

DAVIS: Planning Agent with Knowledge Graph-Powered Inner Monologue

Anonymous ACL submission

Abstract

Designing a generalist scientific agent capable of performing tasks in laboratory settings to assist researchers has become a key goal in recent Artificial Intelligence (AI) research. Unlike everyday tasks, scientific tasks are inherently more delicate and complex, requiring agents to possess a higher level of reasoning ability, structured and temporal understanding of their environment, and a strong emphasis on safety. Existing approaches often fail to address these multifaceted requirements. To tackle these challenges, we present DAVIS¹. Unlike traditional retrieval-augmented generation (RAG) approaches, DAVIS incorporates structured and temporal memory, which enables model-based planning. Additionally, DAVIS implements an agentic, multi-turn retrieval system, similar to a human’s inner monologue, allowing for a greater degree of reasoning over past experiences. DAVIS demonstrates substantially improved performance on the ScienceWorld benchmark comparing to previous approaches on 8 out of 9 elementary science subjects. In addition, DAVIS’s World Model demonstrates competitive performance on the famous HotpotQA and MusiqueQA dataset for multi-hop question answering. To the best of our knowledge, DAVIS is the first RAG agent to employ an interactive retrieval method in a RAG pipeline.

1 Introduction

A core focus of current Artificial Intelligence (AI) research is the development of artificial agents capable of autonomously performing human tasks with high decision-making autonomy (Ahn et al., 2022; Zhao et al., 2024; Wang et al., 2024; Putta et al., 2024). While Reinforcement Learning (RL) has traditionally been used to create goal-oriented agents in Markovian environments (Mnih et al., 2013; Schrittwieser et al., 2020; Hafner et al.,

¹All code and prompts are available at [Anonymous Github](#)

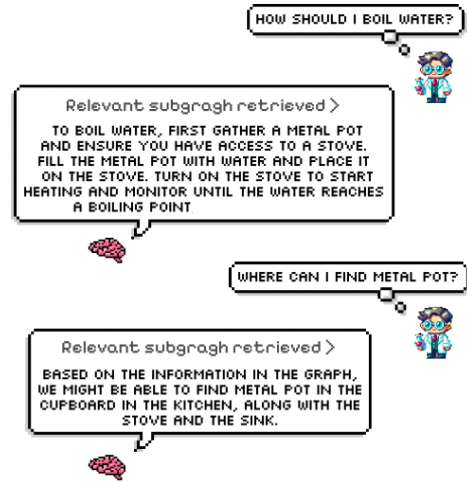


Figure 1: Visualization of DAVIS’s inner monologue during decision-making. The agent uses its World Model to retrieve relevant subgraphs from a Temporal Knowledge Graph (TKG) for reasoning.

2020), it often suffers from sample inefficiency, limited generalizability, and poor interpretability, making real-world deployment challenging (Dulac-Arnold et al., 2019). Recently, large language models (LLMs) (Radford et al.; Touvron et al., 2023) have revolutionized the creation of autonomous agents by leveraging natural language understanding to enhance interpretability and generalization. These LLM-based agents have shown great promise in critical domains such as healthcare (Qiu et al., 2024) and scientific research (Schmidgall et al., 2025) by mimicking human decision-making processes and enabling more intuitive reasoning and actions.

Several approaches have enhanced agentic reasoning and decision-making. SwiftSage (Lin et al., 2023) emulates the fast and slow thinking of humans with fine-tuned language models for planning. SayCan (Ahn et al., 2022) decomposes tasks into subgoals, while ReAct (Yao et al., 2023)²

²SwiftSage, Reflexion, SayCan, and ReAct are used under

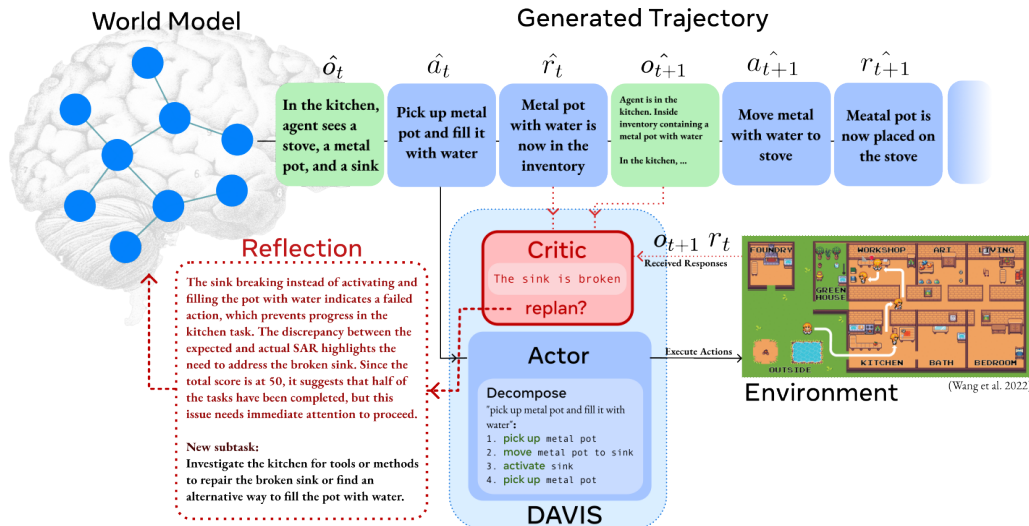


Figure 2: Overview of DAVIS’s decision-making process. The World Model generates a feasible course of actions, which are translated by the actor and executed sequentially by the agent in the environment. The critic detects discrepancies between expected and actual outcomes, identify failures, and suggest replanning.

integrates reasoning into execution. RAG-based systems like Reflexion (Shinn et al., 2023) and RAP³ (Kagaya et al., 2024) retrieve past experiences via semantic search, but their unstructured memory limits multi-hop reasoning and causal understanding. These systems retrieve static information rather than engaging in agentic, multi-turn retrieval, preventing dynamic adaptation.

Humans do not retrieve past knowledge statically; instead, we actively reflect, question our understanding, and refine our knowledge through internal dialogues. Inspired by this, we introduce DAVIS, an agentic multi-turn retrieval system that mirrors human cognition by enabling iterative interactions between the agent and its memory during the planning stage. DAVIS actively engages with its World Model (WM), a temporal knowledge graph-based QA system, to refine its understanding before execution. DAVIS engages in conversation with its WM to retrieve past experiences, evaluate actions, identify gaps, and optimize strategies.

DAVIS proves to be effective for iterative reasoning within scientific domains. Specifically, DAVIS outperforms 4 other baselines (Ahn et al., 2022; Kagaya et al., 2024; Yao et al., 2023; Shinn et al., 2023) on 8 out of 9 science subjects in the ScienceWorld (Wang et al., 2022) environment.⁴ DAVIS’s WM achieves competitive performance on the HotpotQA (Yang et al., 2018) and MusiqueQA (Trivedi

et al., 2022) dataset⁵. Our contributions can be summarized as follows:

- We introduce **DAVIS**, an agentic reasoning framework that leverages multi-turn retrieval and self-reflection to improve decision-making.
- Unlike static retrieval methods, DAVIS leverages a structured temporal knowledge graph memory system to enable multi-hop reasoning and causal understanding.
- Empirical evaluations show that DAVIS outperforms prior agentic reasoning models across scientific benchmarks, demonstrating superior planning and execution.

2 Background & Related Work

2.1 LLM agentic systems

Recent advancements in LLM-based agentic systems have drawn heavily from human decision-making processes and generally fall into two paradigms: direct interaction via chain-of-thought (CoT) reasoning or Retrieval-Augmented Generation (RAG).

The first paradigm involves agents interacting directly with their environment using CoT reasoning (Yao et al., 2023; Ahn et al., 2022; Lin et al., 2023). Chain-of-Thought prompting (Wei et al., 2023) enables large language models to decompose complex tasks into smaller, interpretable reasoning steps.

MIT license

³RAP is used under MIT license

⁴ScienceWorld is used under Apache 2.0 license

⁵The HotpotQA dataset is distributed under the CC BY-SA 4.0 license. The MusiqueQA dataset is distributed under CC BY 4.0 license

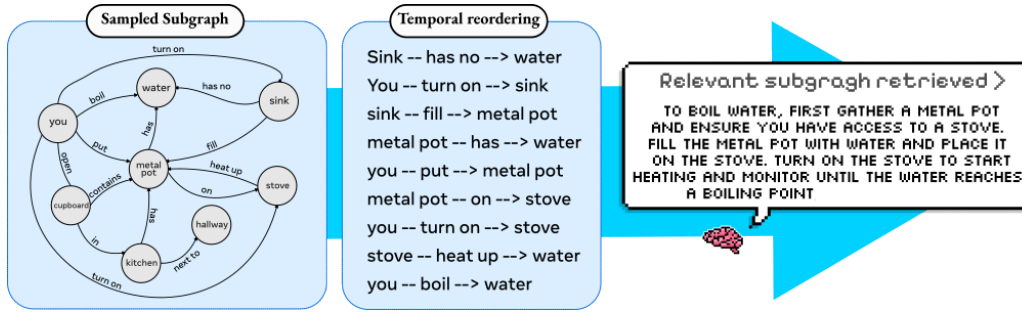


Figure 3: DAVIS’s retrieval and reasoning process. Left: subgraph with relevant entities and their relationships. Middle: temporal reordering of the retrieved information to establish a coherent sequence of actions. Right: DAVIS generates a structured and interpretable response.

117 However, CoT-based systems lack robust memory
 118 for long-term learning and adaptability across multiple
 119 tasks. The absence of memory has been linked
 120 to increased hallucination and stochasticity in task
 121 planning (Guerreiro et al., 2023), posing risks in
 122 domains like scientific research.

123 The second paradigm, RAG-based systems, inte-
 124 grates retrieval mechanisms with generative capa-
 125 bilities, enabling agents to access relevant external
 126 knowledge during task execution. In the Minecraft
 127 domain, extensive work has been done on RAG-
 128 based agents, with JARVIS-1 (Wang et al., 2023b)
 129 and Voyager (Wang et al., 2023a) representing the
 130 state-of-the-art. Since Minecraft is one of the most
 131 popular video games in the world, these agents
 132 leverage the extensive in-domain knowledge of
 133 LLMs but face significant limitations in scientific
 134 environments, where tasks often involve unknown
 135 skills and cannot rely on pre-existing knowledge.
 136 A more general and iterative approach involving
 137 multiple trials is necessary in such cases.

138 Reflexion (Shinn et al., 2023) and RAP (Kagaya
 139 et al., 2024) represent recent advances in agentic
 140 reasoning, using memory logs or semantic retrieval
 141 to guide decisions—Reflexion through reflective
 142 trial histories and RAP through nearest-neighbor
 143 search of past experiences. While these systems
 144 address some shortcomings of chain-of-thought
 145 (CoT) prompting, they rely heavily on unstructured
 146 vector databases, which scatter information and
 147 hinder multi-hop and causal reasoning. Addition-
 148 ally, they lack the capacity for temporal reasoning
 149 and iterative refinement. Also, neither approach in-
 150 corporates internal validation or model-based plan-
 151 ning, which limits their ability to make deliberate
 152 and accurate decisions. Thus the need for hybrid
 153 systems that combine structured memory, iterative
 154 retrieval, and internal planning. DAVIS addresses

155 this gap by integrating a temporal knowledge graph
 156 with agentic, multi-turn reasoning and critic-driven
 157 reflection, offering a more robust framework for
 158 complex scientific environments.

2.2 Graph Question Answering (Graph QA) 159

160 Graph Question Answering (Graph QA) sys-
 161 tems have become effective tools for structured
 162 reasoning and information retrieval. GraphReader
 163 (Li et al., 2024) constructs a graph from docu-
 164 ment chunks and deploys an agent for exploration.
 165 HOLMES (Panda et al., 2024) extracts relevant docu-
 166 ments, builds an entity-document graph, prunes it,
 167 and uses cosine similarity for answers. GraphRAG
 168 (Edge et al., 2024) generates an entity knowledge
 169 graph, pregenerates community summaries, and
 170 synthesizes responses. By encoding knowledge in
 171 a graph format, these systems excel at multi-hop
 172 reasoning over interconnected concepts, making
 173 them particularly valuable for domains that require
 174 relational understanding, such as scientific research.
 175 Unlike unstructured vector-based retrieval systems,
 176 Graph QA systems enable iterative retrieval, al-
 177 lowing agents to retrieve information, reason over
 178 it, and perform subsequent queries based on the
 179 refined context.

3 DAVIS 180

181 DAVIS adopts a model-based planning approach
 182 (Sutton and Barto, 1998), where the agent uses a
 183 World Model (WM) as an internal representation
 184 of its surrounding environment.

3.1 Problem Formulation 185

186 We define the planning problem for DAVIS in
 187 a textual environment as a Partially Observable
 188 Markov Decision Process (POMDP):

$$\mathcal{P} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \Omega, \mathcal{O}, \gamma) \quad 189$$

Algorithm 1 Planning with Retrieval-Augmented World Model

Input: τ, \mathcal{R} **Parameters:** L, k **Output:** τ

- 1: **for** $t = 1$ to L **do**
- 2: $\hat{s}_t \leftarrow f(\tau)$ ▷ State estimation
- 3: $\hat{a}_t \leftarrow \pi(\hat{s}_t, \mathcal{R}, k)$
- 4: $\tau \leftarrow \tau \cup \hat{a}_t$
- 5: $\hat{o}_{t+1}, \hat{r}_{t+1} \leftarrow \text{TRANSITION}(\hat{s}_t, \hat{a}_t)$ ▷ Algorithm 2
- 6: $\tau \leftarrow \tau \cup \{\hat{o}_{t+1}, \hat{r}_{t+1}\}$
- 7: **if** (τ) violates safety constraints (optional) **then**
- 8: Alert supervisor
- 9: **end if**
- 10: **end for**
- 11: **return** τ

In this formulation, \mathcal{S} denotes the set of true environment states, which are not directly observable. \mathcal{A} represents the set of available actions. $\mathcal{T}(s_{t+1} | s_t, a_t)$ is the state transition probability function, modeling the dynamics of the environment. $\mathcal{R}(s_t, a_t)$ is the reward function, specifying the immediate reward received after taking action a_t in state s_t . Ω is the set of possible observations. $\mathcal{O}(o_{t+1} | s_{t+1}, a_t)$ is the observation probability function, defining the likelihood of observing o_{t+1} given the new state s_{t+1} and action a_t . $\gamma \in [0, 1)$ is the discount factor, determining the present value of future rewards.

Since the true state s_t is not directly observable, the agent maintains a belief state b_t , which is a probability distribution over all possible states, representing the agent’s estimate of the environment’s state at time t . The belief state is updated based on the agent’s actions and received observations. The agent selects an action $a_t \in \mathcal{A}$ based on its current belief state, following a policy π :

$$a_t = \pi(b_t)$$

After executing the action a_t , the agent receives a reward $r_t = \mathcal{R}(s_t, a_t)$ and transitions to a new state s_{t+1} according to the transition function \mathcal{T} . The objective of the agent is to find an optimal policy π^* that maximizes the expected cumulative discounted reward over time:

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid \pi \right]$$

3.2 World Model (WM)

The World Model (WM) of DAVIS is represented as a Temporal Knowledge Graph (TKG), constructed through a combination of Stanford CoreNLP⁶ (Manning et al., 2014) for coreference

⁶We used default hyperparameters provided by the Stanza package for CoreNLP

resolution and LLM prompting for knowledge extraction. In textual environments, where state representations are conveyed in natural language, constructing an effective WM requires methods that can process and represent textual information efficiently and accurately.

State representation methods in text-based environments include text encoding techniques using recurrent neural networks (Narasimhan et al., 2015, He et al., 2016, Hausknecht et al., 2020), transformers (Kim et al., 2022), and knowledge graph (KG) representations (Ammanabrolu and Hausknecht, 2020). KGs offer structured and interpretable representations without requiring extensive training. Ammanabrolu and Riedl’s (2021) framed KG construction in text-based games as a question-answering problem, where agents identified objects and their attributes. This approach demonstrated that higher-quality KGs led to improved control policies. DAVIS extends this concept to Temporal Knowledge Graphs, incorporating time-sensitive information to model dynamic environment changes. Temporal reasoning is critical in such settings, and as noted in (Lee et al., 2023), LLMs are highly effective in extrapolating TKGs using in-context learning.

Let $G_t = (\mathcal{E}, \mathcal{R}, \mathcal{T})$ denote the Temporal Knowledge Graph (TKG) at time t , where \mathcal{E} is the set of entities at t , \mathcal{R} is the set of relations representing relationships between entities at, and \mathcal{T} is the set of timestamps associated with each relation e_i .

During training, when DAVIS executes an action a_t and receives the subsequent observation o_{t+1} , the transition is stored as:

$$(o_t \parallel a_t \parallel o_{t+1})$$

We prompted an LLM to summarize the concatenated transition and applied Stanford CoreNLP for coreference resolution. The resolved text is then analyzed to extract entities V_i and relations tuples using LLM-based parsing.

Each extracted tuple (v_i, e_j, v_k, τ) is added to the TKG, where the timestamp τ records the time at which the fact was introduced:

$$G_{t+1} = G_t \cup \{(v_i, e_j, v_k, \tau)\}$$

3.3 Retrieval-Augmented Model Approximation

As demonstrated in Lee et al.’s (2023), LLMs excel at recognizing temporal patterns and extrapolating future events based on past data. DAVIS leverages this capability to approximate future states

and rewards. For example, if sufficient past data indicates that opening a cupboard often reveals a kettle, the LLM can infer such transitions purely from learned patterns without requiring explicit pre-programmed rules. Unlike prior works (Kagaya et al., 2024; Shinn et al., 2023) that rely on vector-based retrieval of experiences, DAVIS employs a more agentic approach. Instead of passively retrieving information, DAVIS engages in a conversational process with its WM, iteratively querying to fill knowledge gaps while retrieving relevant subgraphs to generate informed responses. The retrieval system is described in Section 3.4.

Although true state s_t is not directly observable as mentioned in Section 3.1, it is theoretically possible to maintain a statistic $f(\tau)$ that approximates the belief state from the trajectory history. The statistic is updated recurrently, and captures all relevant information necessary for optimal decision-making (Nguyen et al., 2021; Åström, 1965). Applying this to DAVIS, we approximate the belief state \hat{b}_t with equation:

$$\hat{b}_t = f(\tau_{t':t}),$$

where $f(\cdot)$ is a prompted LLM that extracts relevant information from the trajectory history, and is updated recurrently with new observations and actions. To further refine decision-making, DAVIS maintains an inner monologue \mathcal{M}_t , a running list of iterative queries and answers exchanged between DAVIS and its WM, as illustrated in Figure 1. This monologue allows the system to dynamically update its WM based on retrieved insights.

DAVIS optimizes its policy while simultaneously learning approximations of the transition and reward models using its WM. The learned functions incorporating the inner monologue are:

$$\text{Policy: } \pi(a_t | \hat{b}_t, \mathcal{M}_t) \quad (1)$$

$$\text{Transition Model: } \hat{T}(o_{t+1} | \hat{b}_t, a_t, \mathcal{M}_t) \quad (2)$$

$$\text{Reward Model: } \hat{\mathcal{R}}(r_t | \hat{b}_t, a_t, \mathcal{M}_t) \quad (3)$$

With the approximated belief \hat{b}_t , DAVIS’s WM estimates the transition and reward models using prior experiences retrieved from a TKG. DAVIS leverages prior experiences to directly inform its policy as defined in Equation (1). This retrieval-driven approximation enables DAVIS to construct an adaptive and context-aware model of the world, allowing for informed decision-making in complex, temporally dependent environments.

Algorithm 2 Transition Prediction

Input: \hat{b}_t, \hat{a}_t, k **Output:** $\hat{o}_{t+1}, \hat{r}_{t+1}$
1: $\mathcal{M} \leftarrow \emptyset$ \triangleright Initialize inner monologue set
2: $i \leftarrow 0$
3: **while** $i < k$ or not predicted **do**
4: $\hat{o}_{t+1}, q \leftarrow \hat{T}(\hat{b}_t, \hat{a}_t, \mathcal{M})$
5: $\hat{r}_{t+1}, q \leftarrow \hat{R}(\hat{b}_t, \hat{a}_t, \mathcal{M})$
6: **if** $q \neq \emptyset$ **then**
7: $\mathcal{M} \leftarrow \mathcal{M} \cup \{(q, \text{graphQA}(q))\}$
8: **end if**
9: $i \leftarrow i + 1$
10: **end while**
11: **return** $\hat{o}_{t+1}, \hat{r}_{t+1}$

3.4 Retrieval System

Given a query q , such as “Where can I find water?”, the WM first narrows its search to relevant entity types such as Person (PER) and Location (LOC). It then selects the two most relevant entities from the available options. Limiting the scope to two entities is computationally efficient and ensures a manageable search space without sacrificing relevant context. The query is then expanded and processed as follows and illustrated in Figure 3:

1. **We iteratively expand** the current list of selected entities by adding their neighbors, forming a maximal subgraph as ignoring temporal information might result in an infeasible path.
2. **We reorder the edges** in the maximal subgraph based on timestamps. This reordering shows the proper sequence of events.
3. **The temporal sequence is then passed to an LLM** as in-context examples for extrapolation and summarization, enabling the LLM to generate a coherent response.

3.5 Planning and Execution with a WM

With the reward model and transition model approximated, we can now plan action trajectories within the WM. Algorithm 1 describes the WM-incorporated planning process of DAVIS.

For plan execution, we employ an actor-critic structure, consisting of two distinct models: the actor R_a and the critic R_c , integrated with the WM architecture. The process is illustrated in Figure 2. Below, we provide a formalized description of each model and its role within DAVIS.

World Model (WM). The primary objective of the WM is to generate a comprehensive plan or trajectory for achieving a specific task within the environment. Given an initial observation estimate

358 \hat{o}_t , the WM generates a predicted trajectory

$$359 \tau_{t:t+L} = \left\{ (\hat{o}_i, \hat{a}_i, \hat{o}_{i+1}, \hat{r}_{i+1}) \right\}_{i=t}^{t+L-1}$$

360 of length L . This trajectory $\tau_{t:t+L}$ is passed to the
361 actor-critic model for execution in the environment.

362 **Actor.** The actor R_a decomposes each high-level
363 action $\hat{a}_t \in \tau$ into executable commands within
364 the given environment domain. It also predicts
365 intermediate state transitions between actions:

$$366 \hat{r}_{t:t+L'} = R_a(\tau_{t:t+L})$$

367 where $L' \geq L$ accounts for the expanded trajec-
368 tory with executable low-level actions. The actor
369 model is prompted with permissible commands in
370 the current environment. After decomposition, the
371 expanded trajectory $\hat{r}_{t:t+L'}$ is executed step-by-step
372 in the environment, producing actual environment
373 responses:

$$374 (o_t, r_t, o_{t+1}) = \mathcal{E}(\hat{a}_t)$$

375 where \mathcal{E} is the environment transition function that
376 maps the executed action \hat{a}_t to the resulting obser-
377 vation o_{t+1} and reward r_t . These results are passed
378 to the critic model.

379 **Critic.** The critic R_c evaluates the actual execu-
380 tion results against the predicted trajectory τ . The
381 comparison is performed through an LLM-based
382 evaluation function, which assesses the semantic
383 consistency between the expected and actual obser-
384 vations. At each timestep t , the critic receives
385 the predicted state transition $(\hat{o}_t, \hat{r}_t, \hat{o}_{t+1})$ and the
386 actual environment response (o_t, r_t, o_{t+1}) obtained
387 from executing \hat{a}_t in the environment.

388 The LLM-based critic compares these compo-
389 nents via a prompted evaluation function R_c :

$$390 \Delta_t = R_c\left((\hat{o}_t, \hat{r}_t, \hat{o}_{t+1}), (o_t, r_t, o_{t+1})\right)$$

391 where Δ_t is a qualitative feedback score represent-
392 ing the level of agreement between the predicted
393 and actual transitions. Based on the LLM’s re-
394 sponse, the critic determines whether replanning
395 is necessary. If the predicted and actual observa-
396 tions deviate significantly, the critic updates the
397 reflection memory \mathfrak{R}_t and triggers replanning:

$$398 \mathfrak{R}_{t+1} = \mathfrak{R}_t \cup \{(o_t, \hat{s}_t, \Delta_t)\}$$

399 Algorithm 1 is then called to replan the new subtask.
400 For instance, if the task is "using the stove to heat

water" and the agent encounters an exception (e.g.,
the stove is broken), the LLM evaluates the excep-
tion, updates \mathcal{M}_t , and suggests a revised subtask
such as "find an alternative heating method."

4 Experiment

4.1 ScienceWorld environment

401 We selected ScienceWorld (Wang et al., 2022) as
402 the primary benchmark to evaluate DAVIS, as it is
403 currently the only environment designed for interac-
404 tive scientific reasoning. It features 30 tasks across
405 9 grade-school science subjects, set in a simulated
406 lab where agents must navigate 8 functional rooms
407 and use scientific tools to complete tasks. Each
408 task includes over 100 variations, some of which
409 significantly alter the setup by masking rooms or
410 removing equipment—requiring strong generaliza-
411 tion and adaptability. The environment demands
412 common sense reasoning, deduction, and procedu-
413 ral knowledge. Scores reflect progress toward task
414 completion (e.g., 75 indicates 75% progress be-
415 fore failure), enabling structured and interpretable
416 evaluation. Full details are in Appendix A.

4.2 Performance

423 We evaluated DAVIS on the ScienceWorld
424 benchmark, comparing its performance against
425 state-of-the-art baseline agents: SayCan, ReAct,
426 Reflexion, and RAP. Baselines were selected based
427 on their competitive performance, available imple-
428 mentations, and relevance to ScienceWorld. The
429 current state-of-the-art method on ScienceWorld,
430 SwiftSage (Lin et al., 2023), was excluded from
431 our replication baselines because discrepancies be-
432 tween the available code and the documented eval-
433 uation methods made direct replication infeasible.
434 For consistency, all baselines were reimplemented
435 to align with the latest ScienceWorld version. For
436 fairness, both RAP and DAVIS utilized memory
437 constructed from five episodes of golden trajec-
438 tories rather than the ReAct-based approach proposed
439 in Kagaya et al.’s (2024). Performance was aver-
440 aged across subjects for comparison, with details
441 on tasks and subjects provided in Tables 4 and 6 in
442 the appendix. Figure 4 shows DAVIS outperform-
443 ing all baselines in 8 out of 9 subjects, achieving an
444 overall average score of 65.06—approximately 1.8
445 times higher than competing methods. Full results
446 for each task, including standard deviations, are in
447 the appendix Table 7.

448 Overall, DAVIS took fewer steps before converg-
449 ing to the final score when compared to SayCan,
450

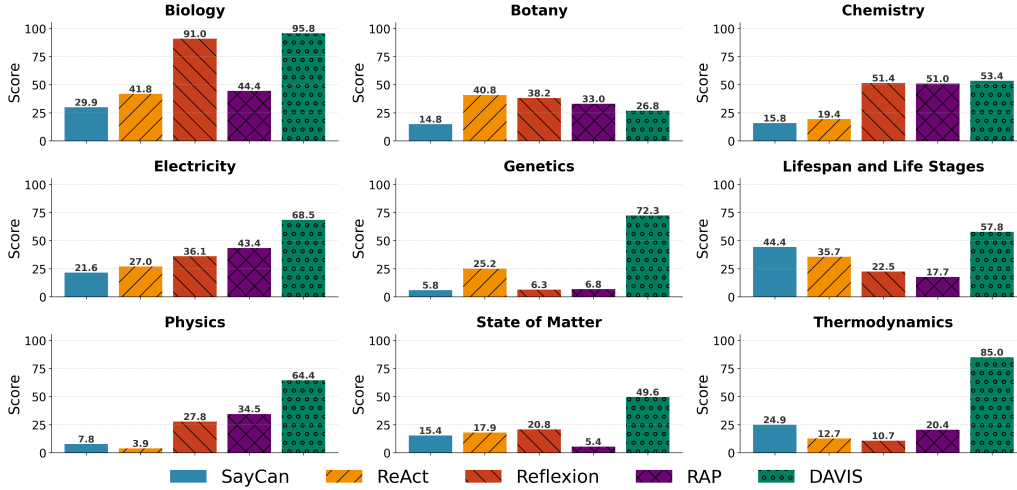


Figure 4: Performance comparison of different agents (SayCan, ReAct, Reflexion, RAP, and DAVIS) across multiple scientific domains. For full results, view table 7 in the Appendix.

| Task | D | D + W |
|------------------------------------|-------|--------|
| Long Tasks | | |
| Melt (1-2) | 3.00 | 70.00 |
| Determine Melting Point Unk. (2-3) | 5.00 | 92.33 |
| Medium Tasks | | |
| Mix Paint Secondary (6-1) | 40.00 | 36.37 |
| Test Conductivity (3-3) | 55.00 | 58.33 |
| Short Tasks | | |
| Lifespan Longest-Lived (7-1) | 66.67 | 100.00 |
| Find Living Thing (4-1) | 25.00 | 100.00 |

Table 1: DAVIS performance with (D + W) and without WM (D)

Table 2: Ablation study: Full model (D+W), w/o Actor (D-A), and w/o Critic (D-C).

| Type | Task | D+W | D-A | D-C |
|--------|------|-------------|-------------|-------------|
| Long | 1-2 | 70 (4.38) | 25 (1.12) | 23.3 (3.20) |
| | 2-3 | 92.3 (1.51) | 79.7 (1.19) | 33.3 (1.29) |
| Medium | 6-1 | 36.4 (3.22) | 100 (1.49) | 86 (2.49) |
| | 3-3 | 58.3 (2.71) | 28 (1.32) | 49.3 (1.72) |
| Short | 7-1 | 83.3 (2.00) | 66.7 (2.00) | 83.3 (2.00) |
| | 4-1 | 100 (2.50) | 44.7 (2.94) | 25 (2.50) |

ReAct, and Reflexion. Compared to RAP, DAVIS was better at transferring knowledge from its training to execution despite differences among variations of the same task. Its World Model allows for multi-hop reasoning and inferences based on past training data.

4.3 Ablation Study

We systematically evaluate the system with individual modules removed—specifically the World Model (WM), Actor, and Critic—and compare their impact on performance using two metrics: (1) average task score (i.e., task progress before timeout) and (2) average number of steps per replanning cycle (*steps/replan*). To ensure representative coverage across task complexity, we select two tasks from each category of task length (short, medium, and long). For each task, results are averaged over three different environment variations.

World Model: Table 1 compares DAVIS with and without its World Model. The WM serves as structured external memory in the form of a tempo-

ral knowledge graph, enabling grounded and long-horizon planning. Without this, the agent has to rely solely on the internal capabilities of the LLM, lacking access to temporal or multi-hop context.

Using the WM consistently improves performance across all tasks, particularly in complex and temporally grounded settings like *Melt* and *Find Living Thing*. This supports our claim that temporal and structured grounding is critical for high-fidelity decision-making in scientific domains.

Actor-Critic. We ablate DAVIS by removing the Actor or Critic module individually. In the no-Actor setup, the World Model directly outputs executable actions, skipping high-level goal decomposition. In the no-Critic setup, reflection and subtask updates are disabled, though replanning triggers remain. Table 2 summarizes the results.

- **No Actor:** The agent struggles to produce valid commands despite access to action formats, resulting in low task scores and near-1.0 steps/replan—indicating constant replanning and poor multi-step coherence. We observed that the

performance in task 6-1 was due to luck, as the agent correctly guessed the critical action ‘focus on’ within the first 5 steps, skipping all the intermediate steps. This behavior was also observed in the D-C model.

- **No Critic:** The agent is able to execute longer action chains without an Actor module, but its lack of introspective feedback (due to the absence of a Critic) limits its ability to recover from errors. While performance differences are minimal on short tasks—typically solvable within one or two replanning cycles—the gap widens on longer tasks that require more adaptive reasoning. Compared to the no-Actor condition, both task performance and steps per replan improve, but remain below those of the full DAVIS system.

The Actor enables structured execution, and the Critic enhances adaptivity. A higher average steps per replan ratio, paired with strong task scores, demonstrates coherent, cost-efficient planning.

4.4 Multi-hop Q&A

We evaluated the performance of DAVIS’s World Model (WM) on the multi-hop QA benchmarks HotpotQA and MusiqueQA using 400 randomly sampled instances, following the evaluation protocol of Li et al.’s (2024). As shown in Table 3, DAVIS (GPT-4o) achieves strong results, surpassing GraphReader and GraphRAG on HotpotQA with an F1 score of 73.8 and a competitive EM of 56.25—approaching the state-of-the-art HOLMES. On MusiqueQA, DAVIS maintains strong performance (F1: 48.5, EM: 33.8), further demonstrating the effectiveness of its structured, temporal memory in reasoning tasks. While HOLMES achieves the highest overall scores, its static hyper-relational graph architecture lacks DAVIS’s ability to support dynamic updates during inference, which is crucial for agents operating in evolving or interactive environments. For each system, we report results using the best-performing language model configuration as documented in the respective original papers. We observe that retrieval-based systems are highly sensitive to the underlying LLM: DAVIS performs better with GPT-4o than with GPT-4-turbo, despite the latter’s generally stronger performance claims. This LLM sensitivity, though observed qualitatively, warrants further study; a systematic analysis of model-architecture alignment is left for future work as it lies beyond the scope of this paper.

Table 3: WM comparison against SotA baselines.

| Method | HotpotQA | | MusiqueQA | |
|------------------------|-------------|-------------|-------------|-------------|
| | EM | F1 | EM | F1 |
| GPT-4o | 46.3 | 64.1 | 19.0 | 34.4 |
| GPT-4-turbo | 44.3 | 60.4 | 20.5 | 34.7 |
| GraphReader (GPT-4) | 55.0 | 70.0 | 38.0 | 47.4 |
| HOLMES (GPT-4) | 66.0 | 78.0 | 48.0 | 58.0 |
| GraphRAG (GPT-4o-mini) | 58.7 | 63.3 | 40.0 | 53.5 |
| DAVIS (GPT-4o) | 56.25 | 73.8 | 33.8 | 48.5 |
| DAVIS (GPT-4-turbo) | 55.25 | 71.0 | 34.0 | 47.1 |

5 Conclusion

DAVIS is an agent designed for scientific interactive reasoning tasks in complex environments. DAVIS represents a novel approach that leverages a structured World Model (WM) in the form of a temporal knowledge graph, enabling iterative retrieval and reasoning over past experiences. This structured representation allows DAVIS to approximate both the transition dynamics and reward models of its environment, facilitating more informed decision-making. DAVIS also uniquely uses an interactive retrieval process, which combines iterative querying with contextual reasoning to fill knowledge gaps and refine understanding. This is augmented by DAVIS’s ability to perform internal planning and validation before interacting with the environment. By engaging in pre-execution deliberation, DAVIS enables clearer inspection of its planned actions, making it easier for human supervisors to review its decision-making process. This transparency facilitates stronger safeguards compared to reinforcement learning (RL) agents, whose policies are often opaque. DAVIS is ideal for scientific tasks that demand precision, adaptability, and strict adherence to experimental protocols.

Evaluations across several scientific domains, including thermodynamics, biology, and physics, demonstrate the efficacy of DAVIS’s structured knowledge representation and retrieval methods. DAVIS significantly outperforms baseline agents by combining robust planning with the capacity for iterative reasoning, enabling it to generalize effectively from demonstrations to new tasks.

6 Limitations

While DAVIS demonstrates strong reasoning capabilities and improved performance over previous agentic approaches, it has several limitations that we will address in future research.

581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629

6.1 High operational cost

DAVIS heavily relies on Large Language Models (LLMs), making it computationally expensive. Due to its careful planning and reasoning process, it sends and receives an average of 43,000 tokens per action, resulting in an estimated cost of \$0.43 per action. For tasks requiring 100 actions, this cost can escalate to \$43 per episode, leading to a total experimental cost of approximately \$3,000 for 90 variations.

6.2 Sensitive to LLM performance

DAVIS’s reasoning and decision-making abilities fluctuate based on the underlying LLM’s performance. Factors such as model version updates, prompt engineering quality, and external API changes can lead to accuracy, consistency, and response time variability. This dependence on LLM stability makes DAVIS susceptible to unexpected performance shifts, which may impact reliability in dynamic or evolving environments.

6.3 Biased Planning & Knowledge Dependence

DAVIS’s decision-making process is heavily influenced by the Temporal Knowledge Graph (TKG), which serves as its structured memory. However, this dependence can lead to biased planning, as DAVIS prioritizes information within the graph. Although efforts were made to increase data diversity by populating the knowledge graph with 150 different ScienceWorld task variations, the model still struggles when encountering novel scenarios or incomplete knowledge. Future work should explore adaptive knowledge integration to mitigate bias.

6.4 Lack of Multimodal Capabilities

DAVIS operates exclusively in textual environments, limiting its applicability as an embodied agent. The absence of visual, auditory, or sensory perception restricts its ability to interact with real-world multimodal tasks. Future research should focus on integrating visual and sensor-based input processing to enhance generalization and deployment in robotic or multimodal AI systems.

References

Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang,

Rosario Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil J. Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy Zeng. 2022. *Do As I Can, Not As I Say: Grounding Language in Robotic Affordances*. *arXiv preprint*. ArXiv:2204.01691 [cs]. 630-640

Prithviraj Ammanabrolu and Matthew Hausknecht. 2020. *Graph Constrained Reinforcement Learning for Natural Language Action Spaces*. *arXiv preprint*. ArXiv:2001.08837 [cs, stat]. 641-644

Prithviraj Ammanabrolu and Mark O. Riedl. 2021. *Learning Knowledge Graph-based World Models of Textual Environments*. *arXiv preprint*. ArXiv:2106.09608 [cs]. 645-648

Gabriel Dulac-Arnold, Daniel Mankowitz, and Todd Hester. 2019. *Challenges of Real-World Reinforcement Learning*. *arXiv preprint*. ArXiv:1904.12901 [cs]. 649-652

Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. *From Local to Global: A Graph RAG Approach to Query-Focused Summarization*. *arXiv preprint*. ArXiv:2404.16130 [cs]. 653-657

Nuno M. Guerreiro, Duarte Alves, Jonas Waldendorf, Barry Haddow, Alexandra Birch, Pierre Colombo, and André F. T. Martins. 2023. *Hallucinations in Large Multilingual Translation Models*. *arXiv preprint*. ArXiv:2303.16104 [cs]. 658-662

Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. 2020. *Dream to Control: Learning Behaviors by Latent Imagination*. *arXiv preprint*. ArXiv:1912.01603 [cs]. 663-666

Matthew Hausknecht, Prithviraj Ammanabrolu, Marc-Alexandre Côté, and Kingdi Yuan. 2020. *Interactive Fiction Games: A Colossal Adventure*. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):7903–7910. Number: 05. 667-671

Ji He, Jianshu Chen, Xiaodong He, Jianfeng Gao, Lihong Li, Li Deng, and Mari Ostendorf. 2016. *Deep Reinforcement Learning with a Natural Language Action Space*. *arXiv preprint*. ArXiv:1511.04636 [cs]. 672-676

Tomoyuki Kagaya, Thong Jing Yuan, Yuxuan Lou, Jayashree Karlekar, Sugiri Pranata, Akira Kinose, Koki Oguri, Felix Wick, and Yang You. 2024. *RAP: Retrieval-Augmented Planning with Contextual Memory for Multimodal LLM Agents*. *arXiv preprint*. ArXiv:2402.03610 [cs]. 677-682

Minsoo Kim, Yeonjoon Jung, Dohyeon Lee, and Seungwon Hwang. 2022. *PLM-based World Models for* 683-684

| | | |
|-----|--|---|
| 685 | Text-based Games . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 1324–1341, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. | |
| 686 | | |
| 687 | | |
| 688 | | |
| 689 | | |
| 690 | Dong-Ho Lee, Kian Ahrabian, Woojeong Jin, Fred Morstatter, and Jay Pujara. 2023. Temporal Knowledge Graph Forecasting Without Knowledge Using In-Context Learning . <i>arXiv preprint</i> . ArXiv:2305.10613 [cs]. | |
| 691 | | |
| 692 | | |
| 693 | | |
| 694 | | |
| 695 | Shilong Li, Yancheng He, Hangyu Guo, Xingyuan Bu, Ge Bai, Jie Liu, Jiaheng Liu, Xingwei Qu, Yangguang Li, Wanli Ouyang, Wenbo Su, and Bo Zheng. 2024. GraphReader: Building Graph-based Agent to Enhance Long-Context Abilities of Large Language Models . <i>arXiv preprint</i> . ArXiv:2406.14550 [cs]. | |
| 696 | | |
| 697 | | |
| 698 | | |
| 699 | | |
| 700 | | |
| 701 | Bill Yuchen Lin, Yicheng Fu, Karina Yang, Faeze Brahman, Shiyu Huang, Chandra Bhagavatula, Prithviraj Ammanabrolu, Yejin Choi, and Xiang Ren. 2023. SwiftSage: A Generative Agent with Fast and Slow Thinking for Complex Interactive Tasks . <i>arXiv preprint</i> . ArXiv:2305.17390 [cs]. | |
| 702 | | |
| 703 | | |
| 704 | | |
| 705 | | |
| 706 | | |
| 707 | Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. 2014. The Stanford CoreNLP Natural Language Processing Toolkit . In <i>Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> , pages 55–60, Baltimore, Maryland. Association for Computational Linguistics. | |
| 708 | | |
| 709 | | |
| 710 | | |
| 711 | | |
| 712 | | |
| 713 | | |
| 714 | | |
| 715 | Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. 2013. Playing Atari with Deep Reinforcement Learning . <i>arXiv preprint</i> . ArXiv:1312.5602 [cs] version: 1. | |
| 716 | | |
| 717 | | |
| 718 | | |
| 719 | | |
| 720 | Karthik Narasimhan, Tejas Kulkarni, and Regina Barzilay. 2015. Language Understanding for Text-based Games using Deep Reinforcement Learning . In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 1–11, Lisbon, Portugal. Association for Computational Linguistics. | |
| 721 | | |
| 722 | | |
| 723 | | |
| 724 | | |
| 725 | | |
| 726 | | |
| 727 | Hai Nguyen, Brett Daley, Xinchao Song, Christopher Amato, and Robert Platt. 2021. Belief-Grounded Networks for Accelerated Robot Learning under Partial Observability . <i>arXiv preprint</i> . ArXiv:2010.09170 [cs]. | |
| 728 | | |
| 729 | | |
| 730 | | |
| 731 | | |
| 732 | Pranoy Panda, Ankush Agarwal, Chaitanya Devaguptapu, Manohar Kaul, and Prathosh Ap. 2024. HOLMES: Hyper-Relational Knowledge Graphs for Multi-hop Question Answering using LLMs . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13263–13282, Bangkok, Thailand. Association for Computational Linguistics. | |
| 733 | | |
| 734 | | |
| 735 | | |
| 736 | | |
| 737 | | |
| 738 | | |
| 739 | | |
| 740 | Pranav Putta, Edmund Mills, Naman Garg, Sumeet Motwani, Chelsea Finn, Divyansh Garg, and Rafael | |
| 741 | | |
| | Rafailov. 2024. Agent Q: Advanced Reasoning and Learning for Autonomous AI Agents . <i>arXiv preprint</i> . ArXiv:2408.07199 [cs]. | 742 743 744 |
| | Jianing Qiu, Kyle Lam, Guohao Li, Amish Acharya, Tien Yin Wong, Ara Darzi, Wu Yuan, and Eric J. Topol. 2024. LLM-based agentic systems in medicine and healthcare . <i>Nature Machine Intelligence</i> , 6(12):1418–1420. Publisher: Nature Publishing Group. | 745 746 747 748 749 750 |
| | Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving Language Understanding by Generative Pre-Training. | 751 752 753 |
| | Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu, Zicheng Liu, and Emad Barsoum. 2025. Agent Laboratory: Using LLM Agents as Research Assistants . <i>arXiv preprint</i> . ArXiv:2501.04227 [cs]. | 754 755 756 757 758 |
| | Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, Timothy Lillicrap, and David Silver. 2020. Mastering Atari, Go, chess and shogi by planning with a learned model . <i>Nature</i> , 588(7839):604–609. Publisher: Nature Publishing Group. | 759 760 761 762 763 764 765 |
| | Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language Agents with Verbal Reinforcement Learning . <i>arXiv preprint</i> . ArXiv:2303.11366 [cs]. | 766 767 768 769 770 |
| | Richard S Sutton and Andrew G Barto. 1998. Reinforcement Learning: An Introduction. | 771 772 |
| | Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models . <i>arXiv preprint</i> . ArXiv:2302.13971 [cs]. | 773 774 775 776 777 778 779 |
| | Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. MuSiQue: Multi-hop Questions via Single-hop Question Composition . <i>arXiv preprint</i> . ArXiv:2108.00573 [cs]. | 780 781 782 783 |
| | Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023a. Voyager: An Open-Ended Embodied Agent with Large Language Models . <i>arXiv preprint</i> . ArXiv:2305.16291 [cs]. | 784 785 786 787 788 |
| | Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. 2022. ScienceWorld: Is your Agent Smarter than a 5th Grader? <i>arXiv preprint</i> . ArXiv:2203.07540 [cs]. | 789 790 791 792 |
| | Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, and Yitao Liang. 2024. Describe, Explain, Plan and Select: Interactive Planning with Large Language Models Enables Open-World Multi-Task Agents . <i>arXiv preprint</i> . ArXiv:2302.01560 [cs]. | 793 794 795 796 797 |

Zihao Wang, Shaofei Cai, Anji Liu, Yonggang Jin, Jinbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng He, Zilong Zheng, Yaodong Yang, Xiaojian Ma, and Yitao Liang. 2023b. [JARVIS-1: Open-World Multi-task Agents with Memory-Augmented Multimodal Language Models](#). *arXiv preprint*. ArXiv:2311.05997 [cs].

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. [Chain-of-Thought Prompting Elicits Reasoning in Large Language Models](#). *arXiv preprint*. ArXiv:2201.11903 [cs].

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. [HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering](#). *arXiv preprint*. ArXiv:1809.09600 [cs].

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. [ReAct: Synergizing Reasoning and Acting in Language Models](#). *arXiv preprint*. ArXiv:2210.03629 [cs].

Haiteng Zhao, Chang Ma, Guoyin Wang, Jing Su, Lingpeng Kong, Jingjing Xu, Zhi-Hong Deng, and Hongxia Yang. 2024. [Empowering Large Language Model Agents through Action Learning](#). *arXiv preprint*. ArXiv:2402.15809 [cs].

K. J Åström. 1965. [Optimal control of Markov processes with incomplete state information](#). *Journal of Mathematical Analysis and Applications*, 10(1):174–205.

A ScienceWorld

ScienceWorld (Wang et al., 2022) is a benchmark designed to evaluate interactive reasoning in digital agents through a realistic laboratory simulation. Developed by the Allen Institute for AI, it provides a text-based environment that emulates scientific experiments, requiring agents to interact with objects, collect observations, and apply reasoning skills to solve tasks. The framework consists of approximately 40,000 lines of SCALA code with a PYTHON interface, following standard RL benchmarking practices.

The ScienceWorld environment consists of 10 interconnected locations (Fig. 5), each populated with up to 200 distinct object types, including scientific instruments, electrical components, biological specimens, substances, and common environmental elements like furniture and books. Agents can interact with objects through a predefined action space of 25 high-level actions, categorized into

| # | Task |
|------|---|
| 1-1 | Changes of State (Boiling) |
| 1-2 | Changes of State (Melting) |
| 1-3 | Changes of State (Freezing) |
| 1-4 | Changes of State (Any) |
| 2-1 | Use Thermometer |
| 2-2 | Measuring Boiling Point (Known) |
| 2-3 | Measuring Boiling Point (Unknown) |
| 3-1 | Create a Circuit |
| 3-2 | Renewable vs Non-Renewable Energy |
| 3-3 | Test Conductivity (Known) |
| 3-4 | Test Conductivity (Unknown) |
| 4-1 | Find a Living Thing |
| 4-2 | Find a Non-Living Thing |
| 4-3 | Find a Plant |
| 4-4 | Find an Animal |
| 5-1 | Grow a Plant |
| 5-2 | Grow a Fruit |
| 6-1 | Mixing (Generic) |
| 6-2 | Mixing Paints (Secondary Colours) |
| 6-3 | Mixing Paints (Tertiary Colours) |
| 7-1 | Identify Longest-Lived Animal |
| 7-2 | Identify Shortest-Lived Animal |
| 7-3 | Identify Longest-Then-Shortest-Lived Animal |
| 8-1 | Identify Life Stages (Plant) |
| 8-2 | Identify Life Stages (Animal) |
| 9-1 | Inclined Planes (Determine Angle) |
| 9-2 | Friction (Known Surfaces) |
| 9-3 | Friction (Unknown Surfaces) |
| 10-1 | Mendelian Genetics (Known Plants) |
| 10-2 | Mendelian Genetics (Unknown Plants) |

Table 4: Tasks in ScienceWorld.

domain-specific operations (e.g., using a thermometer, measuring conductivity) and general interactions (e.g., moving, opening containers, picking up items). At each step, approximately 200,000 possible action-object combinations exist, though only a subset is relevant based on the context.

ScienceWorld tasks are designed to assess scientific reasoning across multiple disciplines. The dataset includes 30 distinct tasks (Table 4), covering a range of experimental procedures and problem-solving scenarios. These tasks are further grouped into 9 science domains (Table 6), including physics, chemistry, biology, and environmental science, allowing for targeted evaluation of an agent’s ability to reason through various scientific concepts, making ScienceWorld a robust benchmark for testing multi-step reasoning in dynamic, interactive environments.

B DAVIS Implementation Details

We utilized GPT-4-turbo for reasoning, GPT-4o for question answering, and LLaMA3-70B for the Knowledge Graph construction pipeline. Agents were run for a maximum of 80 steps per task. All

| Hyperparameter | Value |
|-----------------------------|---------------------|
| Maximum Steps per Task | 100 |
| Simplification Level | Easy |
| Knowledge Graph Pipeline | LLaMA3-70B-Instruct |
| Reasoning Model | GPT-4-Turbo |
| Maximum QA Turns | 5 |
| Predicted Trajectory Length | 5 |

Table 5: Hyperparameter settings for DAVIS.

873 RAG-based agents were initialized with five vari-
874 ations, a total of 150 variations, of rollouts using
875 the golden trajectory for training, while three ran-
876 domly sampled test variations, a total of 90 varia-
877 tions, were drawn from the ScienceWorld test set.
878 In contrast, all CoT agents were evaluated directly
879 on the randomly drawn test set as intended.

880 All experiments were conducted on a system
881 equipped with a NVIDIA RTX 3060 GPU, an
882 AMD Ryzen 9 7900X CPU, 64GB RAM, running
883 Ubuntu 23.04 with Python 3.11.0. The full table
884 of hyperparameters and settings for DAVIS is pro-
885 vided in Table 5. Full results is available in table 7,
886 and all code and prompts are available in the at-
887 tached repository.



Figure 5: The ScienceWorld environment

| Subject | Description | Tasks |
|--------------------------|--|-------------------------|
| Matter | Agents perform experiments to change the state of various materials, such as transforming ice to water or water to steam | 1-1, 1-2, 1-3, 1-4 |
| Thermodynamics | Agents conduct experiments involving temperature manipulation, such as heating or cooling objects. | 2-1, 2-2, 2-3 |
| Electricity | Agents relocate to a workshop and construct electrical circuits to achieve specific objectives. | 3-1, 3-2, 3-3, 3-4 |
| Biology | Agents relocate to a garden and identify animals based on various queries. | 4-1, 4-2, 4-3, 4-4 |
| Botany | Agents relocate to a greenhouse and perform tasks such as growing plants or observing their growth. | 5-1, 5-2 |
| Chemistry | Agents engage in standard chemistry tasks, such as mixing substances to create new compounds | 6-1, 6-2, 6-3 |
| Lifespan and Life Stages | Agents observe and report the life stages of plants and animals, such as germination, flowering, or molting. | 7-1, 7-2, 7-3, 8-1, 8-2 |
| Physics | Agents use physics knowledge to measure angles or explore physical properties of materials | 9-1, 9-2, 9-3 |
| Genetics | Agents identify genetic traits of plants, such as dominant or recessive characteristics, based on observations. | 10-1, 10-2 |

Table 6: Description of subjects and corresponding tasks in ScienceWorld. Each subject represents a unique domain of inquiry, with tasks designed to evaluate agents' reasoning, planning, and execution capabilities in diverse scientific scenarios.

| Task | SayCan | | ReAct | | Reflexion | | RAP | | DAVIS | |
|---------------------------------|--------|------|-------|------|-----------|------|-------|-------|--------|------|
| | Mean | Std | Mean | Std | Mean | Std | Mean | Std | Mean | Std |
| State of Matter | 15.42 | | 17.92 | | 20.83 | | 5.42 | | 49.58 | |
| 1-1 (L) | 1.67 | 1.5 | 2.67 | 2.5 | 27.67 | 41.0 | 13.33 | 20.55 | 25.67 | 19.6 |
| 1-2 (L) | 23.33 | 40.4 | 25.67 | 40.2 | 1.00 | 1.7 | 1.67 | 2.89 | 70.00 | 0.0 |
| 1-3 (L) | 3.33 | 5.8 | 19.33 | 25.3 | 19.33 | 25.3 | 6.67 | 5.78 | 32.00 | 27.7 |
| 1-4 (L) | 33.33 | 57.7 | 24.00 | 39.0 | 35.33 | 56.0 | 0.00 | 0.00 | 70.67 | 0.6 |
| Thermodynamics | 24.89 | | 12.67 | | 10.67 | | 20.44 | | 85.00 | |
| 2-1 (M) | 6.00 | 3.0 | 4.00 | 3.5 | 9.00 | 0.0 | 30.33 | 47.43 | 83.00 | 29.4 |
| 2-2 (M) | 7.67 | 0.6 | 6.33 | 0.6 | 17.33 | 18.8 | 8.67 | 15.02 | 79.67 | 35.2 |
| 2-3 (L) | 61.00 | 48.3 | 27.67 | 39.3 | 5.67 | 0.6 | 22.33 | 20.40 | 92.33 | 13.3 |
| Electricity | 21.58 | | 27.00 | | 36.08 | | 43.42 | | 68.50 | |
| 3-1 (S) | 30.33 | 40.4 | 30.33 | 40.4 | 23.33 | 34.5 | 39.00 | 33.05 | 82.33 | 15.7 |
| 3-2 (M) | 22.67 | 26.4 | 19.33 | 29.3 | 14.33 | 20.6 | 35.33 | 27.31 | 68.67 | 27.1 |
| 3-3 (M) | 23.33 | 27.5 | 5.00 | 5.0 | 39.00 | 34.5 | 38.00 | 35.03 | 58.33 | 2.9 |
| Biology | 29.92 | | 41.83 | | 91.00 | | 44.42 | | 95.83 | |
| 4-1 (S) | 11.33 | 9.8 | 17.00 | 0.0 | 72.33 | 47.9 | 61.00 | 38.1 | 100.00 | 0.0 |
| 4-2 (S) | 36.00 | 34.8 | 58.33 | 28.9 | 100.00 | 0.0 | 19.33 | 9.8 | 83.33 | 14.4 |
| 4-3 (S) | 22.33 | 4.6 | 75.00 | 0.0 | 91.67 | 14.4 | 58.33 | 36.0 | 100.00 | 0.0 |
| 4-4 (S) | 50.00 | 43.3 | 17.00 | 0.0 | 100.00 | 0.0 | 39.00 | 38.1 | 100.00 | 0.0 |
| Botany | 14.83 | | 40.83 | | 38.17 | | 33.00 | | 26.83 | |
| 5-1 (L) | 16.67 | 14.4 | 9.00 | 3.6 | 3.67 | 4.6 | 50.00 | 73.99 | 35.67 | 2.9 |
| 5-2 (L) | 13.00 | 4.6 | 72.67 | 47.3 | 72.67 | 47.3 | 16.00 | 13.89 | 18.00 | 6.2 |
| Chemistry | 15.78 | | 19.44 | | 51.44 | | 51.00 | | 53.44 | |
| 6-1 (M) | 16.67 | 11.5 | 23.33 | 11.5 | 56.67 | 37.9 | 53.33 | 5.78 | 36.67 | 5.8 |
| 6-2 (S) | 26.33 | 2.3 | 20.67 | 18.0 | 83.33 | 28.9 | 22.67 | 21.60 | 53.67 | 40.5 |
| 6-3 (M) | 4.33 | 2.3 | 14.33 | 5.1 | 14.33 | 7.5 | 77.00 | 0.00 | 70.00 | 0.0 |
| Lifespan and Life Stages | 44.40 | | 35.67 | | 22.47 | | 17.67 | | 57.80 | |
| 7-1 (S) | 75.00 | 43.3 | 66.67 | 28.9 | 50.00 | 0.0 | 16.67 | 28.86 | 100.00 | 0.0 |
| 7-2 (S) | 83.33 | 28.9 | 66.67 | 28.9 | 33.33 | 14.4 | 16.67 | 28.86 | 83.33 | 28.9 |
| 7-3 (S) | 33.00 | 0.0 | 22.00 | 19.1 | 22.33 | 9.2 | 5.67 | 9.81 | 83.00 | 0.0 |
| 8-1 (S) | 13.33 | 6.1 | 15.00 | 22.6 | 4.00 | 4.0 | 38.00 | 25.98 | 2.67 | 2.3 |
| 8-2 (S) | 17.33 | 4.6 | 8.00 | 0.0 | 2.67 | 4.6 | 11.33 | 9.81 | 20.00 | 0.0 |
| Physics | 7.78 | | 3.89 | | 27.78 | | 34.48 | | 64.44 | |
| 9-1 (L) | 5.00 | 5.0 | 0.00 | 0.0 | 36.67 | 54.8 | 30.00 | 30.00 | 76.67 | 40.4 |
| 9-2 (L) | 6.67 | 7.6 | 11.67 | 12.6 | 8.33 | 2.9 | 30.00 | 0.00 | 60.00 | 34.6 |
| 9-3 (L) | 11.67 | 16.1 | 0.00 | 0.0 | 38.33 | 53.5 | 43.44 | 23.28 | 56.67 | 37.9 |
| Genetics | 5.83 | | 25.17 | | 6.33 | | 6.83 | | 72.33 | |
| 10-1 (L) | 6.00 | 9.5 | 39.00 | 53.5 | 6.33 | 9.2 | 3.33 | 5.78 | 100.00 | 0.0 |
| 10-2 (L) | 5.67 | 9.8 | 11.33 | 9.8 | 6.33 | 9.2 | 10.33 | 10.50 | 44.67 | 47.9 |

Table 7: Full results on ScienceWorld. The average score for each category is displayed in the grey bar on the same row as the category label.