed to craft {list(action['args']['obj'].keys()
1162                     ) [0]} due to insufficient {material}."
1163                 success = False
1164                 suggestion = f"Gather more {material} to craft {list(action['
1165                     args']['obj'].keys()) [0]}."
1166                 break
1167             else:
1168                 feedback = f"Successfully crafted {list(action['args']['obj'].keys
1169                     ()) [0]}."
1170
1171         return feedback, success, suggestion
1172
1173 def Rule_2_mine(state, action):
1174     feedback = ""
1175     success = True
1176     suggestion = ""
1177
1178     if action["name"] == "mine":
1179         obj = list(action["args"]["obj"].keys()) [0]
1180         tool = action["args"]["tool"]
1181
1182         # Check if the tool is appropriate for the object being mined

```

```
1175     question = f"Is the tool '{tool}' appropriate for mining '{obj}' in
1176     Minecraft? Only reply True or False."
1177     is_tool_appropriate = LLM_request(question)
1178
1179     if is_tool_appropriate == "False":
1180         feedback = f"The tool '{tool}' is not appropriate for mining '{obj}
1181         }'."
1182         print(feedback)
1183         success = False
1184         suggestion = f"Use an appropriate tool for mining '{obj}'."
1185     else:
1186         feedback = f"The tool '{tool}' is appropriate for mining '{obj}'."
1187         print(feedback)
1188         success = True
1189
1190     return feedback, success, suggestion
1191
1192 def Rule_7_mine(state, action):
1193     feedback = ""
1194     success = True
1195     suggestion = ""
1196
1197     if action["name"] == "mine":
1198         tool = action["args"]["tool"]
1199         if tool and tool not in state["inventory"]:
1200             feedback = f"Action failed: Tool '{tool}' is not in the inventory."
1201             success = False
1202             suggestion = f"Please ensure you have the '{tool}' in your
1203             inventory before mining."
1204         else:
1205             feedback = "Action succeeded: Tool is present in the inventory."
1206             success = True
1207             suggestion = ""
1208
1209     return feedback, success, suggestion
1210
1211 def Rule_2_gather(state, action):
1212     feedback = ""
1213     success = True
1214     suggestion = ""
1215
1216     if action["name"] == "gather":
1217         sky_light_level = state["location_stats"]["sky_light_level"][0]
1218         if sky_light_level < 1.0:
1219             feedback = "Action failed: sky light level is less than 1.0."
1220             success = False
1221             suggestion = "Wait until the sky light level is higher."
1222         else:
1223             feedback = "Action succeeded."
1224             success = True
1225             suggestion = ""
1226     else:
1227         feedback = "Action succeeded."
1228         success = True
1229         suggestion = ""
1230
1231     return feedback, success, suggestion
```

```

1222
1223 def Rule_1_fight(state, action):
1224     # Extract action name and arguments
1225     action_name = action.get("name")
1226     action_args = action.get("args", {})
1227
1228     # Initialize feedback, success, and suggestion
1229     feedback = ""
1230     success = True
1231     suggestion = ""
1232
1233     # Rule 1: For action 'fight', check if the 'tool' is present in 'inventory'
1234     # or 'equipment'
1235     if action_name == "fight":
1236         tool = action_args.get("tool")
1237         if tool:
1238             inventory = state.get("inventory", {})
1239             equipment = state.get("equipment", {})
1240             if tool not in inventory and tool not in equipment:
1241                 feedback = f"Action '{action_name}' failed: Tool '{tool}' is
1242                 not present in inventory or equipment."
1243                 success = False
1244                 suggestion = f"Ensure the tool '{tool}' is available in either
1245                 inventory or equipment before attempting to fight."
1246             else:
1247                 feedback = f"Action '{action_name}' succeeded: Tool '{tool}' is
1248                 available."
1249         else:
1250             feedback = f"Action '{action_name}' failed: No tool specified."
1251             success = False
1252             suggestion = "Specify a tool to use for the fight action."
1253     else:
1254         feedback = f"Action '{action_name}' is not restricted by the rule."
1255     return feedback, success, suggestion

```

### Code-based Rules for ALFWorld

```

1254
1255 def Rule_1_clean(state, action):
1256     if action['name'] == 'clean':
1257         obj_to_clean = action['args']['obj']
1258         item_in_hand = state['item_in_hand']['item_name']
1259         if obj_to_clean != item_in_hand:
1260             feedback = f"Action failed: {obj_to_clean} is not in hand."
1261             success = False
1262             suggestion = f"Please take {obj_to_clean} in hand before cleaning."
1263             return feedback, success, suggestion
1264         feedback = "Action executed successfully."
1265         success = True
1266         suggestion = ""
1267         return feedback, success, suggestion
1268
1269 def Rule_3_clean(state, action):
1270     if action["name"] == "clean":
1271         tool = action["args"]["tool"]

```

```

1269         if tool not in state["reachable_locations"]:
1270             feedback = f"Action failed: The tool '{tool}' is not reachable."
1271             success = False
1272             suggestion = f"Make sure the tool '{tool}' is in the list of
1273                 reachable locations."
1274             return feedback, success, suggestion
1275         else:
1276             feedback = "Action succeeded: The tool is reachable."
1277             success = True
1278             suggestion = ""
1279             return feedback, success, suggestion
1280     else:
1281         feedback = "Action succeeded: The action type is not 'clean'."
1282         success = True
1283         suggestion = ""
1284         return feedback, success, suggestion
1285
1286 def Rule_5_clean(state, action):
1287     if action['name'] == 'clean':
1288         current_position = state['current_position']['location_name']
1289         tool_location = action['args']['tool']
1290         if current_position != tool_location:
1291             feedback = f"Action 'clean' failed: You are not at the tool
1292                 location ({tool_location})."
1293             success = False
1294             suggestion = f"Move to the tool location ({tool_location}) before
1295                 cleaning."
1296             return feedback, success, suggestion
1297         # If the action is not 'clean' or the rule conditions are met
1298         feedback = "Action executed successfully."
1299         success = True
1300         suggestion = ""
1301         return feedback, success, suggestion
1302
1303 def Rule_2_take(state, action):
1304     feedback = ""
1305     success = True
1306     suggestion = ""
1307     if action["name"] == "take":
1308         if state["item_in_hand"]["item_name"] is not None:
1309             feedback = "Action failed: Agent is already holding an item."
1310             success = False
1311             suggestion = "You may heat, put, cool the item in hand directly
1312                 without removing the other items in target location/container."
1313         else:
1314             feedback = "Action succeeded: Agent is not holding any item."
1315             success = True
1316             suggestion = ""
1317     else:
1318         feedback = "Action succeeded: Action type is not 'take'."
1319         success = True
1320         suggestion = ""
1321     return feedback, success, suggestion
1322
1323 def Rule_4_take(state, action):

```

```

1316     if action['name'] == 'take':
1317         current_location = state['current_position']['location_name']
1318         source_location = action['args']['source']
1319         if current_location != source_location:
1320             feedback = "Action failed: Agent is not at the location of the item
1321             ."
1322             success = False
1323             suggestion = f"Move to {source_location} before taking the item."
1324             return feedback, success, suggestion
1325         # If the action is not 'take', it is considered successful
1326         feedback = "Action executed successfully."
1327         success = True
1328         suggestion = ""
1329         return feedback, success, suggestion
1330
1331 def Rule_3_open(state, action):
1332     if action['name'] == 'open':
1333         target = action['args']['target']
1334         current_position = state['current_position']['location_name']
1335
1336         if current_position != target:
1337             feedback = f"Action 'open' failed: You are not at the target
1338             location '{target}'."
1339             success = False
1340             suggestion = f"Move to '{target}' before trying to open it."
1341             return feedback, success, suggestion
1342         else:
1343             feedback = f"Action 'open' succeeded: You are at the target
1344             location '{target}'."
1345             success = True
1346             suggestion = ""
1347             return feedback, success, suggestion
1348         else:
1349             feedback = "Action succeeded: The action type is not 'open'."
1350             success = True
1351             suggestion = ""
1352             return feedback, success, suggestion
1353
1354 def Rule_2_put(state, action):
1355     if action['name'] == 'put':
1356         item_in_hand = state['item_in_hand']['item_name']
1357         item_to_put = action['args']['obj']
1358         if item_in_hand != item_to_put:
1359             feedback = f"Action failed: The item '{item_to_put}' is not in hand
1360             ."
1361             success = False
1362             suggestion = f"Please ensure you have '{item_to_put}' in hand
1363             before attempting to put it."
1364             return feedback, success, suggestion
1365         # If the action is not 'put', it is considered successful
1366         feedback = "Action executed successfully."
1367         success = True
1368         suggestion = ""
1369         return feedback, success, suggestion
1370

```

```

1363 def Rule_1_use(state, action):
1364     if action['name'] == 'use':
1365         obj = action['args']['obj']
1366         current_location = state['current_position']['location_name']
1367         # Check if the object is in the current location
1368         if obj not in state['items_in_locations'].get(current_location, []):
1369             feedback = f"Action failed: {obj} is not at the current position {
1370                 current_location}."
1371             success = False
1372             suggestion = f"Move to the location where {obj} is present or bring
1373                 {obj} to the current location."
1374             return feedback, success, suggestion
1375         # If the action is not 'use', it is considered successful
1376         feedback = "Action executed successfully."
1377         success = True
1378         suggestion = ""
1379         return feedback, success, suggestion
1380
1381 def Rule_1_heat(state, action):
1382     feedback = ""
1383     success = True
1384     suggestion = ""
1385     if action["name"] == "heat":
1386         tool = action["args"]["tool"]
1387         current_position = state["current_position"]["location_name"]
1388         if current_position != tool:
1389             feedback = f"Action failed: The tool '{tool}' is not at the current
1390                 position '{current_position}'."
1391             success = False
1392             suggestion = f"Move to the location of the tool '{tool}' before
1393                 attempting to heat."
1394         else:
1395             feedback = "Action succeeded: The tool is at the current position."
1396             success = True
1397             suggestion = ""
1398     else:
1399         feedback = "Action succeeded: The action type is not 'heat'."
1400         success = True
1401         suggestion = ""
1402     return feedback, success, suggestion
1403
1404 def Rule_5_heat(state, action):
1405     if action["name"] == "heat":
1406         item_in_hand = state["item_in_hand"]["item_name"]
1407         item_to_heat = action["args"]["obj"]
1408
1409         if item_in_hand != item_to_heat:
1410             feedback = f"Action failed: You are trying to heat {item_to_heat}
1411                 but you are holding {item_in_hand}."
1412             success = False
1413             suggestion = f"Hold {item_to_heat} before trying to heat it."
1414             return feedback, success, suggestion
1415
1416     feedback = "Action executed successfully."
1417     success = True

```

Table 5: Techtree Task Details

Task Level	Tasks
Wooden	wooden sword, wooden pickaxe, wooden axe, wooden hoe, wooden shovel
Stone	stone sword, stone pickaxe, stone axe, stone hoe, stone shovel
Iron	iron sword, iron pickaxe, iron axe, iron hoe iron shovel, iron boots, iron chestplate, iron helmet, iron leggings
Golden	golden sword, golden pickaxe, golden axe, golden hoe golden shovel, golden boots, golden chestplate, golden helmet, golden leggings
Diamond	diamond sword, diamond pickaxe, diamond axe, diamond hoe diamond shovel, diamond boots, diamond chestplate, diamond helmet, diamond leggings
Redstone	redstone block, redstone clock, redstone compass, redstone dispenser, redstone dropper redstone piston, redstone torch, redstone repeater, redstone detector rail, redstone activator rail

```

1424     suggestion = ""
1425     return feedback, success, suggestion
1426
1427 def Rule_1_go_to(state, action):
1428     if action['name'] == 'go_to':
1429         current_location = state['current_position']['location_name']
1430         target_location = action['args']['target']
1431
1432         if current_location == target_location:
1433             feedback = f"Action failed: You are already at {target_location}."
1434             success = False
1435             suggestion = "Try moving to a different location."
1436             return feedback, success, suggestion
1437
1438         # If the action is not 'go_to' or the target location is different from the
1439         # current location
1440         feedback = "Action executed successfully."
1441         success = True
1442         suggestion = ""
1443         return feedback, success, suggestion

```

## E EXPERIMENT DETAILS

### E.1 MINECRAFT

**Task Details.** We used the 'Tech Tree' series tasks in MineDojo. Minecraft presents a structured progression system involving various levels of tools and armor, each with unique properties and increasing difficulty to unlock. To advance through these levels, the agent must develop and apply systematic, compositional skills to navigate the technology tree. Tasks are structured into six technology tiers: *wood*, *stone*, *iron*, *gold*, *diamond*, and *redstone* with each level presenting a higher degree of difficulty. And each level contains a certain number of tasks, which is shown in Table 5. Additionally, we only do tasks in overworld, so tasks that require materials from Nether and End to complete are disregarded (e.g. redstone observer, redstone lamp, redstone comparator).

**Method Setup.** We utilize GPT-4o as the backend for our method. To rigorously assess the agent's performance, we initialize it in the "open-ended" mode—the most challenging and interactive environment

available, analogous to survival mode. In this setting, the agent starts with an empty inventory and randomized seeds for both the environment and its starting position, requiring it to strategize effectively. Unlike creative or adventure modes, the agent must contend with the dynamic generation of hostile mobs, introducing additional complexity and difficulty. Starting without any resources, the agent is forced to actively mine materials and craft essential items to progress, testing its planning, adaptability, and problem-solving skills.

We select four tasks from each task level to serve as the testing set and the remaining tasks to construct the training set. For the level with a limited number of tasks, such as *Wood* and *Stone*, we add additional tasks from Optimus-1 (Li et al., 2024b) to ensure sufficient diversity for rule learning in the training process. Finally, we have a total of 30 training tasks and 24 testing tasks.

Within MPC framework, reward is assigned as follows: a reward of +1 if the world model  $f_{wm}$  predicts the transition will be successful (`action_result = True`), and 0 if it predicts failure (`action_result = False`). The world model provides feedback to the agent, enabling the agent to refine its plan based on the state prior to the failed action and the received feedback. This iterative process continues until the task is successfully completed within the planning phase.

For buffered trajectories (e.g., GITM (Zhu et al., 2023b)), we adopted the original settings by storing successful task trajectories. During planning, we search this buffer for the trajectory most similar to the current task and include it in the prompt as a reference.

## E.2 ALFWORLD

**Task Details.** ALFWorld is a virtual environment designed as a text-based simulation where agents perform tasks by interacting with a simulated household. The environment includes six distinct task types, each requiring the agent to accomplish a high-level objective, such as placing a cooled lettuce on a countertop. Agents use text commands to navigate and manipulate objects in the virtual space, for example, issuing instructions like "go to countertop 1," "take lettuce 1 from countertop 1," or "cool lettuce 1 with fridge 1." The visual observations from the agent’s point of view are converted into natural language descriptions before being delivered to the agent. The agent’s state is represented by the cumulative history of these observations. Success is measured by the completion the specified task goal.

**Method Setup.** We conducted rule learning on the training set, with the resulting rules presented in Appendix D. Since tasks in ALFWorld require agents to continuously gather information from the environment, and our learned rules focus on capturing the dynamic of the environment, we adopted a one-step MPC. This method evaluates whether the agent’s current action aligns with the environment’s dynamic patterns based on its state information. Additionally, to enhance rule discovery, we developed scripts to convert the natural language dialogue history and action information into a structured JSON format, as illustrated in Appendix C.2. We utilize GPT-3.5-Instruct as our backbone model.

## E.3 EXPERIMENT DESIGN FOR EFFECTIVENESS OF RULE LEARNING

We conduct 3 training tasks per iteration over a total of 10 iterations during training. After each iteration, the model, equipped with latest learned rules or buffered trajectories, is tested on the testing set.

The cover rate quantifies the extent to which the rules derived from the rule learning process address the LLM’s failed predictions. Specifically, it represents the probability that mispredicted transitions by the LLM are correctly handled by the learned rules.

To assess the alignment between the LLM-based world model and the actual environment, we first identify transitions where the LLM fails to make accurate predictions. This is achieved by utilizing an unaligned LLM world model—one without any rules—to generate predictions for trajectories obtained from the test set. The discrepancies between the predicted states  $\hat{s}_{t+1}$  and the actual states  $s_{t+1}$  are compiled into a dataset



of mispredicted transitions. These mispredictions highlight areas where the LLM world model does not align with the environment’s dynamics.

Subsequently, the learned rules at each iteration are evaluated against the mispredicted transitions dataset to determine their effectiveness in correcting these mispredictions. If a rule successfully predicts the outcome of a previously mispredicted transition, it demonstrates that the rule effectively addresses the LLM’s failure in that instance. The cover rate is then calculated as the ratio of correctly addressed mispredictions to the total number of mispredicted transitions:

$$\text{Cover Rate} = \frac{\text{Number of Mispredictions Addressed by Rules}}{\text{Total Number of Mispredicted Transitions}} \quad (9)$$

Furthermore, as depicted in Figure 5, predictions and rules are categorized by action types—*craft*, *mine*, *gather*, and *fight*—allowing the cover rate to be calculated for each action category individually. A higher cover rate indicates that the rule learning process effectively enhances the alignment of the LLM world model with the environment, thereby improving the overall accuracy and reliability of the agent’s planning.

## F GREEDY ALGORITHM

We implement the following Algorithm 1 to solve the maximum set cover problem 8.

---

### Algorithm 1 Greedy Algorithm for Maximum Set Cover Problem

---

```

1: Input:
2:  $\mathcal{D}^{\text{incorrect}} = \{\delta_1^{\text{incorrect}}, \delta_2^{\text{incorrect}}, \dots, \delta_n^{\text{incorrect}}\}$  ▷ Set of incorrect transitions to cover
3:  $R^{\text{code}} = \{R_1^{\text{code}}, R_2^{\text{code}}, \dots, R_m^{\text{code}}\}$  ▷ Set of rules covering subsets of transitions
4:  $a_{ij}$ : Indicator matrix where  $a_{ij} = 1$  if  $\delta_j^{\text{incorrect}} \in R_i^{\text{code}}$ , otherwise  $a_{ij} = 0$ 
5: Output: Set of selected rules  $R_{\text{selected}}$ 
6: Initialize  $R_{\text{selected}} \leftarrow \emptyset$ 
7: Initialize  $\mathcal{D}_{\text{covered}} \leftarrow \emptyset$  ▷ Set of covered transitions
8: Initialize  $x_i \leftarrow 0$  for all  $i \in \{1, \dots, m\}$  ▷ Rule selection indicators
9: Initialize  $y_j \leftarrow 0$  for all  $j \in \{1, \dots, n\}$  ▷ Transition coverage indicators
10: while  $\mathcal{D}_{\text{covered}} \neq \mathcal{D}^{\text{incorrect}}$  do
11:   For each rule  $R_i^{\text{code}} \in R^{\text{code}}$ , compute:
12:      $\text{gain}(R_i) = |(\{\delta_j^{\text{incorrect}} \mid a_{ij} = 1\} \setminus \mathcal{D}_{\text{covered}})|$ 
13:   Select the rule  $R_{i^*}^{\text{code}}$  with the largest gain, i.e.,
14:      $i^* = \arg \max_i \text{gain}(R_i)$ 
15:   if  $\max \text{gain}(R_i) = 0$  then ▷ Terminate if no rule can cover any additional transitions
16:     Break
17:   end if
18:   Add  $R_{i^*}^{\text{code}}$  to  $R_{\text{selected}}$ 
19:   Update  $\mathcal{D}_{\text{covered}} \leftarrow \mathcal{D}_{\text{covered}} \cup \{\delta_j^{\text{incorrect}} \mid a_{i^*j} = 1\}$ 
20:   Set  $x_{i^*} \leftarrow 1$  ▷ Mark rule  $R_{i^*}^{\text{code}}$  as selected
21:   For each  $\delta_j^{\text{incorrect}}$  covered by  $R_{i^*}^{\text{code}}$ , set  $y_j \leftarrow 1$ 
22: end while
23: Return  $R_{\text{selected}}$ 

```

---

## G CODE-BASED RULES VERIFICATION LOGIC

**Criteria for Determining Whether a Rule is Active.** We apply rules only in scenarios where they are relevant to the current transition, defining this relevance as the rule being "active" or "effective." A rule is considered active if its outcome of the current transition yields the specific outcome it is designed to address. Specifically, in our framework, transition success is represented as `True`, and failure is represented as `False`.

- **Rules Designed to Identify Successes:**  
A rule intended to detect successes is considered active when it evaluates the current transition and returns `True`.
- **Rules Designed to Identify Failures:**  
A rule intended to detect failures is considered active when it evaluates the current transition and returns `False`.

In essence, a rule is active when its outcome aligns with the type of outcome it is meant to assess (either success or failure). This ensures that rules are applied appropriately and only influence the LLM world model's predictions when relevant to the specific circumstances of the transition.

**Determining Whether a Rule is Correct or Incorrect** When a rule is active, if it makes an incorrect judgment—predicting success when the transition actually fails or vice versa—the rule is considered invalid and is removed from the rule set. Transitions where the rule is not applicable—referred to as "inactive" or "dormant"—are excluded from the evaluation process.

## H LIMITATION AND FUTURE WORK

Currently, our rule learning framework generates simple rules that primarily assess whether actions align with environment dynamics (i.e., rules for transitions). Future research should explore advanced reasoning methods that enable LLMs to derive more abstract rules, such as those governing entire planning processes. Furthermore, many embodied environments exhibit stochastic dynamics, where actions have probabilistic outcomes. For example, resource gathering at night in Minecraft often fails due to hostile creatures but can sometimes succeed. Our current rule learning process cannot handle such randomness, typically classifying these scenarios as failures. Addressing this limitation by enabling rules to account for stochastic dynamics is a promising research direction, potentially leading to more accurate and reliable world models.