

Agent Planning with World Knowledge Model

Shuofei Qiao^{♣*}, Runnan Fang^{♣*}, Ningyu Zhang^{♣†}, Yuqi Zhu[♣], Xiang Chen[♣],
Shumin Deng[♣], Yong Jiang[◇], Pengjun Xie[◇], Fei Huang[◇], Huajun Chen^{♣†}

[♣]Zhejiang University [♣]National University of Singapore, NUS-NCS Joint Lab [◇]Alibaba Group
[♡]Zhejiang Key Laboratory of Big Data Intelligent Computing
{shuofei, zhangningyu}@zju.edu.cn

Abstract

Recent endeavors towards directly using large language models (LLMs) as agent models to execute interactive planning tasks have shown commendable results. Despite their achievements, however, they still struggle with brainless trial-and-error in global planning and generating hallucinatory actions in local planning due to their poor understanding of the “real” physical world. Imitating humans’ *mental* world knowledge model which provides global prior knowledge before the task and maintains local dynamic knowledge during the task, in this paper, we introduce *parametric* **World Knowledge Model (WKM)** to facilitate agent planning. Concretely, we steer the agent model to self-synthesize knowledge from both expert and sampled trajectories. Then we develop WKM, providing prior *task knowledge* to guide the global planning and dynamic *state knowledge* to assist the local planning. Experimental results on three complex real-world simulated datasets with three state-of-the-art open-source LLMs, Mistral-7B, Gemma-7B, and Llama-3-8B, demonstrate that our method can achieve superior performance compared to various strong baselines. Other interesting findings include: 1) our instance-level task knowledge can generalize better to unseen tasks, 2) weak WKM can guide strong agent model planning, and 3) unified WKM training has promising potential for further development³.

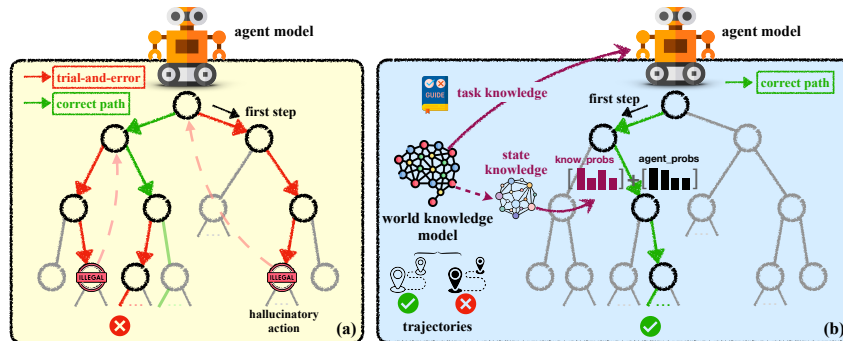


Figure 1: Traditional agent planning vs. Agent planning with world knowledge model.

1 Introduction

The remarkable advances in Large Language Models (LLMs) have witnessed a rapid development of various natural language processing tasks [25, 16, 28, 47, 60, 33]. Recently, multiple attempts that

* Equal Contribution.

† Corresponding Author.

³The code is available at <https://github.com/zjunlp/WKM>.

directly exploit LLMs as agent models to address physical world planning tasks have demonstrated promising achievements [54, 57, 56, 34, 38, 64, 44]. However, as most state-of-the-art LLMs are autoregressive models trained with next-token prediction, they lack the ability to essentially understand the real world, leading to generating hallucinatory actions and performing brainless trial-and-error in the environment as shown in Figure 1(a).

In contrast to LLMs, humans possess a mental knowledge model about the physical world [1, 18, 17, 30]. When facing a specific task, they will first briefly rehearse the entire process in mind using their rich prior knowledge before performing mindless actions. We call this kind of knowledge *global task knowledge* (a.k.a. environment/task commonsense). In addition, during the task procedure, the mental world knowledge model will constantly maintain a kind of local *state knowledge*, representing humans’ cognition of the current world state. For example, imagine you are in a room and your task is to put a clean egg in microwave. The *task knowledge* may refer to The egg is most likely in the fridge ... The workflows are: 1) locate and take the egg; 2) clean the egg using sinkbasin ... The *state knowledge* possibly refers to My task is to ... I have found and taked the egg ... Next I should ... The absence of world knowledge can lead to blind trial-and-error in the early planning stages when environmental information is limited. Conversely, in later stages when information is redundant, it can easily result in a confused cognition of the current world state and generate hallucinatory actions.

The process by which humans handle planning tasks reminds us to develop a parametric **World Knowledge Model (WKM)** to facilitate agent planning. As humans typically acquire knowledge from expertise and practical experience, we build WKM based on knowledge learned from both expert and explored trajectories. Specifically, we first steer the agent model to synthesize task knowledge from the comparison between expert and sampled trajectories. Then we prompt it to summarize state knowledge for each planning step from expert trajectories and combine the previous and next actions to build a state knowledge base. Lastly, we integrate the generated knowledge into expert trajectories and train a WKM. The agent model needs to be retrained to adapt to the task knowledge. Note our agent and knowledge model are both trained with LoRA [12] sharing the same backbone.

During the planning phase, we use the WKM to provide global prior task knowledge and maintain local dynamic state knowledge for the agent model as shown in Figure 1(b). The task knowledge will be concatenated in natural language form following the specific task to guide the agent model’s trial-and-error. At each planning step, to prevent the occurrence of hallucinatory actions, we utilize the generated state knowledge as the query to conduct k NN retrieval from the pre-built state knowledge base. We then use the constraints from the previous action, the probabilities of the retrieved next actions, and the probabilities from the agent model to make a weighted prediction for the next action.

We evaluate our method on three real-world simulated planning tasks: ALFWorld [41], WebShop [53], and ScienceWorld [50] with three state-of-the-art open-source LLMs: Mistral-7B [16], Gemma-7B [24], and Llama-3-8B [25]. Empirical results demonstrate that our method achieves superior performance compared to various strong baselines on both seen and unseen tasks. Moreover, further analytical results show that 1) our WKM can effectively reduce blind trial-and-error and hallucinatory actions, 2) our model-generated instance-level knowledge can generalize better to unseen tasks, 3) weak-guide-strong is feasible, 4) multi-task unified WKM possesses strong potential, and 5) explicit state knowledge will hurt the performance of agent planning.

2 Preliminaries

We mainly focus on interactive tasks with partial observations from environments. Following the task formulation in [44], the problem can be viewed as a Partially Observable Markov Decision Process (POMDP): $(\mathcal{U}, \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T})$. The instruction space \mathcal{U} defines the task and its corresponding regulations. \mathcal{S} is the state space, \mathcal{O} is the observation space, and \mathcal{A} is the action space. $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ defines the transition function, which we assume to be given by the environments. It is noticed that \mathcal{U} , \mathcal{A} , and \mathcal{O} are subspaces of the natural language space in the language agent scenarios.

Based on the above, the historical trajectory h_t that consists of a list of actions and observations at time t can be represented as:

$$h_t = (u, a_0, o_0, a_1, o_1, \dots, a_t, o_t), \tag{1}$$

where $u \in \mathcal{U}$ is the task instruction and $a \in \mathcal{A}$, $o \in \mathcal{O}$ are the action and the observation. Given a task, the language agent with parameter θ serves as the policy model π_θ responsible for generating

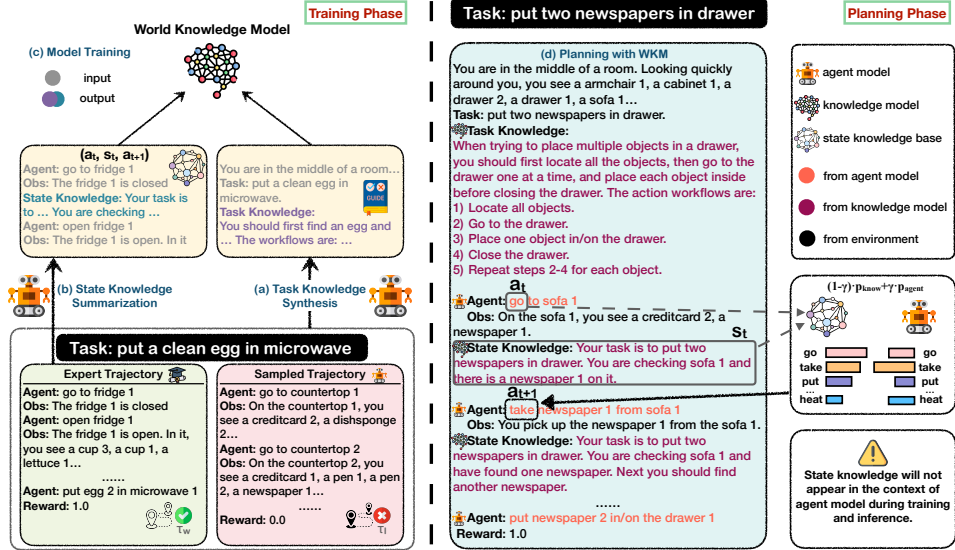


Figure 2: **Overview of our WKM.** We train a world knowledge model on the knowledge synthesized by the agent model itself from both expert and explored trajectories, providing prior task knowledge to guide global planning and dynamic state knowledge to assist local planning.

the action a_{t+1} based on h_t at each time step $t + 1$:

$$a_{t+1} \sim \pi_{\theta}(\cdot|h_t). \quad (2)$$

Specifically, $a_0 \sim \pi_{\theta}(\cdot|u)$ is generated according to the task instruction u . The whole trajectory τ concludes when the task is completed or exceeds the maximum time steps. Then the production of the entire trajectory with time length n can be modeled as:

$$\pi_{\theta}(\tau|u) = \prod_{t=0}^n \pi_{\theta}(a_{t+1}|h_t)\pi_{\theta}(a_0|u). \quad (3)$$

Ultimately, the final reward $r(u, \tau) \in [0, 1]$ representing the task completion rate is calculated. Note that we follow a REACT-style [54] trajectory that includes rationales before each action. We use a to represent the action with rationales for convenience.

World Knowledge Model. *World knowledge model* serves as humans’ mental cognition of the physical environment, more intricate than the *word knowledge model* which LLM-powered agent models are trained to be [61, 10, 52, 13]. Our “world” here refers to the simulated environment of the task. Based on the static environment of the task and the dynamic changes during interaction with the agent, we define world knowledge as a combination of prior global knowledge and dynamic local knowledge, corresponding to the blind trial-and-error problem in global planning and the hallucinatory action issue in local planning in traditional agent models, respectively. To attain precise and efficient agent planning, we develop a *parametric* WKM to simulate the *mental* WKM of humans.

3 Method

As shown in Figure 2, we steer the agent model to self-synthesize the *task knowledge* from the comparison of expert and sampled trajectories (§3.1). Then we prompt the agent model to self-summarize the *state knowledge* based on historical behavior and construct a state knowledge base (§3.2). The generated knowledge will be integrated into the expert trajectories for training the WKM. After the training process (§3.3), we augment the agent model with the world knowledge model to achieve effective and accurate planning (§3.4).

3.1 Task Knowledge Synthesis

The *task knowledge* serves as the prior knowledge to guide the agent model’s global planning and prevent it from dropping into blind trial-and-error.

Experienced Agent Exploration. We primarily acquire task knowledge through the comparison of preference trajectories (*chosen vs. rejected*). In order to improve the quality of rejected trajectories and obtain more targeted task knowledge, we employ an experienced agent for exploration. Firstly, we train a vanilla language model with expert trajectories⁴ from the training set to obtain an experienced agent. Subsequently, the experienced agent explores the training set tasks again to generate rejected trajectories. Our purpose is to extract superior task knowledge that cannot be acquired solely through supervised fine-tuning on chosen trajectories, thus further effectively boosting the agent’s capabilities.

Self Knowledge Synthesis. With the expert trajectories as the chosen ones and the trajectories sampled from the experienced agent as the rejected ones, we prompt the agent model itself to synthesize the task knowledge. Supposing \mathcal{K} is the task knowledge space:

$$\kappa \sim \pi_{\theta}(\cdot | \rho_{\text{TaskKnow}}, u, \tau_w, \tau_l), \quad (4)$$

where $\kappa \in \mathcal{K}$ is the task knowledge, ρ_{TaskKnow} stands for the prompt to instruct the task knowledge extraction, and τ_w, τ_l are the chosen and rejected trajectories respectively. Note that given the same task u , τ_w and τ_l always satisfy $r(u, \tau_w) = 1 \geq r(u, \tau_l)$. Even when $r(u, \tau_w) = r(u, \tau_l)$, we still consider trajectories sampled from the experienced agent as rejected ones. This is because expert trajectories often have shorter step lengths, enabling the agent to learn more knowledge of efficient planning. For detailed prompts of task knowledge synthesis, please refer to Appendix I.1.

3.2 State Knowledge Summarization

The *state knowledge* serves as the dynamic knowledge to constrain the agent model’s local planning and prevent it from generating hallucinatory actions. We prompt the agent model to self-summarize state knowledge at each planning step based on the expert trajectories to guarantee quality. For detailed prompts of state knowledge summarization, please refer to Appendix I.2. Supposing the prompt used to summarize state knowledge is $\rho_{\text{StateKnow}}$ and the state knowledge $s \in \mathcal{S}$ is a part of the state space \mathcal{S} , the generation of state knowledge at time t can be represented as:

$$s_t \sim \pi_{\theta}(\cdot | \rho_{\text{StateKnow}}, h_t). \quad (5)$$

State Knowledge Base Construction. To avoid confusion caused by excessive additional information, instead of explicitly concatenating the state knowledge to the context, we construct a state knowledge base for retrieval (we analyze in §4.3 how explicit state knowledge may affect the performance of agent model). We combine the state knowledge s_t with the previous action a_t and next action a_{t+1} from the expert trajectory to form a action-state-action triplet (a_t, s_t, a_{t+1}) . After iterating through all expert trajectories, we obtain a State Knowledge Base $\mathcal{B} = \{(s, a_{\text{pre}}, a_{\text{next}})^{(i)}\}_{i=1}^{|\mathcal{B}|}$, where $a_{\text{pre}} = a_t, a_{\text{next}} = a_{t+1}$, and $|\mathcal{B}|$ is the size of the state knowledge base.

3.3 Model Training

We integrate the generated world knowledge into expert trajectories and train a world knowledge model. The agent model needs to be re-trained to adapt to the incorporation of task knowledge. Note that our agent model and knowledge model are both trained with LoRA **sharing the same backbone**. We list the examples of training data for both the agent model and WKM in Appendix E.

Agent Model Training. Given the expert trajectories dataset $\mathcal{D} = \{(u, \kappa, \tau_w)^{(i)}\}_{i=1}^{|\mathcal{D}|}$ with **task knowledge** κ generated in §3.1, we train the agent model to follow the task knowledge to generate actions. Under an auto-regressive manner, the loss of the agent model can be formulated as:

$$\mathcal{L}_{\text{agent}}(\pi_{\theta}) = -\mathbb{E}_{\tau_w \sim \mathcal{D}}[\pi_{\theta}(\tau_w | u, \kappa)] \quad (6)$$

Suppose $\mathcal{X} = (x_1, x_2, \dots, x_{|\mathcal{X}|})$ is the token sequence of the trajectory τ_w , we have:

$$\pi_{\theta}(\tau_w | u, \kappa) = -\sum_{j=1}^{|\mathcal{X}|} (\mathbb{1}(x_j \in \mathcal{A}) \times \log \pi_{\theta}(x_j | u, \kappa, x_{<j})). \quad (7)$$

Here $\mathbb{1}(x_j \in \mathcal{A})$ is the indicator function to mask tokens unrelated to actions. Please note that τ_w here **does not include** the state knowledge mentioned in §3.2.

⁴For details on how to collect expert trajectories, please refer to Appendix A.

World Knowledge Model Training. The main difference in the training data between the agent and knowledge model is **the added state knowledge**. Given the expert trajectories dataset with both task and state knowledge $\mathcal{D}' = \{(u, \kappa, \tau'_w)^{(i)}\}_{i=1}^{|\mathcal{D}'|}$ where $\tau'_w = (a_0, o_0, s_0, \dots, a_n, o_n, s_n)$, the loss of the knowledge model π_ϕ can be formulated as:

$$\mathcal{L}_{\text{know}}(\pi_\phi) = -\mathbb{E}_{\kappa, \tau'_w \sim \mathcal{D}'} [\pi_\phi(\kappa|u) \pi_\phi(\tau'_w|u, \kappa)] \quad (8)$$

Suppose $\mathcal{X}' = (x'_1, x'_2, \dots, x'_{|\mathcal{X}'|})$ is the token sequence of the expert trajectory with state knowledge τ'_w and $\mathcal{Y} = (y_1, y_2, \dots, y_{|\mathcal{Y}|})$ represents the token sequence of the task knowledge κ , we have:

$$\pi_\phi(\kappa|u) = -\sum_{i=1}^{|\mathcal{Y}|} \log \pi_\phi(y_i|u, y_{<i}) \quad (9)$$

$$\pi_\phi(\tau'_w|u, \kappa) = -\sum_{j=1}^{|\mathcal{X}'|} (\mathbb{1}(x'_j \in \mathcal{S}) \times \log \pi_\phi(x'_j|u, \kappa, x'_{<j})), \quad (10)$$

where $\mathbb{1}(x_j \in \mathcal{S})$ is the indicator function to mask tokens unrelated to state knowledge.

3.4 Agent Planning with World Knowledge Model

At inference time, the agent model plans on the evaluation tasks with the aid of the world knowledge model. We redefine the historical trajectory $h_t = (u, \kappa, a_0, o_0, a_1, o_1, \dots, a_t, o_t)$. Given a specific task instruction u , the knowledge model first generates the task knowledge $\kappa \sim \pi_\phi(\cdot|u)$, then the agent model starts planning. Assuming the available action set $\mathcal{A}_u \subseteq \mathcal{A}$ for the task u is $(\alpha_u^{(1)}, \alpha_u^{(2)}, \dots, \alpha_u^{(|\mathcal{A}_u|)})$, at any time $t \geq 0$, instead of directly generating a next action $a_{t+1} \in \mathcal{A}_u$ based on h_t , we first employ the world knowledge model to generate the current state knowledge $s_t \sim \pi_\phi(\cdot|h_t)$ and leverage s_t to query the state knowledge base $\mathcal{B} = \{(s, a_{\text{pre}}, a_{\text{next}})^{(i)}\}_{i=1}^{|\mathcal{B}|}$. With the state knowledge as the $k \in \mathcal{Y}$, we retrieve \mathcal{N} nearest triplets **from where** $a_{\text{pre}} = a_t$ based on semantic similarity and collect the corresponding next actions a_{next} . We count the probability of each action $p_{\text{know}}(\alpha_u^{(i)}) = \frac{\mathcal{N}_i}{\mathcal{N}}$, where \mathcal{N}_i is the occurrence number of action $\alpha_u^{(i)}$ in all the collected a_{next} . Therefore, we get the probability acquired from the state knowledge base:

$$P_{\text{know}}(\mathcal{A}_u) = (p_{\text{know}}(\alpha_u^{(1)}), p_{\text{know}}(\alpha_u^{(2)}), \dots, p_{\text{know}}(\alpha_u^{(|\mathcal{A}_u|)})), \quad \sum_{i=1}^{|\mathcal{A}_u|} p_{\text{know}}(\alpha_u^{(i)}) = 1. \quad (11)$$

Afterward, we sample the probability distribution of the first token for each action $\alpha_u^{(i)}, 1 \leq i \leq |\mathcal{A}_u|$ from the last layer of the agent model and apply a `softmax` function to normalize the probability distribution. We define the probability acquired from the agent model as:

$$P_{\text{agent}}(\mathcal{A}_u) = (p_{\text{agent}}(\alpha_u^{(1)}), p_{\text{agent}}(\alpha_u^{(2)}), \dots, p_{\text{agent}}(\alpha_u^{(|\mathcal{A}_u|)})), \quad \sum_{i=1}^{|\mathcal{A}_u|} p_{\text{agent}}(\alpha_u^{(i)}) = 1. \quad (12)$$

Finally, we determine the next action by combining the above two probabilities:

$$a_{t+1} = \arg \max_{\alpha_u^{(i)} \in \mathcal{A}_u, 1 \leq i \leq |\mathcal{A}_u|} (\gamma \cdot p_{\text{agent}}(\alpha_u^{(i)}) + (1 - \gamma) \cdot p_{\text{know}}(\alpha_u^{(i)})), \quad (13)$$

where γ is the hyperparameter that controls the proportion of $P_{\text{agent}}(\mathcal{A}_u)$. Based on the above, we enhance the agent planning by global guidance from task knowledge and local constraints from state knowledge generated by our WKM. Due to the WKM and retrieval, the inference stage incurs additional time overhead compared to the pure agent model. The approximate ratio is around 2.5:1.

4 Experiments

4.1 Experimental Settings

Datasets and Metrics. We evaluate our method on three real-world simulated planning datasets: **ALFWorld** [41], **WebShop** [53], and **ScienceWorld** [50]. ALFWorld and ScienceWorld include

Table 1: **Main Results.** The best results are marked in **bold** and the second-best results are marked with underline. All the prompt-based baselines (\odot) are evaluated under one-shot prompting and all the fine-tuning-based baselines (\bullet) are trained through LoRA. **Red** represents the changes of WKM relative to the optimal results in the baselines. WKM and agent model are different LoRAs sharing the same backbone.

Backbone	Method	ALFWorld		WebShop	ScienceWorld	
		Seen	Unseen		Seen	Unseen
GPT-3.5-Turbo GPT-4	\odot REACT	8.57	5.97	44.37	15.41	13.99
		44.29	38.05	62.76	67.32	65.09
Mistral-7B	\odot REACT	7.86	5.22	14.63	20.72	17.65
	\odot Reflexion	11.56	6.00	16.64	21.07	18.11
	\bullet NAT	64.43	68.96	61.01	57.12	50.79
	\bullet ETO	66.84	<u>71.43</u>	<u>64.09</u>	58.17	<u>51.85</u>
	\bullet KNOWAGENT	<u>70.44</u>	70.72	61.28	<u>59.32</u>	47.24
	WKM	73.57 +3.13	76.87 +5.44	65.48 +1.39	62.12 +2.80	53.62 +1.77
Gemma-7B	\odot REACT	6.43	2.24	5.93	3.58	3.51
	\odot Reflexion	7.14	2.99	7.71	4.94	3.93
	\bullet NAT	67.86	65.88	55.82	47.63	44.98
	\bullet ETO	66.43	<u>68.66</u>	<u>62.67</u>	<u>50.44</u>	47.84
	\bullet KNOWAGENT	<u>69.29</u>	67.60	58.80	48.55	45.28
	WKM	70.71 +1.42	70.40 +1.74	63.75 +1.08	53.68 +3.24	49.24 +1.40
Llama-3-8B	\odot REACT	2.86	3.73	19.32	24.76	22.66
	\odot Reflexion	4.29	4.48	22.73	27.23	25.41
	\bullet NAT	60.71	59.70	61.60	55.24	48.76
	\bullet ETO	64.29	<u>64.18</u>	<u>64.57</u>	57.90	<u>52.33</u>
	\bullet KNOWAGENT	<u>66.71</u>	62.69	64.40	<u>58.67</u>	49.18
	WKM	68.57 +1.86	65.93 +1.75	66.64 +2.07	60.12 +1.55	54.75 +2.42

unseen tasks to evaluate the agent’s generalization ability. The reward of ALFWorld is binary 0 or 1, indicating whether the agent has completed the task or not. WebShop and ScienceWorld provide dense rewards from 0 to 1 to measure the completion level of the task. For all the datasets, we apply **average reward** as the final metrics. Please refer to Appendix B for detailed dataset information.

Models and Baselines. We evaluate on three state-of-the-art open-source models: 1) **Mistral-7B** [16], the Mistral-7B-Instruct-v0.2 version. 2) **Gemma-7B** [24], the Gemma-1.1-7B-it version. 3) **Llama-3-8B** [25], the Meta-Llama-3-8B-Instruct version. We compare our method with two prompt-based baselines: **REACT** [54] and **Reflexion** [40]. Besides, we adopt two strong baselines that introduce rejected trajectories into the training process to learn from experience: **NAT** [49], learn from rejected trajectories through SFT, and **ETO** [44], learn from rejected trajectories through DPO [36]. Moreover, we compare with a knowledge-augmented planning method **KNOWAGENT**. We also include **ChatGPT** (gpt-3.5-turbo-0125) [27] and **GPT-4** (gpt-4-32K-0613) [28] for comparison. All the prompt-based baselines are tested under one-shot and all the fine-tuning-based baselines are trained with LoRA [12]. Please refer to Appendix C for baselines and re-producing details.

Training and Inference Setups. We fine-tune the proposed approach with LoRA [12] using the LlamaFactory [62] framework. During training, the model is tuned after finishing the entire trajectory rather than each step of action. The learning rate is $1e-4$ and the sequence length is 2048 for all the models. The training epoch is 3 and the batch size is 32. We adopt the AdamW optimizer [22] with a cosine learning scheduler. During inference, we apply the embedding layer of WKM as the encoder and use the **cosine similarity** between sentences for retrieval. The number of retrieved action-state-action triplets \mathcal{N} is set to 3000 and the $P_{\text{agent}}(\mathcal{A}_u)$ weight γ is set to $\{0.4, 0.5, 0.7\}$. All the training and inference experiments are conducted on 8 NVIDIA V100 32G GPUs within 12 hours. Please refer to Appendix D for detailed hyperparameters used in our paper.

4.2 Results

Main Results. As shown in Table 1, **for prompt-based baselines** on open-source models, both REACT and Reflexion exhibit poor performance, far behind our method and fine-tuning-based baselines on various datasets. GPT-3.5-Turbo performs ordinarily on two datasets other than WebShop, and it even falls behind Mistral-7B and Llama-3-8B’s REACT performance on ScienceWorld. However, GPT-4 exhibits strong performance across various datasets. Nevertheless, our approach, through

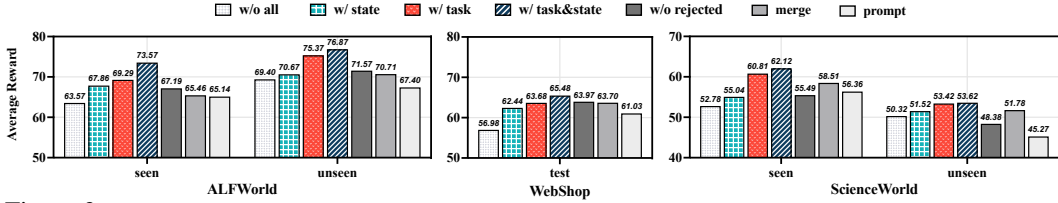


Figure 3: **Ablation Study** on Mistral-7B. **w/o all** means the vanilla experienced agent model training with pure expert trajectories. **w/ state** is testing agent model with only state knowledge base constraints. **w/ task** stands for guiding agent model with only task knowledge. **w/ task&state** is our WKM with both task knowledge guidance and state knowledge constraints. **w/o rejected** means synthesizing task knowledge solely through expert trajectories. **merge** stands for training WKM and the agent model together with one single model. **prompt** means using few-shot prompts to replace the WKM for providing knowledge.

Table 2: **Average Steps**. The maximum number of steps in ALFWorld and WebShop is 40 and 10. In ScienceWorld, the number of steps ranges from 10 to 120 depending on the task type, with an average of around 40.

Method	ALFWorld		WebShop	ScienceWorld	
	Seen	Unseen		Seen	Unseen
NAT	23.27	23.42	4.08	20.18	21.21
ETO	19.82	22.29	3.99	24.13	26.35
KNOWAGENT	18.51	24.56	4.01	21.06	24.74
WKM	17.66	17.92	3.97	18.74	19.59

Table 3: **Hallucinatory Action Rates** on ALFWorld. We calculate the proportion of trajectories containing invalid actions regardless of their correctness.

Method	ALFWorld	
	Seen	Unseen
NAT	45.71%	50.00%
ETO	34.29%	36.57%
KNOWAGENT	33.57%	44.78%
WKM	32.86%	29.85%

LoRA training alone, surpasses GPT-4 on ALFWorld (44.29→73.57 on seen, 38.05→76.87 on unseen) and WebShop (62.76→66.64). **For fine-tuning-based baselines**, both NAT and ETO fall behind our method, implying that just integrating world knowledge for agent models is worth more than further fussy SFT or DPO on negative examples. Our method also performs better than KNOWAGENT which brings human-designed fixed action knowledge and long action paths into trajectories. This suggests the effectiveness of our WKM which is responsible for generating instance-level task knowledge and maintaining implicit action constraints. Furthermore, KNOWAGENT’s performance on unseen tasks is not as impressive as on seen tasks, while WKM can keep its advantage. This phenomenon also demonstrates the generalization ability of WKM.

Approach Ablations. As shown in Figure 3, taking Mistral-7B as an example, we decompose the key components of WKM to examine the roles of the task and state knowledge separately. In a macro view, removing each module results in a clear drop in the agent’s performance, which validates the power of our world knowledge. Furthermore, the improvement through task knowledge (*w/ task*) is more pronounced than that through state knowledge (*w/ state*), suggesting the necessity of global prior knowledge for agent planning. A more micro observation reveals that the impact of state knowledge is more significant on seen tasks compared to unseen tasks, while the influence of task knowledge is sustainable across seen and unseen tasks. This may be attributed that although our real-time state knowledge is generated by WKM, the state knowledge base is built on the training set, which may weaken generalization to some extent. Additionally, to validate our motivation of allowing the agent to learn task knowledge from both expert and generated trajectories, we exclude the rejected trajectories during the synthesis of task knowledge, instructing the agent model to synthesize knowledge solely based on the chosen trajectories. The results (*w/o rejected*) demonstrate that learning from the contrast between chosen and rejected trajectories is more effective than learning from chosen examples alone. This procedure is a little similar to DPO, but we achieve it through knowledge augmentation rather than directly converting it into a loss calculation between chosen and rejected trajectories. Additional results can further evident that training a WKM separately performs better than training one single model together with the agent model as well as using few-shot prompts to replace WKM for providing knowledge.

4.3 Analysis

World knowledge can mitigate blind trial-and-error and reduce hallucinatory actions. We compare the number of planning steps for each dataset between three strong baselines and WKM and calculate the average steps of each method. As depicted in Figure 9 (in Appendix F), WKM

demonstrates the ability to complete a significant proportion of tasks using the shortest trajectory, indicating that guidance from world knowledge can effectively reduce the agent’s blind trial-and-error in the environment. Taking a further perspective from an average standpoint in Table 2, it can be observed that WKM exhibits lower average planning steps compared to other baselines. As ALFWorld can respond to invalid actions, in Table 3, we count the percentage of hallucinatory actions that occurred in trajectories from ALFWorld for each method. The results confirm the effectiveness of our world knowledge model to decrease hallucinatory actions. Furthermore, it is worth noting that most baselines show a prominent increase in the average number of steps and percentage of invalid actions when transitioning from seen tasks to unseen tasks, but WKM can still maintain a relatively low level. This reflects laterally that our world knowledge can still effectively guide the agent model on unseen tasks, highlighting the knowledge generalization brought by the world knowledge model. To see how our world knowledge works, please refer to our case study in Appendix H.

Our instance-level knowledge can generalize better to unseen tasks.

To further explore the benefit of using a knowledge model to generate instance-level task knowledge, we carefully survey the task knowledge generated by our WKM and abstract it into dataset-level knowledge for each dataset. Then we retrain the agent model to adapt to new dataset-level knowledge⁵. As illustrated in Figure 4, we compare the performance of dataset-level knowledge with our instance-level task knowledge (WKM *w/o state*) on ALFWorld and ScienceWorld. It can be observed that our model-generated instance-level knowledge not only surpasses human-designed knowledge on seen tasks but also exhibits even more remarkable performance on unseen tasks, with the improvement in performance on unseen tasks significantly greater than that on seen tasks. This phenomenon straightly reflects the strong generalization ability of our knowledge model compared to rigidly designed knowledge by humans.

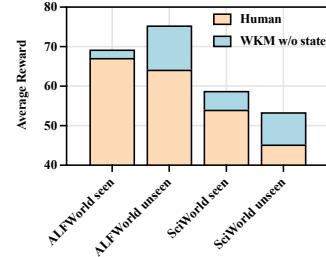


Figure 4: Performance of human-designed dataset-level knowledge compared to WKM generated instance-level knowledge.

Table 4: **Weak-guide-strong.** The knowledge model here is based on Mistral-7B.

Backbone	Method	ALFWorld	
		Seen	Unseen
GPT-3.5-Turbo	REACT	8.57	5.97
	WKM w/o state	12.86	8.96
GPT-4	REACT	44.29	38.05
	WKM w/o state	50.71	47.01

Weak knowledge model guides strong agent model planning.

In our main experiments, the knowledge model and agent model are based on the same backbone. Here, we explore on ALFWorld what will happen if we use a weak knowledge model to guide a strong agent model. We choose Mistral-7B as the backbone of the knowledge model and ChatGPT and GPT-4 as the agent model. Since we cannot get the token distribution from OpenAI API, we

only apply task knowledge to the agent model. As exhibited in Table 4, the results of both ChatGPT and GPT-4 show distinct advances after being guided by the Mistral-7B world knowledge model, indicating the weak world knowledge model also contains knowledge that the strong model may lack. In the era of LLMs, this inspires us with a new agent learning paradigm: **weak-guide-strong**. Due to its lightweight nature, the weak knowledge model can flexibly adjust its parameters based on the needs of the agent model, which can address the difficulty of large agent models in adapting to new environments through fine-tuning.

Unified World Knowledge Model Training.

We mix the world knowledge collected from all three datasets and jointly train one single world knowledge model to investigate the effect of multi-task world knowledge learning. Figure 5 illustrates the relative performance comparison between multi-task WKM and various baselines, from which we can observe that multi-task WKM not only does not lead to performance degradation but also exhibits visible improvements compared to single-task WKM, especially on WebShop and ScienceWorld. Similar to [57, 58, 3] which endeavor to train a unified agent model and achieve strong generalization ability to held-out tasks, this observation inspires us with the potential of training a unified world knowledge model that can be applied to help

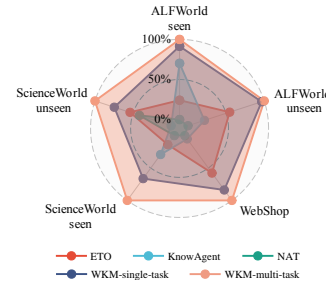


Figure 5: Relative performance of multi-task WKM compared to various baselines.

⁵Detailed manually designed dataset-level knowledge prompt can be found in Appendix I.3

various held-in agent models and also generalize to guide held-out agent models. A more daring idea is whether a unified agent model combined with a unified world knowledge model is the key to Artificial General Intelligence (AGI).

Explicit state knowledge will hurt the planning performance.

To demonstrate the rationality of our choice to construct a state knowledge base, we explore the effect of directly incorporating state knowledge into the context of the agent model (we retrain the agent model to follow both the task and state knowledge), as shown in Figure 6. The performance of explicit state knowledge is far inferior to our approach of retrieving from a state knowledge base and utilizing probabilistic constraints. It even performs worse than when we remove state knowledge and only include task knowledge. This clearly indicates that blindly extending prompts with a large amount of explicit natural language feedback is lose-more-than-gain for agent planning, and implicit knowledge constraints may be sometimes more prudent.

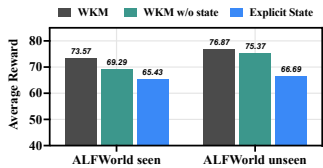


Figure 6: Performance of explicit state knowledge.

Case Study. In Figure 10 (Appendix H), we list the trajectories of ETO and our WKM within the same task in ALFWORLD to illustrate how world knowledge functions. **The rationales before each action have been omitted to guarantee a clear illustration.** The task is to clean some soapbar and put it in cabinet. Initially, ETO blindly searches for the soapbar in the countertop and cabinet, introducing a lot of irrelevant information and unnecessary context. In the later stages of planning, ETO experiences the hallucination and executes the put action after close the cabinet, causing the environment to become unrecognizable and resulting in a collapse. On the contrary, guided by task knowledge, WKM directly identified the possible locations of the soapbar and successfully found it in the first attempt. Subsequently, WKM efficiently completed the task with precision, adhering to the constraints of state knowledge.

5 Related Work

LLM Agents. LLMs have emerged as a promising avenue towards unlocking the potential of Artificial General Intelligence, offering robust support for the development of agent systems [48, 51, 8, 63]. Existing works in this field mainly focuses on agent planning [14, 21, 54, 42], external tools harnessing [39, 23, 43, 29, 32, 35, 46], code generation [45, 21, 31, 11], etc. Recently, there has been an increasing focus on endowing open-source LLMs with agent functionalities through fine-tuning [2, 57, 56, 38, 44, 49]. However, these approaches rely on blindly fitting the probabilities of tokens to learn planning, without having an intimate cognition of the environment. The lack of knowledge can lead to the agent blindly attempting trial-and-error and generating hallucinatory actions.

Knowledge Augmented Agent Planning. Planning [15] is a crucial capability for intelligent agents to accomplish real-world tasks, often requiring agents to possess rich knowledge and environmental commonsense. Few works have explored the field of knowledge-augmented agent planning. [14, 61, 5] utilize the rich parametric knowledge stored in pre-trained language models to assist agent planners. [7, 20, 59, 64] design structured or natural language knowledge to regulate the actions. However, the above studies require the manual design of fixed prompt templates or task procedures, making it challenging to transfer across different task environments. [63, 55, 6] propose the automation of knowledge generation using language models. However, their knowledge either consists of only global workflow or only local action principles. In contrast, we train our world knowledge model both on global task knowledge and local state knowledge to assist agent planning, and these knowledge sources are derived from the model’s self-summary rather than hand-curated.

LLM-based World Model. World model and agent model often co-occur in the domain of reinforcement learning and robotics [13, 9, 19, 37, 26, 4]. With LLMs commonly deemed as the most powerful intelligent machines constructed by humans thus far, the LLM-backed world models have been proposed [61, 10, 13]. In our paper, we attempt to self-synthesize world knowledge and train to obtain a world knowledge model. However, we consider our model to be a world **knowledge model** rather than a **world model** based on the reason that our model is temporarily unable to utilize search algorithms (e.g. MCTS) in conjunction with the agent model to make predictions about the world and we leave this for our future work.

6 Conclusion and Future Work

In this paper, we strive to develop a parametric world knowledge model (WKM) to augment language agent model planning. Our WKM can generate prior task knowledge to guide global planning as well as dynamic state knowledge to regulate local planning. Our extensive results show that our world knowledge can work on both GPT-4 and state-of-the-art open-source models and achieve superior performance compared to various strong baselines. Analytical experiments validate that our WKM can 1) reduce brainless trial-and-error and invalid actions, 2) generalize better to unseen tasks, 3) achieve weak-guide-strong, and 4) be effectively extended to unified world knowledge training. Potential future directions include: 1) building a unified world knowledge model, 2) learning to predict the world like a world model, 3) applying to multi-modal agent planning, etc.

Limitations

Despite our best efforts, this paper still has some limitations: 1) Our primary intention behind designing the WKM is to compensate for the lack of world knowledge in the agent model. However, determining what a language model knows and doesn't know has been an ongoing challenge that remains unresolved. 2) It is widely acknowledged that world knowledge extends beyond textual representations. While our world knowledge is currently limited to textual information, exploring multi-modal world knowledge models is indeed one of our important future tasks. 3) Our world knowledge model cannot dynamically update with the changes of the world and feedback from the agent. 4) Generating world knowledge can introduce additional inference overhead.

Acknowledgments and Disclosure of Funding

We would like to express our great gratitude to the anonymous reviewers for their kind comments. This work was supported by the National Natural Science Foundation of China (No. 62206246, No. NSFCU23B2055, No. NSFCU19B2027), the Fundamental Research Funds for the Central Universities (226-2023-00138), Zhejiang Provincial Natural Science Foundation of China (No. LGG22F030011), Yongjiang Talent Introduction Programme (2021A-156-G), CIPSC-SMP-Zhipu Large Model Cross-Disciplinary Fund, Ningbo Science and Technology Special Projects under Grant No. 2023Z212, Information Technology Center and State Key Lab of CAD&CG, Zhejiang University, and NUS-NCS Joint Laboratory (A-0008542-00-00). We gratefully acknowledge the support of Zhejiang University Education Foundation Qizhen Scholar Foundation.

References

- [1] Robert Eamon Briscoe. Mental imagery and the varieties of amodal perception. *Pacific Philosophical Quarterly*, 92(2):153–173, 2011.
- [2] Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao. Fireact: Toward language agent fine-tuning. *CoRR*, abs/2310.05915, 2023.
- [3] Zehui Chen, Kuikun Liu, Qiuchen Wang, Wenwei Zhang, Jiangning Liu, Dahua Lin, Kai Chen, and Feng Zhao. Agent-flan: Designing data and methods of effective agent tuning for large language models. *CoRR*, abs/2403.12881, 2024.
- [4] Anna Dawid and Yann LeCun. Introduction to latent variable energy-based models: A path towards autonomous machine intelligence. *CoRR*, abs/2306.02572, 2023.
- [5] Yan Ding, Xiaohan Zhang, Saeid Amiri, Nieqing Cao, Hao Yang, Andy Kaminski, Chad Esselink, and Shiqi Zhang. Integrating action knowledge and llms for task planning and situation handling in open worlds. *Auton. Robots*, 47(8):981–997, 2023.
- [6] Yao Fu, Dong-Ki Kim, Jaekyeom Kim, Sungryull Sohn, Lajanugen Logeswaran, Kyunghoon Bae, and Honglak Lee. Autoguide: Automated generation and selection of state-aware guide-lines for large language model agents. *CoRR*, abs/2403.08978, 2024.

- [7] Jian Guan, Wei Wu, Zujie Wen, Peng Xu, Hongning Wang, and Minlie Huang. AMOR: A recipe for building adaptable modular knowledge agents through process feedback. *CoRR*, abs/2402.01469, 2024.
- [8] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges. *CoRR*, abs/2402.01680, 2024.
- [9] David Ha and Jürgen Schmidhuber. World models. *CoRR*, abs/1803.10122, 2018.
- [10] 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. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 8154–8173. Association for Computational Linguistics, 2023.
- [11] Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, and Chenglin Wu. Metagpt: Meta programming for multi-agent collaborative framework. *CoRR*, abs/2308.00352, 2023.
- [12] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022.
- [13] Zhiting Hu and Tianmin Shu. Language models, agent models, and world models: The LAW for machine reasoning and planning. *CoRR*, abs/2312.05230, 2023.
- [14] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 9118–9147. PMLR, 2022.
- [15] Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. Understanding the planning of LLM agents: A survey. *CoRR*, abs/2402.02716, 2024.
- [16] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b. *CoRR*, abs/2310.06825, 2023.
- [17] Philip N Johnson-Laird. Mental models and human reasoning. *Proceedings of the National Academy of Sciences*, 107(43):18243–18250, 2010.
- [18] Philip Nicholas Johnson-Laird. *Mental models: Towards a cognitive science of language, inference, and consciousness*. Harvard University Press, 1983.
- [19] Lukasz Kaiser, Mohammad Babaeizadeh, Piotr Milos, Blazej Osinski, Roy H. Campbell, Konrad Czechowski, Dumitru Erhan, Chelsea Finn, Piotr Kozakowski, Sergey Levine, Afroz Mohiuddin, Ryan Sepassi, George Tucker, and Henryk Michalewski. Model based reinforcement learning for atari. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- [20] Zelong Li, Wenyue Hua, Hao Wang, He Zhu, and Yongfeng Zhang. Formal-llm: Integrating formal language and natural language for controllable llm-based agents. *CoRR*, abs/2402.00798, 2024.

- [21] Lajanugen Logeswaran, Yao Fu, Moontae Lee, and Honglak Lee. Few-shot subgoal planning with language models. In Marine Carpuat, Marie-Catherine de Marneffe, and Iván Vladimir Meza Ruíz, editors, *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, pages 5493–5506. Association for Computational Linguistics, 2022.
- [22] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- [23] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [24] Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussonot, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, and et al. Gemma: Open models based on gemini research and technology. *CoRR*, abs/2403.08295, 2024.
- [25] Meta. Introducing meta llama 3: The most capable openly available llm to date, 2024. <https://ai.meta.com/blog/meta-llama-3/>.
- [26] Thomas M. Moerland, Joost Broekens, Aske Plaat, and Catholijn M. Jonker. Model-based reinforcement learning: A survey. *Found. Trends Mach. Learn.*, 16(1):1–118, 2023.
- [27] OpenAI. Chatgpt: Optimizing language models for dialogue, 2022. <https://openai.com/blog/chatgpt/>.
- [28] OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023.
- [29] Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large language model connected with massive apis. *CoRR*, abs/2305.15334, 2023.
- [30] RT Pramod, Michael Cohen, Kirsten Lydic, Josh Tenenbaum, and Nancy Kanwisher. Evidence that the brain’s physics engine runs forward simulations of what will happen next. *Journal of Vision*, 20(11):1521–1521, 2020.
- [31] Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. Communicative agents for software development. *CoRR*, abs/2307.07924, 2023.
- [32] Shuofei Qiao, Honghao Gui, Huajun Chen, and Ningyu Zhang. Making language models better tool learners with execution feedback. *CoRR*, abs/2305.13068, 2023.
- [33] Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang, and Huajun Chen. Reasoning with language model prompting: A survey. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 5368–5393. Association for Computational Linguistics, 2023.
- [34] Shuofei Qiao, Ningyu Zhang, Runnan Fang, Yujie Luo, Wangchunshu Zhou, Yuchen Eleanor Jiang, Chengfei Lv, and Huajun Chen. AUTOACT: automatic agent learning from scratch via self-planning. *CoRR*, abs/2401.05268, 2024.
- [35] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language models to master 16000+ real-world apis. *CoRR*, abs/2307.16789, 2023.

- [36] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [37] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, Timothy P. Lillicrap, and David Silver. Mastering atari, go, chess and shogi by planning with a learned model. *Nat.*, 588(7839):604–609, 2020.
- [38] Weizhou Shen, Chenliang Li, Hongzhan Chen, Ming Yan, Xiaojun Quan, Hehong Chen, Ji Zhang, and Fei Huang. Small llms are weak tool learners: A multi-llm agent. *CoRR*, abs/2401.07324, 2024.
- [39] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving AI tasks with chatgpt and its friends in hugging face. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [40] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: language agents with verbal reinforcement learning. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [41] Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew J. Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [42] Chan Hee Song, Brian M. Sadler, Jiaman Wu, Wei-Lun Chao, Clayton Washington, and Yu Su. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pages 2986–2997. IEEE, 2023.
- [43] Yifan Song, Weimin Xiong, Dawei Zhu, Cheng Li, Ke Wang, Ye Tian, and Sujian Li. Rest-gpt: Connecting large language models with real-world applications via restful apis. *CoRR*, abs/2306.06624, 2023.
- [44] Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. Trial and error: Exploration-based trajectory optimization for LLM agents. *CoRR*, abs/2403.02502, 2024.
- [45] Haotian Sun, Yuchen Zhuang, Lingkai Kong, Bo Dai, and Chao Zhang. Adaplaner: Adaptive planning from feedback with language models. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [46] Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, and Le Sun. Toolalpaca: Generalized tool learning for language models with 3000 simulated cases. *CoRR*, abs/2306.05301, 2023.
- [47] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, and et al. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023.
- [48] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Jirong Wen. A survey on large language model based autonomous agents. *Frontiers Comput. Sci.*, 18(6):186345, 2024.

- [49] Renxi Wang, Haonan Li, Xudong Han, Yixuan Zhang, and Timothy Baldwin. Learning from failure: Integrating negative examples when fine-tuning large language models as agents. *CoRR*, abs/2402.11651, 2024.
- [50] Ruoyao Wang, Peter A. Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. Science-world: Is your agent smarter than a 5th grader? In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 11279–11298. Association for Computational Linguistics, 2022.
- [51] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huan, and Tao Gui. The rise and potential of large language model based agents: A survey. *CoRR*, abs/2309.07864, 2023.
- [52] Jiannan Xiang, Tianhua Tao, Yi Gu, Tianmin Shu, Zirui Wang, Zichao Yang, and Zhiting Hu. Language models meet world models: Embodied experiences enhance language models. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10-16, 2023*, 2023.
- [53] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- [54] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- [55] Yining Ye, Xin Cong, Shizuo Tian, Jiannan Cao, Hao Wang, Yujia Qin, Yaxi Lu, Heyang Yu, Huadong Wang, Yankai Lin, Zhiyuan Liu, and Maosong Sun. Proagent: From robotic process automation to agentic process automation. *CoRR*, abs/2311.10751, 2023.
- [56] Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. Lumos: Learning agents with unified data, modular design, and open-source llms. *CoRR*, abs/2311.05657, 2023.
- [57] Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. Agenttuning: Enabling generalized agent abilities for llms. *CoRR*, abs/2310.12823, 2023.
- [58] Jianguo Zhang, Tian Lan, Rithesh Murthy, Zhiwei Liu, Weiran Yao, Juntao Tan, Thai Hoang, Liangwei Yang, Yihao Feng, Zuxin Liu, Tulika Awalganekar, Juan Carlos Niebles, Silvio Savarese, Shelby Heinecke, Huan Wang, and Caiming Xiong. Agentohana: Design unified data and training pipeline for effective agent learning. *CoRR*, abs/2402.15506, 2024.
- [59] Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. Expel: LLM agents are experiential learners. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan, editors, *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2024, February 20-27, 2024, Vancouver, Canada*, pages 19632–19642. AAAI Press, 2024.
- [60] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models. *CoRR*, abs/2303.18223, 2023.

- [61] Zirui Zhao, Wee Sun Lee, and David Hsu. Large language models as commonsense knowledge for large-scale task planning. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [62] Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, and Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. *CoRR*, abs/2403.13372, 2024.
- [63] Wangchunshu Zhou, Yuchen Eleanor Jiang, Long Li, Jialong Wu, Tiannan Wang, Shi Qiu, Jintian Zhang, Jing Chen, Ruipu Wu, Shuai Wang, Shiding Zhu, Jiyu Chen, Wentao Zhang, Ningyu Zhang, Huajun Chen, Peng Cui, and Mrinmaya Sachan. Agents: An open-source framework for autonomous language agents. *CoRR*, abs/2309.07870, 2023.
- [64] Yuqi Zhu, Shuofei Qiao, Yixin Ou, Shumin Deng, Ningyu Zhang, Shiwei Lyu, Yue Shen, Lei Liang, Jinjie Gu, and Huajun Chen. Knowagent: Knowledge-augmented planning for llm-based agents. *CoRR*, abs/2403.03101, 2024.

A Expert Trajectories Collection

We mainly use the expert trajectories with a REACT-style [54] collected from [44]:

1. **ALFWorld** [41]. The dataset provides human-annotated trajectories.
2. **WebShop** [53]. Except for human-annotated trajectories, GPT-4 is also applied to generate trajectories with a reward larger than 0.7 being reserved.
3. **ScienceWorld** [50]. The dataset offers heuristic algorithms to search golden trajectories for each sub-task.

Since the original golden trajectories do not contain rationales, GPT-4 is further leveraged to generate the corresponding information.

B Dataset Information

We evaluate our method on three real-world simulated agent planning datasets: ALFWorld [41], WebShop [53], and ScienceWorld [50].

1. **ALFWorld** is a household dataset requiring the agent to navigate through the room and manipulate objects. Except for seen tasks, ALFWorld also includes unseen tasks to evaluate the agent’s generalization ability. The reward of ALFWorld is binary 0 or 1, indicating whether the agent has completed the task or not.
2. **WebShop** is an online shopping dataset in a website environment. It provides dense final rewards from 0 to 1 to measure the completion level of the task.
3. **ScienceWorld** is a scientific reasoning dataset which is at the level of a standard elementary school science curriculum. It also possesses both seen and unseen parts and a dense reward function from 0 to 1.

For all the datasets, we apply **average reward** as the final metrics. Table 5 illustrates the statistics of each dataset.

Table 5: Dataset statistics.

Dataset	Train	Text-Seen	Text-Unseen
ALFWorld	3,119	140	134
WebShop	1,824	200	-
ScienceWorld	1,483	194	211

C Compared Baselines

Here we detailedly introduce the baselines we compare with and our re-produce details.

1. **REACT** [54]. The first approach incorporates Chain-of-Thought (CoT) prompting in agent planning tasks with a format of Thought-Action-Observation loop. In our paper, we apply one-shot prompting for REACT⁶.
2. **Reflexion** [40]. A strong prompt-based baseline reinforces agent planning with verbal feedback. Manually designed prompts are used to enable the agent to reflect on the historical trajectory and re-plan based on the feedback. In our paper, we utilize one-shot prompting for reflection and select the first reflect iteration as our result due to limited context⁷.
3. **NAT** [49]. NAT includes negative trajectories by employing different prompts during agent fine-tuning. When evaluating, only positive prompts are used to encourage the language

⁶<https://github.com/ysmyth/ReAct>

⁷<https://github.com/noahshinn/reflexion>

agent to generate correct trajectories. As it also follows the REACT-style format, we directly use the default positive and negative prompts and train with LoRA in our paper⁸.

4. **ETO** [44]. Another baseline includes negative trajectories during agent training. The method contains two training phases, of which the first phase is behavior cloning which fine-tunes the agent on expert trajectories, and the second phase is learning from failures which further fine-tunes the agent through Direct Preference Optimization (DPO) [36]. In our paper, we remove the one-shot prompt for fairness and retain all the default hyperparameters proposed in ETO except for LoRA training⁹.
5. **KNOWAGENT** [64]. KNOWAGENT is a knowledge-augmented agent planning baseline that applies action knowledge in the prompt and maintains an action path in the context during planning to constrain the agent’s action. We directly use the default prompt mentioned in KNOWAGENT for ALFWorld and carefully extend it to WebShop and ScienceWorld by following a similar format¹⁰.

All the prompt-based baselines are tested under one-shot and all the fine-tuning-based baselines are trained with LoRA [12].

D Hyperparameters

The detailed hyperparameters we use during training and inference are shown in Table 6. We employ identical hyperparameters for different models. The temperature of the agent model is set to 0.0 when conducting exploration and 0.5 when introduced into WKM. The temperature of WKM is set to 0.0 for all the time. The $P_{\text{agent}}(\mathcal{A}_u)$ weight γ is set to 0.4 for ALFWorld, 0.5 for WebShop, and 0.7 for ScienceWorld.

Table 6: Detailed hyperparameters used in our paper.

Name	Value
lora r	8
lora alpha	16
lora dropout	0.05
lora target modules	q_proj, v_proj
cutoff len	2048
epochs	3
batch size	32
batch size per device	4
gradient accumulation steps	2
learning rate	1e-4
warmup ratio	0.03
temperature	0.0, 0.5
retrieved state knowledge \mathcal{N}	3000
$P_{\text{agent}}(\mathcal{A}_u)$ weight γ	0.4, 0.5, 0.7

E Training Data Example

We show the training data example for agent model and world knowledge model in Figure 7 and Figure 8. During training, transparent tokens do not contribute to the calculation of the loss. Only the parts of the model responsible for output (agent model: purple, knowledge world model: green), participate in the loss calculation.

⁸<https://github.com/Reason-Wang/NAT>

⁹<https://github.com/Yifan-Song793/ETO>

¹⁰<https://github.com/zjunlp/KnowAgent>

Training Data Example for Agent Model

Task Instruction: You are in the middle of a room. Looking quickly around you, you see a bed 1, a desk 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer 1, a dresser 1, a garbagecan 1, a laundryhamper 1, and a shelf 1. Your task is to: examine the book with the desklamp.

Task Knowledge: When you cannot find the object required for the task under the desklamp, first check all the drawers and shelves in the room, then revisit the previous locations and double-check before giving up and looking in the garbage can or laundry hamper. The action workflows are as follows: 1. For each drawer and shelf, go to the location and open it to check if the object is inside. 2. If the object is not found, close the drawer or shelf and move on to the next one. 3. After checking all the drawers and shelves, revisit the previous locations and double-check if the object might have been overlooked. 4. If the object is still not found, look in the garbage can or laundry hamper as a last resort.

Thought: The task requires me to examine a book with the desklamp. First, I need to locate the book.

Action: go to bed 1

Observation: On the bed 1, you see a book 1, a laptop 1, and a pillow 1.

Thought: I see a book on the bed. I should pick it up to examine it with the desklamp.

Action: take book 1 from bed 1

Observation: You pick up the book 1 from the bed 1.

Thought: Now that I have the book, I need to find the desklamp to examine the book under its light.

Action: go to desk 1

Observation: On the desk 1, you see an alarm clock 1, a bowl 1, a CD 2, a mug 3, a mug 2, a mug 1, a pen 1, and a pencil 2.

Thought: The desklamp isn't on the desk. I need to look around the room to find where the desklamp is located.

Action: go to dresser 1

Observation: On the dresser 1, you see a cellphone 1, a desklamp 1, and a keychain 2.

Thought: I've found the desklamp on the dresser. I should use it to examine the book.

Action: use desklamp 1

Figure 7: Training Data Example for Agent Model.

Training Data Example for World Knowledge Model

Task Instruction: You are in the middle of a room. Looking quickly around you, you see a bed 1, a desk 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer 1, a dresser 1, a garbagecan 1, a laundryhamper 1, and a shelf 1. Your task is to: examine the book with the desklamp.

Task Knowledge: When you cannot find the object required for the task under the desklamp, first check all the drawers and shelves in the room, then revisit the previous locations and double-check before giving up and looking in the garbage can or laundry hamper. The action workflows are as follows: 1. For each drawer and shelf, go to the location and open it to check if the object is inside. 2. If the object is not found, close the drawer or shelf and move on to the next one. 3. After checking all the drawers and shelves, revisit the previous locations and double-check if the object might have been overlooked. 4. If the object is still not found, look in the garbage can or laundry hamper as a last resort.

Thought: The task requires me to examine a book with the desklamp. First, I need to locate the book.

Action: go to bed 1

Observation: On the bed 1, you see a book 1, a laptop 1, and a pillow 1.

State Knowledge: Your task is to examine a book with the desklamp. You have not found the desklamp yet, and have only checked the bed 1.

Thought: I see a book on the bed. I should pick it up to examine it with the desklamp.

Action: take book 1 from bed 1

Observation: You pick up the book 1 from the bed 1.
State Knowledge: Your task is to examine a book using the desk lamp. You have located the book on bed 1 and picked it up.
Thought: Now that I have the book, I need to find the desk lamp to examine the book under its light.
Action: go to desk 1
Observation: On the desk 1, you see an alarm clock 1, a bowl 1, a CD 2, a mug 3, a mug 2, a mug 1, a pen 1, and a pencil 2.
State Knowledge: Your task is to examine a book with the desk lamp. You have located the book on bed 1 and picked it up. Next, You need to find the desk lamp to examine the book under its light.
Thought: The desk lamp isn't on the desk. I need to look around the room to find where the desk lamp is located.
Action: go to dresser 1
Observation: On the dresser 1, you see a cellphone 1, a desk lamp 1, and a keychain 2.
State Knowledge: Your task is to examine a book with the desk lamp. You have located the book on the bed and picked it up, now you find a desk lamp on a dresser.

Figure 8: Training Data Example for World Knowledge Model.

F Win Rate of Planning Steps

See Figure 9.

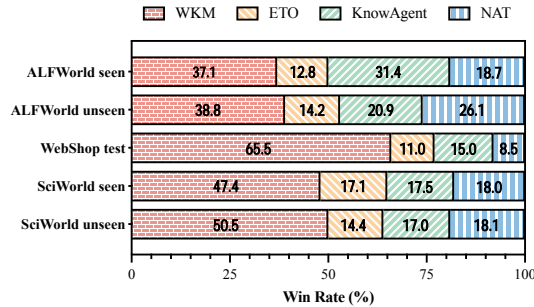


Figure 9: **Win Rate of Planning Steps.** We choose the method with the shortest steps for each task and calculate the proportion.

G Impact of γ

In fact, the ratio γ can be viewed as a signal to reflect whether knowledge or planning is more important for a task. To understand which part of the output action has the most significant impact, we further analyze $\gamma = 0$ (fully trust state knowledge base) and $\gamma = 1$ (fully trust agent model, equivalent to remove state knowledge in Figure 3). The empirical results can be seen in Table 7. It can be observed that state knowledge primarily serves as a constraint to alleviate hallucinated actions for the agent model. However, when we fully trust it ($\gamma = 0$), its lack of generalization significantly harms the performance of the agent model.

H Case Study

In Figure 10, we list the trajectories of ETO and our WKM within the same task in ALFWorld to illustrate how world knowledge functions. **The rationales before each action have been omitted to guarantee a clear illustration.** The task is to clean some soapbar and put it in cabinet. Initially, ETO blindly searches for the soapbar in the countertop and cabinet, introducing a

Table 7: **Impact of γ** . In fact, the ratio γ can be viewed as a signal to reflect whether knowledge or planning is more important for a task. To understand which part of the output action has the most significant impact, we further analyze $\gamma = 0$ (fully trust state knowledge base) and $\gamma = 1$ (fully trust agent model, equivalent to remove state knowledge in Figure 3). It can be observed that state knowledge primarily serves as a constraint to alleviate hallucinated actions for the agent model. However, when we fully trust it ($\gamma = 0$), its lack of generalization significantly harms the performance of the agent model.

Method	ALFWorld		WebShop	ScienceWorld	
	Seen	Unseen		Seen	Unseen
$\gamma = 0$	1.58	0.00	25.83	18.69	15.37
$\gamma = 1$	69.29	75.37	63.68	60.81	53.42
WKM	73.57	76.87	65.48	62.12	53.62

lot of irrelevant information and unnecessary context. In the later stages of planning, ETO experiences the hallucination and executes the `put` action after `close the cabinet`, causing the environment to become unrecognizable and resulting in a collapse. On the contrary, guided by task knowledge, WKM directly identified the possible locations of the `soapbar` and successfully found it in the first attempt. Subsequently, WKM efficiently completed the task with precision, adhering to the constraints of state knowledge.

I Prompts

In this section, we illustrate all the prompts used in our paper. Part of our prompts are taken from AutoGuide [6].

I.1 Task Knowledge Synthesis Prompt

See Figure 11.

I.2 State Knowledge Summarization Prompt

See Figure 12.

I.3 Dataset-Level Knowledge Prompt

See Figure 13.

J Ethics Statement

This research was conducted following the ethical standards and best practices. All our experiments use publicly available datasets (as detailed in Appendix B), avoiding ethical concerns related to privacy, confidentiality, or misuse of personal biological information. However, despite our best efforts, it is not avoidable if someone maliciously modifies the world knowledge model to contradict the world’s knowledge and leads the agent to engage in unethical behavior.

ETO

Task Instruction: You are in the middle of a room. Looking quickly around you, you see a cabinet 4, a cabinet 3, a cabinet 2, a cabinet 1, a countertop 1, a garbagecan 1, a handtowelholder 2, a handtowelholder 1, a sinkbasin 2, a sinkbasin 1, a toilet 1, a toiletpaperhanger 1, and a towelholder 1. Your task is to: clean some soapbar and put it in cabinet.

Blind trial-and-error begins.

Agent: go to countertop 1

Observation: On the countertop 1, you see a candle 3, a candle 1, a soapbottle 1, a spraybottle 2, and a toiletpaper 2.

Agent: go to cabinet 1

Observation: The cabinet 1 is closed.

Agent: open cabinet 1

Observation: You open the cabinet 1. The cabinet 1 is open. In it, you see a soapbottle 2.

Agent: close cabinet 1

Observation: You close the cabinet 1.

Blind trial-and-error ends.

Agent: go to toilet 1

Observation: On the toilet 1, you see a soapbar 1, and a spraybottle 3.

Agent: take soapbar 1 from toilet 1

Observation: You pick up the soapbar 1 from the toilet 1.

Agent: go to sinkbasin 1

Observation: On the sinkbasin 1, you see nothing.

Agent: clean soapbar 1 with sinkbasin 1

Observation: You clean the soapbar 1 using the sinkbasin 1.

Agent: go to cabinet 1

Observation: The cabinet 1 is closed.

Agent: open cabinet 1

Observation: You open the cabinet 1. The cabinet 1 is open. In it, you see a soapbottle 2.

Agent: close cabinet 1

Observation: You close the cabinet 1.

Hallucinatory action.

Agent: put soapbar 1 in/on cabinet 1

Observation: Nothing happens.

Agent: check cabinet 1

Observation: Nothing happens.

Agent: go to exit

Observation: Nothing happens.

Agent: exit

Observation: Nothing happens.

.....

Reward: 0.0

WKM

Task Instruction: ...

Task Knowledge: When the task is to clean and place an object, you should first locate the object. The possible locations of the soapbar are: toilet, sinkbasin. Once the object is found, you should clean it using the sinkbasin, then put it in the cabinet. The action workflows are as follows: 1. Go to the place where the object is located. 2. Take the object from the place. 3. Clean the object using the sinkbasin. 4. Go to the cabinet. 5. Put the object in/on the cabinet.

Agent: go to toilet 1

Observation: On the toilet 1, you see a soapbar 1, and a spraybottle 3.

State Knowledge: Your task is to clean some soapbar and put it in cabinet. Now you are checking toilet 1 and there is a soapbar 1.

Agent: take soapbar 1 from toilet 1

Observation: You pick up the soapbar 1 from the toilet 1.

State Knowledge: Your task is to clean some soapbar and put it in cabinet. You have checked the toilet 1 and found the soapbar 1.

Agent: go to sinkbasin 1

Observation: On the sinkbasin 1, you see nothing.

State Knowledge: Your task is to clean some soapbar and put it in cabinet. You have checked the toilet 1 and found the soapbar 1. Now you are at the sinkbasin.

Agent: clean soapbar 1 with sinkbasin 1

Observation: You clean the soapbar 1 using the sinkbasin 1.

State Knowledge: Your task is to clean some soapbar and put it in cabinet. You have cleaned the soapbar 1 using the sinkbasin 1.

Agent: go to cabinet 1

Observation: The cabinet 1 is closed.

State Knowledge: Your task is to clean some soapbar and put it in cabinet. You have cleaned the soapbar 1 using the sinkbasin 1. Now you are checking cabinet 1.

Agent: open cabinet 1

Observation: You open the cabinet 1. In it, you see a soapbottle 2.

State Knowledge: Your task is to clean some soapbar and put it in cabinet. You have cleaned the soapbar 1 and opened the cabinet 1.

Agent: put soapbar 1 in/on cabinet 1

Reward: 1.0

Figure 10: Case study for our WKM.

Prompt for Task Knowledge Synthesis

Task Knowledge

Prompt for Synthesis: I will provide you with an analysis of both a successful trajectory and an explored trajectory for the same task. By comparing the two, we can identify the key factors that contribute to success. Based on this analysis, you need to generate task-related task knowledge to help increase the success rate of future endeavors.

Success Trajectory: **Success_T**
Explored Trajectory: **Explored_T**

The task knowledge should specify what to do in what task. Here is a task knowledge example:

Task Knowledge Example

You should make your answer concise. Put your answer in this format: Task Knowledge: When ... you should (or should not) ... The action workflows are: ...

Figure 11: Prompt for Task Knowledge Synthesis.

Prompt for State Knowledge Synthesis

State Knowledge

Prompt for Synthesis: You'll get a segment of a trajectory of a text-based task. Your task is to generate a brief and general state knowledge of the task state now, following "State Knowledge: ". Keep it wise and general for the same task. Here is an example:

State Knowledge Example

Now it's your turn. Here is the trajectory :

Trajectory

Make sure your output is within 128 tokens.
Put your answer in this format: State Knowledge: ...

Figure 12: Prompt for State Knowledge Summarization.

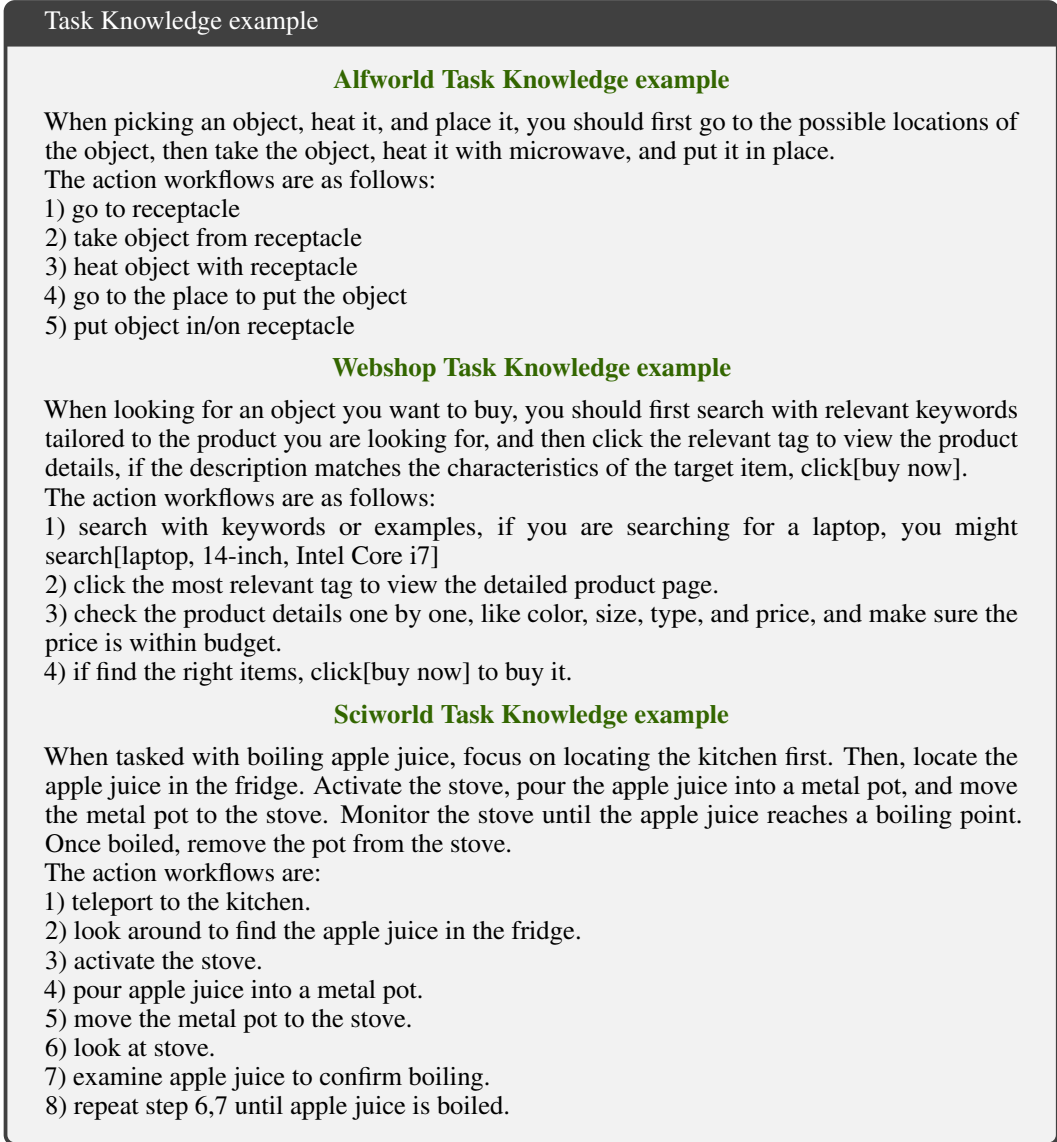


Figure 13: Dataset-Level Task Knowledge Examples.

NeurIPS Paper Checklist

1. Claims

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

Answer: [Yes]

Justification: The main claims made in the abstract and introduction do accurately reflect our paper's contributions and scope.

Guidelines:

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

2. Limitations

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

Answer: [Yes]

Justification: See Section Limitations.

Guidelines:

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

3. Theory Assumptions and Proofs

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

Answer: [NA]

Justification: Our paper does not contain theoretical results.

Guidelines:

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

4. Experimental Result Reproducibility

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

Answer: [Yes]

Justification: See Section 4 and Appendix Hyperparameters.

Guidelines:

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

5. Open access to data and code

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

Answer: [Yes]

Justification: We include our source code and data in the supplemental material submission, and we outline the data generation procedure, the evaluation protocol, the training regime, and everything else necessary for reproduction either in the main body of the paper or in the appendix.

Guidelines:

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

6. Experimental Setting/Details

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

Answer: [Yes]

Justification: See Section 4.1 and Appendix.

Guidelines:

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

7. Experiment Statistical Significance

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

Answer: [Yes]

Justification: We include the average planning steps of each method.

Guidelines:

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

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

8. Experiments Compute Resources

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

Answer: [Yes]

Justification: See Section 4.1.

Guidelines:

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

9. Code Of Ethics

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

Answer: [Yes]

Justification: We have carefully checked the NeurIPS Code of Ethics and ensure our code aligns with the request.

Guidelines:

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

10. Broader Impacts

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

Answer: [Yes]

Justification: See Appendix I.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

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

Answer: [NA]

Justification: Our paper does not have high risk PLMs or datasets.

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes]

Justification: We have cited all the corresponding papers.

Guidelines:

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

- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. **New Assets**

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

Answer: [NA]

Justification: Our paper does not release new assets.

Guidelines:

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

14. **Crowdsourcing and Research with Human Subjects**

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

Answer: [NA]

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

Guidelines:

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

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

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

Answer: [NA]

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

Guidelines:

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