

Optimus-1 : Hybrid Multimodal Memory Empowered Agents Excel in Long-Horizon Tasks

Zaijing Li^{1,2}, Yuquan Xie¹, Rui Shao^{1*}, Gongwei Chen¹, Dongmei Jiang², Liqiang Nie^{1*}

¹Harbin Institute of Technology, Shenzhen

²Peng Cheng Laboratory

{lzj14011, xieyuquan20016, rshaojimmy, nieliqiang}@gmail.com

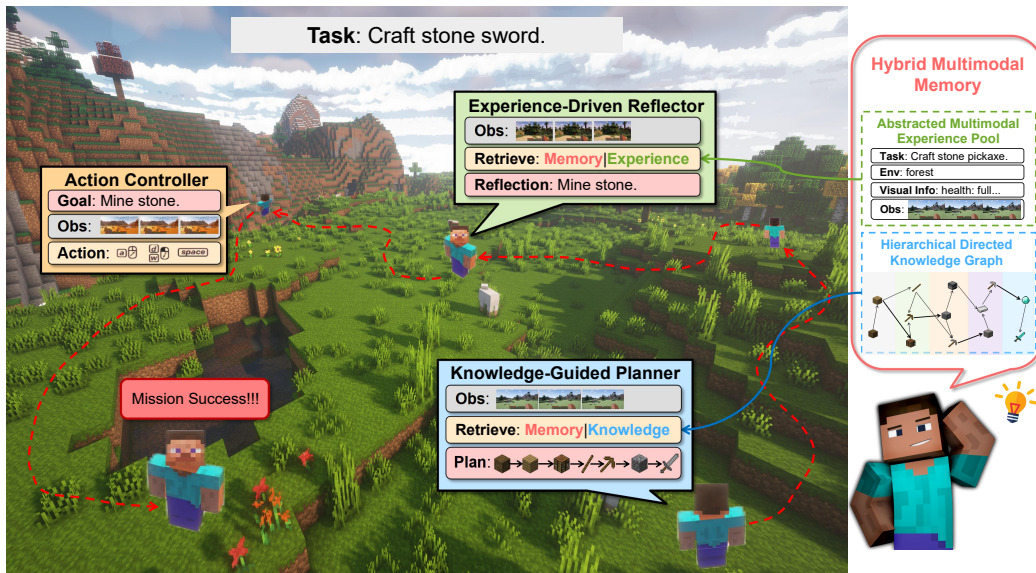


Figure 1: An illustration of Optimus-1 performing long-horizon tasks in Minecraft. Given the task “Craft stone sword”, Knowledge-Guided Planner incorporates knowledge from Hierarchical Directed Knowledge Graph into planning, then Action Controller executes these planning sequences step-by-step. During the execution of the task, the Experience-Driven Reflector is periodically activated and retrieve experience from Abstracted Multimodal Experience Pool to make reflection.

Abstract

Building a general-purpose agent is a long-standing vision in the field of artificial intelligence. Existing agents have made remarkable progress in many domains, yet they still struggle to complete long-horizon tasks in an open world. We attribute this to the lack of necessary world knowledge and multimodal experience that can guide agents through a variety of long-horizon tasks. In this paper, we propose a **Hybrid Multimodal Memory** module to address the above challenges. It **1)** transforms knowledge into **Hierarchical Directed Knowledge Graph** that allows agents to explicitly represent and learn world knowledge, and **2)** summarises historical information into **Abstracted Multimodal Experience Pool** that provide agents

*Corresponding authors

with rich references for in-context learning. On top of the Hybrid Multimodal Memory module, a multimodal agent, Optimus-1, is constructed with dedicated **Knowledge-guided Planner** and **Experience-Driven Reflector**, contributing to a better planning and reflection in the face of long-horizon tasks in Minecraft. Extensive experimental results show that Optimus-1 significantly outperforms all existing agents on challenging long-horizon task benchmarks, and exhibits near human-level performance on many tasks. In addition, we introduce various Multimodal Large Language Models (MLLMs) as the backbone of Optimus-1. Experimental results show that Optimus-1 exhibits strong generalization with the help of the Hybrid Multimodal Memory module, outperforming the GPT-4V baseline on various tasks. Please see the project page at <https://cybertronagent.github.io/Optimus-1.github.io/>.

1 Introduction

Optimus Prime faces complex tasks alongside humans in Transformers to protect the peace of the planet. Creating an agent [43, 13] like Optimus that can perceive, plan, reflect, and complete long-horizon tasks in an open world has been a longstanding aspiration in the field of artificial intelligence [22, 35, 36]. Early research developed simple policy through reinforcement learning [7] or imitation learning [1, 25]. A lot of work [46, 49] have utilized Large Language Models (LLMs) as action planners for agents, generating executable sub-goal sequences for low-level action controllers. Further, recent studies [51, 32] employed Multimodal Large Language Models (MLLMs) [4, 38, 55] as planner and reflector. Leveraging the powerful instruction-following and logical reasoning capabilities of (Multimodal) LLMs [24], LLM-based agents have achieved remarkable success across multiple domains [14, 9, 10, 54]. Nevertheless, the ability of these agents to complete long-horizon tasks still falls significantly short of human-level performance.

According to relevant studies [27, 41, 45], the human ability to complete long-horizon tasks in an open world relies on long-term memory storage, which is divided into knowledge and experience. The storage and utilization of knowledge and experience play a crucial role in guiding human behavior and enabling humans to adapt flexibly to their environments in order to accomplish long-horizon tasks. Inspired by this theory, we summarize the challenges faced by current agents as follows:

Insufficient Exploration of Structured Knowledge: Structured knowledge, encompassing open world rules, object relationships, and interaction methods with the environment, is essential for agents to complete complex tasks [33, 43]. However, MLLMs such as GPT-4V * lack sufficient knowledge in Minecraft. Existing agents [1, 25, 7] only learn dispersed knowledge from video data and are unable to efficiently represent and learn this structured knowledge, rendering them incapable of performing complex tasks.

Lack of Multimodal Experience: Humans derive successful strategies and lessons from information on historical experience [8, 31], which assists them in tackling current complex tasks. In a similar manner, agents can benefit from in-context learning with experience demonstrations [42, 53]. However, existing agents [46, 50, 32] only consider unimodal information, which prevents them from learning from multimodal experience as humans do.

To address the aforementioned challenges, we propose **Hybrid Multimodal Memory** module that consists of **Hierarchical Directed Knowledge Graph** (HDKG) and **Abstracted Multimodal Experience Pool** (AMEP). For HDKG, we map the logical relationships between objects into a directed graph structure, thereby transforming knowledge into high-level semantic representations. HDKG efficiently provides the agent with the necessary knowledge for task execution, without requiring any parameter updates. For AMEP, we dynamically summarize and store the multimodal information (*e.g.*, environment, agent state, task plan, video frames, etc.) from the agent’s task execution process, ensuring that historical information contains both a global overview and local details. Different from the method of directly storing successful cases as experience [51], AMEP considers both successful and failed cases as references. This innovative approach of incorporating failure cases into in-context learning significantly enhances the performance of the agent.

*<https://openai.com/index/gpt-4v-system-card/>

On top of the Hybrid Multimodal Memory module, we construct a multimodal composable agent, **Optimus-1**. As shown in Figure 1, Optimus-1 consists of Knowledge-Guided Planner, Experience-Driven Reflector, and Action Controller. To enhance the ability of agents to cope with complex environments and long-horizon tasks, Knowledge-Guided Planner incorporates visual observation into the planning phase, leveraging HDKG to capture the knowledge needed. This allows the agent to efficiently transform tasks into executable sub-goals. Action Controller takes the sub-goal and the current observation as inputs and generates low-level actions, interacting with the game environment to update the agent’s state. In open-world complex environments, agents are prone to be erroneous when performing long-horizon tasks. To address this, we propose Experience-Driven Reflector, which is periodically activated to retrieve relevant multimodal experiences from AMEP. This encourages the agent to reflect on its current actions and refine the plan.

We validate the performance of Optimus-1 in Minecraft, a popular open-world game environment. Experimental results show that Optimus-1 exhibits remarkable performance on long-horizon tasks, representing up to 30% improvement over existing agents. Moreover, we introduce various Multimodal Large Language Models (MLLMs) as the backbone of Optimus-1. Experimental results show that Optimus-1 has a 2 to 6 times performance improvement with the help of Hybrid Multimodal Memory, outperforming powerful GPT-4V baseline on lots of long-horizon tasks. Additionally, we verified that the plug-and-play Hybrid Multimodal Memory can drive Optimus-1 to incrementally improve its performance in a self-evolution manner. The extensive experimental results show that Optimus-1 makes a major step toward a general agent with a human-like level of performance. Main contributions of our paper:

- We propose **Hybrid Multimodal Memory** module which is composed of HDKG and AMEP. HDKG helps the agent make the planning of long-horizon tasks efficiently. AMEP provides refined historical experience and guides the agent to reason about the current situation state effectively.
- On top of the Hybrid Multimodal Memory module, we construct **Optimus-1**, which consists of **Knowledge-Guided Planner**, **Experience-Driven Reflector**, and Action Controller. Optimus-1 outperforms all baseline agents on long-horizon task benchmarks, and exhibits capabilities close to the level of human players.
- Driven by Hybrid Multimodal Memory, various MLLM-based Optimus-1 have demonstrated 2 to 6 times performance improvement, demonstrating the generalization of Hybrid Multimodal Memory.

2 Optimus-1

In this section, we first elaborate on how to implement the Hybrid Multimodal Memory in Sec 2.1. As a core innovation, it plays a crucial role in enabling Optimus-1 to execute long-horizon tasks. Next, we give an overview of Optimus-1 framework (Sec 2.2), which consists of Hybrid Multimodal Memory, Knowledge-Guided Planner, Experience-Driven Reflector, and Action Controller. Finally, we introduce a non-parametric learning approach to expand the hybrid multimodal memory (Sec 2.3), thereby enhancing the success rate of task execution for Optimus-1.

2.1 Hybrid Multimodal Memory

In order to endow agent with a long-term memory storage mechanism [27, 45], we propose the Hybrid Multimodal Memory module, which consists of Abstracted Multimodal Experience Pool (AMEP) and Hierarchical Directed Knowledge Graph (HDKG).

2.1.1 Abstracted Multimodal Experience Pool

Relevant studies [23, 28, 17, 15] highlight the importance of historical information for agents completing long-horizon tasks. Minedojo [7] and Voyager [46] employed unimodal storage of historical information. Jarvis-1 [51] used a multimodal experience mechanism that stores task planning and visual information without summarization, posing challenges to storage capacity and retrieval speed. To address this issue, we propose AMEP, which aims to dynamically summarize all multimodal information during task execution. It preserves the integrity of long-horizon data while enhancing storage and retrieval efficiency.

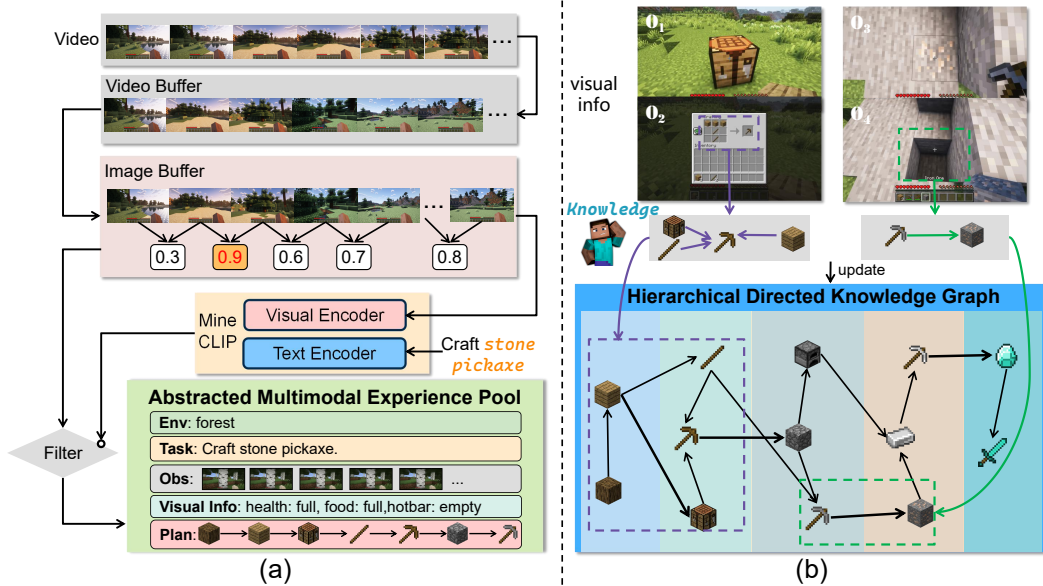
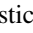
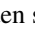
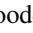
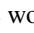


Figure 2: **(a)** Extraction process of multimodal experience. The frames are filtered through video buffer and image buffer, then MineCLIP [7] is employed to compute the visual and sub-goal similarities and finally they are stored in Abstracted Multimodal Experience Pool. **(b)** Overview of Hierarchical Directed Knowledge Graph. Knowledge is stored as a directed graph, where its nodes represent objects, and directed edges point to materials that can be crafted by this object.

Specifically, as depicted in Figure 2, to conduct the static visual information abstraction, the video stream captured by Optimus-1 during task execution is first input to a video buffer, filtering the stream at a fixed frequency of 1 frame per second. Based on the filtered video frames, to further perform a dynamic visual information abstraction, these frames are then fed into an image buffer with a window size of 16, where the image similarity is dynamically computed and final abstracted frames are adaptively updated. To align such abstracted visual information with the corresponding textual sub-goal, we then utilize MineCLIP [7], a pre-trained video-text alignment model, to calculate their multimodal correlation. When this correlation exceeds a threshold, the corresponding image buffer and textual sub-goal are saved as multimodal experience into a pool. Finally, we further incorporate environment information, agent initial state, and plan generated by Knowledge-Guided Planner, into such a pool, which forms the AMEP. In this way, we consider the multimodal information of each sub-goal, and summarise it to finally compose the multimodal experience of the given task.

2.1.2 Hierarchical Directed Knowledge Graph

In Minecraft, mining and crafting represent a complex knowledge network crucial for effective task planning. For instance, crafting a diamond sword  requires two diamonds  and one wooden stick , while mining diamonds requires an iron pickaxe , which involving further materials and steps. Such knowledge is essential for an agent’s ability to perform long-horizon complex tasks. Instead of implicit learning through fine-tuning [32, 58], we propose HDKG, which transforms knowledge into a graph representation. It enables the agent to perform explicit learning by retrieving information from the knowledge graph.

As shown in the Figure 2, we transform knowledge into a graph $\mathcal{D}(\mathcal{V}, \mathcal{E})$, where nodes set \mathcal{V} represent objects, and directed edges set \mathcal{E} point to nodes that can be crafted by this object. An edge $e \in \mathcal{E}$ in the \mathcal{D} can be represented as $e = (u, v)$, where $u, v \in \mathcal{V}$. The directed graph efficiently stores and updates knowledge. For a given object x , retrieving the corresponding node allows extraction of a sub-graph $\mathcal{D}_j(\mathcal{V}_j, \mathcal{E}_j) \in \mathcal{D}$, where nodes set \mathcal{V}_j and edges set \mathcal{E}_j can be formulated as:

$$\mathcal{V}_j = \{v \in \mathcal{V} \mid x\}, \quad \mathcal{E}_j = \{e = (u, v) \in \mathcal{E} \mid u \in \mathcal{V}_j \cup v \in \mathcal{V}_j\}, \quad (1)$$

Then by topological sorting, we can get all the materials and their relationships needed to complete the task. This knowledge is provided to the Knowledge-Guided Planner as a way to generate a more

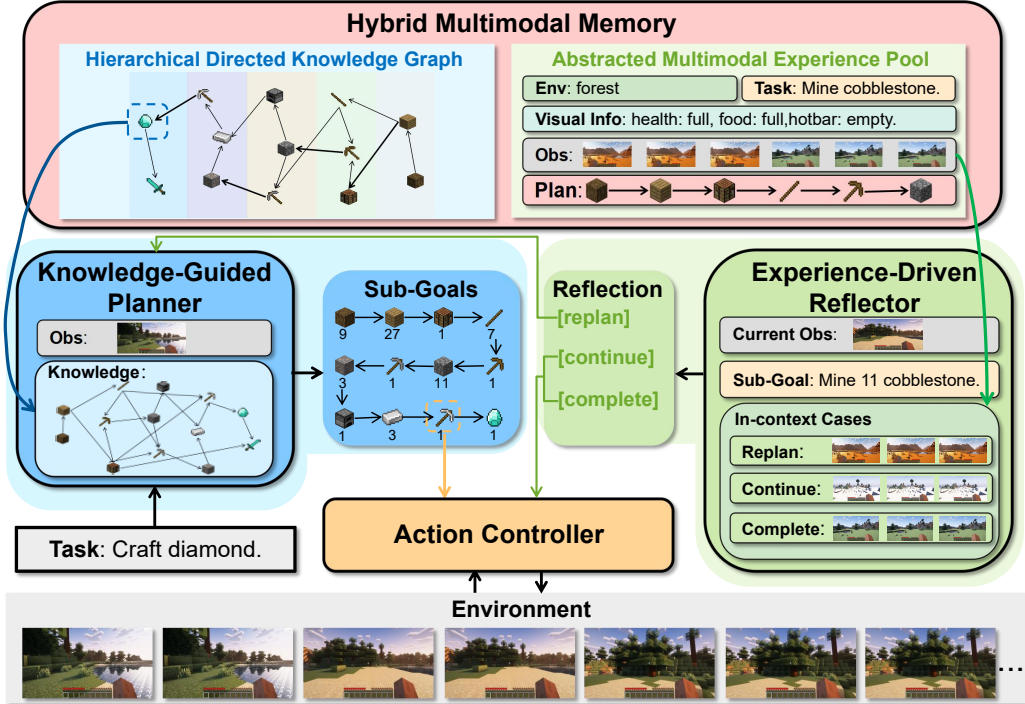


Figure 3: Overview framework of our Optimus-1. Optimus-1 consists of Knowledge-Guided Planner, Experience-Driven Reflector, Action Controller, and Hybrid Multimodal Memory architecture. Given the task “craft stone sword”, Optimus-1 incorporates the knowledge from HDKG into Knowledge-Guided Planning, then Action Controller generates low-level actions. Experience-Driven Reflector is periodically activated to introduce multimodal experience from AMEP to determine if the current task can be executed successfully. If not, it will ask the Knowledge-Guided Planner to refine the plan.

reasonable sequence of sub-goals. With HDKG, we can significantly enhance the world knowledge of the agent in a train-free manner.

2.2 Optimus-1: Framework

Relevant studies indicate that the human brain is essential for planning and reflection, while the cerebellum controls low-level actions, both crucial for complex tasks [39, 40]. Inspired by this, we divide the structure of Optimus-1 into Knowledge-Guided Planner, Experience-Driven Reflector, and Action Controller. In a given game environment with a long-horizon task, the Knowledge-Guided Planner senses the environment, retrieves knowledge from HDKG, and decomposes the task into executable sub-goals. The action controller then sequentially executes these sub-goals. During execution, the Experience-Driven Reflector is activated periodically, leveraging historical experience from AMEP to assess whether Optimus-1 can complete the current sub-goal. If not, it instructs the Knowledge-Guided Planner to revise its plan. Through iterative interaction with the environment, Optimus-1 ultimately completes the task.

Knowledge-Guided Planner. Open-world environments vary greatly, affecting task execution. Previous approaches [50] using LLMs for task planning failed to consider the environment, leading to the failure of tasks. For example, an agent in a cave aims to catch fish. It lacks visual information to plan conditions on the current situation, such as “leave the cave and find a river”. Therefore, we integrate environmental information into the planning stage. Unlike Jarvis-1 [51] and MP5 [32], which convert observation to textual descriptions, Optimus-1 directly employs observation as visual conditions to generate environment-related plans, *i.e.*, sub-goal sequences. This results in more comprehensive and reasonable planning. More importantly, Knowledge-Guided Planner retrieves the knowledge needed to complete the task from HDKG, allowing task planning to be done once, rather than generating the next step in each iteration. Given the task t , observation o , the sub-goals

sequence $g_1, g_2, g_3, \dots, g_n$ can be formulated as:

$$g_1, g_2, g_3, \dots, g_n = p_\theta(o, t, p_\eta(t)), \quad (2)$$

where n is the number of sub-goals, p_η denotes sub-graph retrieved from HDKG, p_θ denotes MLLM. In this paper, we employ OpenAI’s GPT-4V as Knowledge-Guided Planner and Experience-Driven Reflector. We also evaluate other alternatives of GPT-4V, such as open-source models like Deepseek-VL [26] and InternLM-XComposer2-VL [6] in Section 3.4.

Action Controller. It takes the sub-goal and the current observation as inputs and then generates low-level actions, which are control signals for the mouse and keyboard. Thus, it can interact with the game environment to update the agent’s state and the observation. The formulation is as follows:

$$a_k = p_\pi(o, g_i), \quad (3)$$

where a_k denotes low-level action at time k , p_π denotes action controller. Unlike generating code [46, 32, 49], generating control actions for the mouse and keyboard [1, 25, 51, 3] more closely resembles human behavior. In this paper, we employ STEVE-1 [25] as our Action Controller.

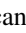
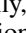

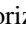
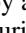
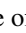
Experience-Driven Reflector. The sub-goals generated by Knowledge-Guided Planner are interdependent. The failure of any sub-goal halts the execution of subsequent ones, leading to overall task failure. Therefore, a reflection module is essential to identify and rectify errors promptly. During task execution, the Experience-Driven Reflector activates at regular intervals, retrieving historical experience from AMEP, and then analyzing the current state of Optimus-1. The reflection results of Optimus-1 are categorized as COMPLETE, CONTINUE, or REPLAN. COMPLETE indicates successful execution, prompting the action controller to proceed to the next sub-goal. CONTINUE signifies ongoing execution without additional feedback. REPLAN denotes failure, requiring the Knowledge-Guided Planner to revise the plan. The reflection r generated by Experience-Driven Reflector can be formulated as:

$$r = p_\theta(o, g_i, p_\epsilon(t)), \quad (4)$$

where p_ϵ denotes multimodal experience retrieved from AMEP. Experimental results in Section 3.3 demonstrate that the Experience-Driven Reflector significantly enhances the success rate of long-horizon tasks.








During task execution, even in cases where task failure necessitates REPLAN, multimodal experiences are stored in AMEP. Thus, during the reflection phase, Optimus-1 can retrieve the most relevant cases from each of the three scenarios COMPLETE, CONTINUE, and REPLAN from AMEP as references. Experimental Results in Section 3.3 demonstrate the effectiveness of this innovative method of incorporating failure cases into in-context learning.

2.3 Non-parametric Learning of Hybrid Multimodal Memory

To implement the Hybrid Multimodal Memory and enhance Optimus-1’s capacity, we propose a non-parametric learning method named “free exploration-teacher guidance”. In the free exploration phase, Optimus-1’s equipment and tasks are randomly initialized, and it explores random environments, acquiring world knowledge through environmental feedback. For example, it learns that “a stone sword  can be crafted with a wooden stick  and two cobblestones ”, storing this in the HDKG. Additionally, successful and failed cases are stored in the AMEP, providing reference experience for the reflection phase. We initialize multiple Optimus-1, and they share the same HDKG and AMEP. Thus the memory is filled up efficiently. After free exploration, Optimus-1 has basic world knowledge and multimodal experience. In the teacher guidance phase, Optimus-1 needs to learn a small number of long-horizon tasks based on extra knowledge. For example, it learns “a diamond sword  is obtained by a stick  and two diamonds  from the teacher, then perform the task “craft diamond sword”. During the teacher guidance phase, Optimus-1’s memory is further expanded and it gains the experience of executing complete long-horizon tasks.

Unlike fine-tuning, this method enhances Optimus-1 incrementally without updating parameters, in a self-evolution manner. Starting with an empty Hybrid Multimodal Memory, Optimus-1 iterates between “free exploration-teacher guidance” learning and unseen task inference. With each iteration, its memory capacity grows, enabling mastery of tasks from easy to hard.

Table 1: Main Result of Optimus-1 on long-horizon tasks benchmark. We report the average success rate (SR), average number of steps (AS), and average time (AT) on each task group, the results of each task can be found in the Appendix F. Lower AS and AT metrics mean that the agent is more efficient at completing the task, while $+\infty$ indicates that the agent is unable to complete the task. Overall represents the average result on the five groups of Iron, Gold, Diamond, Redstone, and Armor.

Group	Metric	GPT-3.5	GPT-4V	DEPS	Jarvis-1	Optimus-1	Human-level
 Wood	SR \uparrow	40.16	41.42	77.01	93.76	98.60	100.00
	AT \downarrow	56.39	55.15	85.53	67.76	47.09	31.08
	AS \downarrow	1127.78	1103.04	1710.61	1355.25	841.94	621.59
 Stone	SR \uparrow	20.40	20.89	48.52	89.20	92.35	100.00
	AT \downarrow	135.71	132.77	138.71	141.50	129.94	80.85
	AS \downarrow	2714.21	2655.47	2574.30	2830.05	2518.88	1617.00
 Iron	SR \uparrow	0.00	0.00	16.37	36.15	46.69	86.00
	AT \downarrow	$+\infty$	$+\infty$	944.61	722.78	651.33	434.38
	AS \downarrow	$+\infty$	$+\infty$	8892.24	8455.51	6017.85	5687.60
 Gold	SR \uparrow	0.00	0.00	0.00	7.20	8.51	17.31
	AT \downarrow	$+\infty$	$+\infty$	$+\infty$	787.37	726.35	557.08
	AS \downarrow	$+\infty$	$+\infty$	$+\infty$	15747.13	15527.07	13141.60
 Diamond	SR \uparrow	0.00	0.00	0.60	8.98	11.61	16.98
	AT \downarrow	$+\infty$	$+\infty$	1296.96	1255.06	1150.98	744.82
	AS \downarrow	$+\infty$	$+\infty$	23939.30	25101.25	23019.64	16237.54
 Redstone	SR \uparrow	0.00	0.00	0.00	16.31	25.02	33.27
	AT \downarrow	$+\infty$	$+\infty$	$+\infty$	1070.42	932.50	617.89
	AS \downarrow	$+\infty$	$+\infty$	$+\infty$	17408.40	12709.99	12357.00
 Armor	SR \uparrow	0.00	0.00	9.98	15.82	19.47	28.48
	AT \downarrow	$+\infty$	$+\infty$	997.59	924.60	824.53	551.30
	AS \downarrow	$+\infty$	$+\infty$	17951.95	16492.96	16350.56	11026.00
Overall	SR \uparrow	0.00	0.00	5.39	16.89	22.26	36.41

3 Experiments

3.1 Experiments Setting

Environment. To ensure realistic gameplay like human players, we employ MineRL [11] with Minecraft 1.16.5 as our simulation environment. The agent operates at a fixed speed of 20 frames per second and only interacts with the environment via low-level action control signals of the mouse and keyboard. For more information about the detailed descriptions of the observation and action spaces, please refer to the **Appendix B**.

Benchmark. We constructed a benchmark of 67 tasks to evaluate the Optimus-1’s ability to complete long-horizon tasks. As illustrated in Table 5, we divide the 67 Minecraft tasks into 7 groups according to recommended categories in Minecraft. Please refer to **Appendix D** for more details.

Baseline. We compare Optimus-1 with various agents, including GPT-3.5², GPT-4V, DEPS [50], and Jarvis-1 [51] on the challenging long-horizon tasks benchmark. In addition, we employed 10 volunteers to perform the same task on the benchmark, and their average performance served as a human-level baseline. Please refer to **Appendix D.2** for more details about human-level baseline. For a more comprehensive comparison, we also report Optimus-1’s performances on the benchmark used by Voyager [46], MP5 [32], and DEPS [50] in the **Appendix F.2**. Note that we initialize Optimus-1 with an empty inventory, while DEPS [50] and Jarvis-1 [51] have tools in their initial state. This makes it more challenging for Optimus-1 to perform the same tasks.

²<https://openai.com/research/gpt-3.5>

Table 2: Ablation study results. We report average success rate (SR) on each task group. P., R., K., E. represent Planning, Reflection, Knowledge, and Zero, Suc., and Fai. represent retrieving from Experience, respectively.

Ablation Setting				Task Group				
P.	R.	K.	E.	Wood	Stone	Iron	Gold	Diamond
				14.29	0.00	0.00	0.00	0.00
✓				42.95	25.67	0.00	0.00	0.00
✓	✓			55.00	47.37	18.11	2.08	1.11
✓	✓		✓	73.53	64.20	24.19	3.08	1.86
✓	✓	✓		92.37	69.63	38.33	3.49	2.42
✓	✓	✓	✓	97.49	94.26	53.33	11.54	9.59

Ablation Setting			Task Group				
Zero	Suc.	Fai.	Wood	Stone	Iron	Gold	Diamond
✓			92.00	79.26	36.32	4.25	3.25
	✓		95.00	84.29	46.98	9.36	7.89
		✓	95.00	81.10	45.47	7.50	6.39
	✓	✓	97.49	94.26	53.33	11.54	9.59

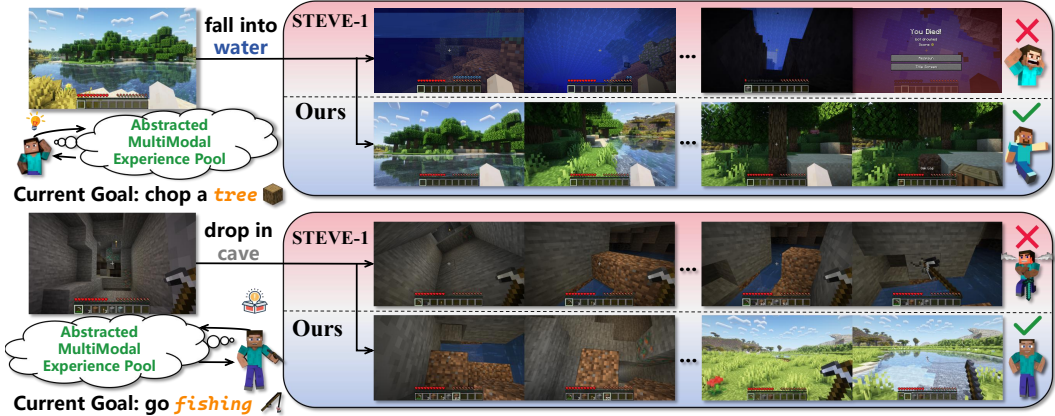


Figure 4: Illustration of the role of reflection mechanism. Without the help of reflective mechanisms, STEVE-1 [25] often gets into trouble and fails to complete the task. While Optimus-1, with the help of the Experience-Driven Reflector, leverages the AMEP to retrieve relevant experience, reflect current situation and correct errors. This improves Optimus-1’s success rate on long-horizon tasks.

Evaluation Metrics. The agent always starts in survival mode, with an empty inventory. We conducted at least 30 times for each task using different world seeds and reported the average success rate to ensure fair and thorough evaluation. Additionally, we add the average steps and average time of completing the task as evaluation metrics.

3.2 Experimental Results

The overall experimental results on benchmark are shown in Table 1, see the accuracy for each task in Appendix F. Optimus-1 has a success rate near 100% on the Wood Group 🪵. Compared with Jarvis-1, Optimus-1 has 29.28% and 53.40% improvement on the Diamond Group 💎 and Redstone Group 🔴, respectively. Optimus-1 achieves the best performance and the shortest elapsed time among all task groups. It reveals the effectiveness and efficiency of our proposed Optimus-1 framework. Moreover, compared with all baselines, Optimus-1 performance was closer (average 5.37% improvement) to human levels on long-horizon task groups.

3.3 Ablation Study

We conduct extensive ablation experiments on 18 tasks, experiment setting can be found in Table 6. As shown in Table 2, we first remove Knowledge-Driven Planner and Experience-Driven Reflector, the performance of Optimus-1 on all task groups drops dramatically. It demonstrates the necessity of Knowledge-Guided Planner and Experience-Driven Reflector modules for performing long-horizon tasks. As for Hybrid Multimodal Memory, we remove HDKG from Optimus-1. Without the help of world knowledge, the performance of Optimus-1 decreased by an average of 20% across all task groups. We then removed AMEP, this resulted in the performance of Optimus-1 decreased by an

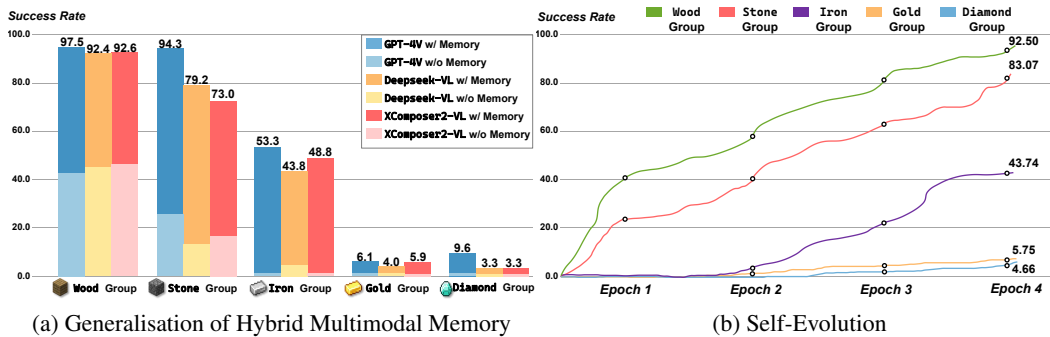


Figure 5: **(a)** With the help of Hybrid Multimodal Memory, various MLLM-based Optimus-1 have demonstrated 2 to 6 times performance improvement. **(b)** Illustration of the change in Optimus-1 success rate on the unseen task over 4 epochs.

average of 12%. Finally, we performed ablation experiments on the way of retrieving cases from AMEP. As shown in Table 3, without retrieving cases from AMEP, the success rate shows an average of 10% decrease across all groups. It reveals that this reflection mechanism, which considers both success and failure cases, has a significant impact on the performance of Optimus-1. To illustrate the role of the reflection mechanism, we have shown some cases in Figure 4.

3.4 Generalization Ability

In this section, we explore an interesting issue: whether generic MLLMs can effectively perform various long-horizon complex tasks in Minecraft using Hybrid Multimodal Memory. As shown in Figure 5, We employ Deepseek-VL [26] and InternLM-XComposer2-VL [6] as Knowledge-Guided Planner and Experience-Driven Reflector. The experimental results show that the original MLLM has low performance on long-horizon tasks due to the lack of knowledge and experience of Minecraft. With the assistance of Hybrid Multimodal Memory, the performance of MLLMs has improved by 2 to 6 times across various task groups, outperforming the GPT-4V baseline. This encouraging result demonstrates the generalization of the proposed Hybrid Multimodal Memory.

3.5 Self-Evolution via Hybrid Multimodal Memory

As shown in Section 2.3, we randomly initialize the Hybrid Multimodal Memory of Optimus-1, then update it multiple times by using the “free exploration-teacher guidance” learning method. We set the epoch to 4, and the number of learning tasks to 160. At each period, Optimus-1 performs free exploration on 150 tasks and teacher guidance learning on the remaining 10 tasks, we then evaluate Optimus-1’s learning ability on the task groups same as ablation study. Experimental results are shown in Figure 5. It reveals that Optimus-1 keeps getting stronger through the continuous expansion of memory during the learning process of multiple periods. Moreover, it demonstrates that MLLM with Hybrid Multimodal Memory can incarnate an expert agent in a self-evolution manner [44].

4 Related Work

4.1 Agents in Minecraft

We summarise the differences of existing Minecraft agents in the **Appendix D.3**. Earlier work [29, 56, 2, 3] introduced policy models for agents to perform simple tasks in Minecraft. MineCLIP [7] used text-video data to train a contrastive video-language model as a reward model for policy, while VPT [1] pre-trained on unlabelled videos but lacked instruction as input. Building on VPT and MineCLIP, STEVE-1 [25] added text input to generate low-level action sequences from human instructions and images. However, these agents struggle with complex tasks due to limitations in instruction comprehension and planning. Recent work [49, 46, 59] incorporated LLMs as planning and reflection modules, but lacked visual information integration for adaptive planning. MP5

[32], MineDremer [58], and Jarvis-1 [51] enhanced situation-aware planning by obtaining textual descriptions of visual information, yet lacked detailed visual data. Optimus-1 addresses these issues by directly using observation as situation-aware conditions in the planning phase, enabling more rational, visually informed planning. Additionally, unlike other agents requiring multiple queries for task refinement, Optimus-1 generates a complete and effective plan in one step with the help of HDKG. This makes Optimus-1 planning more efficient.

4.2 Memory in Agents

In the agent-environment interaction process, memory is key to achieving experience accumulation [21], environment exploration [16], and knowledge abstraction [57]. There are two forms to represent memory content in LLM-based agents: textual form [17, 15, 30] and parametric form [5, 28, 47, 20]. In textual form, the information is explicitly retained and recalled by natural languages. In parametric form, the memory information [37] is encoded into parameters and implicitly influences the agent’s actions. Recent work [48, 52, 12] has explored the long-term visual information storage [18, 19] and summarisation in MLLM. Our proposed hybrid multimodal memory module is plug-and-play and can provide world knowledge and multimodal experience for Optimus-1 efficiently.

5 Conclusion

In this paper, we propose Hybrid Multimodal Memory module, which consists of two parts: HDKG and AMEP. HDKG provides the necessary world knowledge for the planning phase of the agent, and AMEP provides the refined historical experience for the reflection phase of the agent. On top of the Hybrid Multimodal Memory, we construct the multimodal composable agent, Optimus-1, in Minecraft. Extensive experimental results show that Optimus-1 outperforms all existing agents on long-horizon tasks. Moreover, we validate that general-purpose MLLMs, based on Hybrid Multimodal Memory and without additional parameter updates, can exceed the powerful GPT-4V baseline. The extensive experimental results show that Optimus-1 makes a major step toward a general agent with a human-like level of performance.

6 Limitation and Future Work

In the framework of Optimus-1, we are dedicated to leverage proposed Hierarchical Directed Knowledge Graph and Abstracted Multimodal Experience Pool can be used to enhance the agent’s ability to plan and reflect. For Action Controller, we directly introduce STEVE-1 [25] as a generator of low-level actions. However, limited by STEVE-1’s ability to follow instructions and execute complex actions, Optimus-1 is weak in completing challenging tasks such as “beat ender dragon” and “build a house”. Therefore, a potential future research direction is to enhance the instruction following and action generation capabilities of action controller.

In addition, most of the work, including Optimus-1, utilize a multimodal large language model for planning and reflection, which then drives an action controller to perform the task. Building an end-to-end vision-language-action agent will be future work.

7 Acknowledgement

This study is supported by National Natural Science Foundation of China (Grant No. 62236003 and 62306090), Shenzhen College Stability Support Plan (Grant No. GXWD20220817144428005), Natural Science Foundation of Guangdong Province of China (Grant No. 2024A1515010147), and Major Key Project of Peng Cheng Laboratory (Grant No. PCL2023A08).

References

- [1] Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654, 2022.
- [2] Shaofei Cai, Zihao Wang, Xiaojian Ma, Anji Liu, and Yitao Liang. Open-world multi-task control through goal-aware representation learning and adaptive horizon prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13734–13744, 2023.
- [3] Shaofei Cai, Bowei Zhang, Zihao Wang, Xiaojian Ma, Anji Liu, and Yitao Liang. Groot: Learning to follow instructions by watching gameplay videos. In *The Twelfth International Conference on Learning Representations*, 2023.
- [4] Gongwei Chen, Leyang Shen, Rui Shao, Xiang Deng, and Liqiang Nie. Lion: Empowering multimodal large language model with dual-level visual knowledge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [5] Nicola De Cao, Wilker Aziz, and Ivan Titov. Editing factual knowledge in language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6491–6506, 2021.
- [6] Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang Zhang, Haodong Duan, Maosong Cao, et al. Internlm-xcomposer2: Mastering free-form text-image composition and comprehension in vision-language large model. *arXiv preprint arXiv:2401.16420*, 2024.
- [7] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35:18343–18362, 2022.
- [8] Mariel K Goddu and Alison Gopnik. The development of human causal learning and reasoning. *Nature Reviews Psychology*, pages 1–21, 2024.
- [9] Maitrey Gramopadhye and Daniel Szafir. Generating executable action plans with environmentally-aware language models. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3568–3575. IEEE, 2023.
- [10] Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. A real-world webagent with planning, long context understanding, and program synthesis. *arXiv preprint arXiv:2307.12856*, 2023.
- [11] William H Guss, Brandon Houghton, Nicholay Topin, Phillip Wang, Cayden Codel, Manuela Veloso, and Ruslan Salakhutdinov. Minerl: A large-scale dataset of minecraft demonstrations. *arXiv preprint arXiv:1907.13440*, 2019.
- [12] Bo He, Hengduo Li, Young Kyun Jang, Menglin Jia, Xuefei Cao, Ashish Shah, Abhinav Shrivastava, and Ser-Nam Lim. Ma-lmm: Memory-augmented large multimodal model for long-term video understanding. *arXiv preprint arXiv:2404.05726*, 2024.
- [13] Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. An embodied generalist agent in 3d world. *arXiv preprint arXiv:2311.12871*, 2023.
- [14] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*, pages 9118–9147. PMLR, 2022.
- [15] Ziheng Huang, Sebastian Gutierrez, Hemanth Kamana, and Stephen MacNeil. Memory sandbox: Transparent and interactive memory management for conversational agents. In *Adjunct Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–3, 2023.

- [16] Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen, Li Fei-Fei, Anima Anandkumar, Yuke Zhu, and Linxi Fan. Vima: General robot manipulation with multimodal prompts. In *NeurIPS 2022 Foundation Models for Decision Making Workshop*, 2022.
- [17] Dacheng Li, Rulin Shao, Anze Xie, Ying Sheng, Lianmin Zheng, Joseph Gonzalez, Ion Stoica, Xuezhe Ma, and Hao Zhang. How long can context length of open-source llms truly promise? In *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*, 2023.
- [18] Xiaojie Li, Shaowei He, Jianlong Wu, Yue Yu, Liqiang Nie, and Min Zhang. Mask again: Masked knowledge distillation for masked video modeling. In *Proceedings of the ACM International Conference on Multimedia*, page 2221–2232. ACM, 2023.
- [19] Xiaojie Li, Jianlong Wu, Shaowei He, Shuo Kang, Yue Yu, Liqiang Nie, and Min Zhang. Fine-grained key-value memory enhanced predictor for video representation learning. In *Proceedings of the ACM International Conference on Multimedia*, page 2264–2274. ACM, 2023.
- [20] Xiaojie Li, Yibo Yang, Xiangtai Li, Jianlong Wu, Yue Yu, Bernard Ghanem, and Min Zhang. Genview: Enhancing view quality with pretrained generative model for self-supervised learning. In *Proceedings of the European Conference on Computer Vision*. Springer, 2024.
- [21] Xiaojie Li, Yibo Yang, Jianlong Wu, Bernard Ghanem, Liqiang Nie, and Min Zhang. Mamba-fscil: Dynamic adaptation with selective state space model for few-shot class-incremental learning. *arXiv preprint arXiv:2407.06136*, 2024.
- [22] Zaijing Li, Fengxiao Tang, Ming Zhao, and Yusen Zhu. Emocaps: Emotion capsule based model for conversational emotion recognition. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1610–1618, 2022.
- [23] Zaijing Li, Ting-En Lin, Yuchuan Wu, Meng Liu, Fengxiao Tang, Ming Zhao, and Yongbin Li. Unisa: Unified generative framework for sentiment analysis. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 6132–6142, 2023.
- [24] Zaijing Li, Gongwei Chen, Rui Shao, Dongmei Jiang, and Liqiang Nie. Enhancing the emotional generation capability of large language models via emotional chain-of-thought. *arXiv preprint arXiv:2401.06836*, 2024.
- [25] Shalev Lifshitz, Keiran Paster, Harris Chan, Jimmy Ba, and Sheila McIlraith. Steve-1: A generative model for text-to-behavior in minecraft. *Advances in Neural Information Processing Systems*, 2023.
- [26] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, et al. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024.
- [27] Simon Makin. The amyloid hypothesis on trial. *Nature*, 559(7715):S4–S4, 2018.
- [28] Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. Fast model editing at scale. In *International Conference on Learning Representations*, 2021.
- [29] Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. Zero-shot task generalization with multi-task deep reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning*, pages 2661–2670. PMLR, 2017.
- [30] Arka Pal, Deep Karkhanis, Manley Roberts, Samuel Dooley, Arvind Sundararajan, and Sid-dartha Naidu. Giraffe: Adventures in expanding context lengths in llms. *arXiv preprint arXiv:2308.10882*, 2023.
- [31] Eileen Parkes. Scientific progress is built on failure. *Nature*, 10, 2019.
- [32] Yiran Qin, Enshen Zhou, Qichang Liu, Zhenfei Yin, Lu Sheng, Ruimao Zhang, Yu Qiao, and Jing Shao. Mp5: A multi-modal open-ended embodied system in minecraft via active perception. *arXiv preprint arXiv:2312.07472*, 2023.

- [33] Maria Abi Raad, Arun Ahuja, Catarina Barros, Frederic Besse, Andrew Bolt, Adrian Bolton, Bethanie Brownfield, Gavin Buttimore, Max Cant, Sarah Chakera, et al. Scaling instructable agents across many simulated worlds. *arXiv preprint arXiv:2404.10179*, 2024.
- [34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [35] Rui Shao, Xiangyuan Lan, Jiawei Li, and Pong C Yuen. Multi-adversarial discriminative deep domain generalization for face presentation attack detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10023–10031, 2019.
- [36] Rui Shao, Tianxing Wu, and Ziwei Liu. Detecting and grounding multi-modal media manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6904–6913, 2023.
- [37] Rui Shao, Tianxing Wu, Jianlong Wu, Liqiang Nie, and Ziwei Liu. Detecting and grounding multi-modal media manipulation and beyond. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [38] Leyang Shen, Gongwei Chen, Rui Shao, Weili Guan, and Liqiang Nie. Mome: Mixture of multimodal experts for generalist multimodal large language models. In *NeurIPS*, 2024.
- [39] Shan H Siddiqi, Konrad P Kording, Josef Parvizi, and Michael D Fox. Causal mapping of human brain function. *Nature reviews neuroscience*, pages 361–375, 2022.
- [40] JF Stein. Role of the cerebellum in the visual guidance of movement. *Nature*, pages 217–221, 1986.
- [41] Tim Stuart, Andrew Butler, Paul Hoffman, Christoph Hafemeister, Efthymia Papalexi, William M Mauck, Yuhan Hao, Marlon Stoeckius, Peter Smibert, and Rahul Satija. Comprehensive integration of single-cell data. *cell*, 177(7):1888–1902, 2019.
- [42] Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Zhengxiong Luo, Yuezhe Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, et al. Generative multimodal models are in-context learners. *arXiv preprint arXiv:2312.13286*, 2023.
- [43] Weihao Tan, Ziluo Ding, Wentao Zhang, Boyu Li, Bohan Zhou, Junpeng Yue, Haochong Xia, Jiechuan Jiang, Longtao Zheng, Xinrun Xu, et al. Towards general computer control: A multimodal agent for red dead redemption ii as a case study. *arXiv preprint arXiv:2403.03186*, 2024.
- [44] Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu Li, Yuchuan Wu, Yongbin Li, Zhi Jin, Fei Huang, Dacheng Tao, and Jingren Zhou. A survey on self-evolution of large language models. *arXiv preprint arXiv:2404.14387*, 2024.
- [45] Deniz Vatanserver, Jonathan Smallwood, and Elizabeth Jefferies. Varying demands for cognitive control reveals shared neural processes supporting semantic and episodic memory retrieval. *Nature communications*, 12(1):2134, 2021.
- [46] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.
- [47] Mengru Wang, Ningyu Zhang, Ziwen Xu, Zekun Xi, Shumin Deng, Yunzhi Yao, Qishen Zhang, Linyi Yang, Jindong Wang, and Huajun Chen. Detoxifying large language models via knowledge editing. *arXiv preprint arXiv:2403.14472*, 2024.
- [48] Shihao Wang, Zhiding Yu, Xiaohui Jiang, Shiyi Lan, Min Shi, Nadine Chang, Jan Kautz, Ying Li, and Jose M Alvarez. Omnidrive: A holistic llm-agent framework for autonomous driving with 3d perception, reasoning and planning. *arXiv preprint arXiv:2405.01533*, 2024.

- [49] Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. *arXiv preprint arXiv:2302.01560*, 2023.
- [50] Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. *arXiv preprint arXiv:2302.01560*, 2023.
- [51] Zihao Wang, Shaofei Cai, Anji Liu, Yonggang Jin, Jinbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng He, Zilong Zheng, Yaodong Yang, et al. Jarvis-1: Open-world multi-task agents with memory-augmented multimodal language models. *arXiv preprint arXiv:2311.05997*, 2023.
- [52] Yuetian Weng, Mingfei Han, Haoyu He, Xiaojun Chang, and Bohan Zhuang. Longvlm: Efficient long video understanding via large language models. *arXiv preprint arXiv:2404.03384*, 2024.
- [53] Xu Yang, Yongliang Wu, Mingzhuo Yang, Haokun Chen, and Xin Geng. Exploring diverse in-context configurations for image captioning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [54] Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent: Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*, 2023.
- [55] Qilang Ye, Zitong Yu, Rui Shao, Xinyu Xie, Philip Torr, and Xiaochun Cao. Cat: Enhancing multimodal large language model to answer questions in dynamic audio-visual scenarios. In *ECCV*, 2024.
- [56] Haoqi Yuan, Chi Zhang, Hongcheng Wang, Feiyang Xie, Penglin Cai, Hao Dong, and Zongqing Lu. Skill reinforcement learning and planning for open-world long-horizon tasks. *arXiv preprint arXiv:2303.16563*, 2023.
- [57] Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. A survey on the memory mechanism of large language model based agents. *arXiv preprint arXiv:2404.13501*, 2024.
- [58] Enshen Zhou, Yiran Qin, Zhenfei Yin, Yuzhou Huang, Ruimao Zhang, Lu Sheng, Yu Qiao, and Jing Shao. Minedreamer: Learning to follow instructions via chain-of-imagination for simulated-world control. *arXiv preprint arXiv:2403.12037*, 2024.
- [59] Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, et al. Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory. *arXiv preprint arXiv:2305.17144*, 2023.

A Broader Impact

With the increasing capability level of Multimodal Large Language Models (MLLM) comes many potential benefits and also risks. On the positive side, we anticipate that the techniques that used to create Optimus-1 could be applied to the creation of helpful agents in robotics, video games, and the web. This plug-and-play architecture that we have created can be quickly adapted to different MLLMs, and the proposed methods also provide a viable solution for other application areas in the agent domain. However, on the negative side, it is imperative to acknowledge the inherent stochastic nature of MLLMs in text generation. If not addressed carefully, this could lead to devastating consequences for society. Prior to deploying MLLMs in conjunction with the Hybrid Multimodal Memory methodology, a comprehensive assessment of their potential risks must be undertaken. We hope that while the stakes are low, works such as ours can improve access to safety research on instruction-following models in multimodal agents domains.

B Minecraft

Minecraft is an extremely popular sandbox video game developed by Mojang Studios³. It allows players to explore a blocky, procedurally generated 3D world with infinite terrain, discover and extract raw materials, craft tools and items, and build structures or earthworks (shown in Figure 6). Minecraft is a valuable and representative environment for evaluating long-horizon tasks, offering greater diversity and complexity compared to other environments. Unlike web/app navigation [54] and embodied manipulation [16], Minecraft is an open world with a complex and dynamic environment (79 biomes, including ocean, plains, forest, desert, etc.). To complete long-horizon tasks, agents must achieve multiple sub-goals (e.g., 15 sub-goals to craft a diamond sword), making the construction of a Minecraft agent quite challenging. Many studies [46, 32, 51] have chosen Minecraft as the environment for validating performance on long-horizon tasks. Extensive experimental results in the paper show that Optimus-1 outperforms all baselines. Therefore, we chose Minecraft as open-world environment to evaluate the ability of agents to perform long-horizon tasks.

B.1 Basic Rules

Biomes. The Minecraft world is divided into different areas called “biomes”. Different biomes contain different blocks and plants and change how the land is shaped. There are 79 biomes in Minecraft 1.16.5, including ocean, plains, forest, desert, etc. Diverse environments have high requirements for the generalization of agents.

Time. Time passes within this world, and a game day lasts for 20 real-world minutes. Nighttime is much more dangerous than daytime: the game starts at dawn, and agents have 10 minutes of game time before nightfall. Hostile or neutral mobs spawn when night falls, and most of these mobs are dangerous, trying to attack agents. How to survive in such a dangerous world is an open problem for Minecraft agents research.

Item. In Minecraft 1.16.5, there are 975 items can be obtained, such as wooden pickaxe 🪓, iron sword ⚔️. Item can be obtained by crafting or destroying blocks or attacking entities. For example, agent can attack cows 🐮 to obtain leather 🍷 and beef 🍖. Agent also can use 1 stick 🪵 and 2 diamonds 💎 to craft diamond sword ⚔️.

Gameplay progress. Progression primarily involves discovering and utilizing various materials and resources, each of which unlocks new capabilities and options. For instance, crafting a wooden pickaxe 🪓 enables the player to mine stone 🪨, which can be used to create a stone pickaxe 🪓 and a furnace 🏠; these, in turn, allow for the mining and smelting of iron ore 🪨. Subsequently, an iron pickaxe 🪓 permits the extraction of diamonds 💎, and a diamond pickaxe ⚔️ can mine virtually any block in the game. Similarly, cultivating different crops allows for the breeding of various animals, each providing distinct resources beyond mere sustenance. Enemy drops also have specific applications, with some being more beneficial than others. By integrating resources from mining, farming, and breeding, players can enchant their equipment. The collection and crafting of materials also facilitate construction, enabling players to build diverse structures. Beyond practical

³<https://www.minecraft.net/en-us/article/meet-mojang-studios>

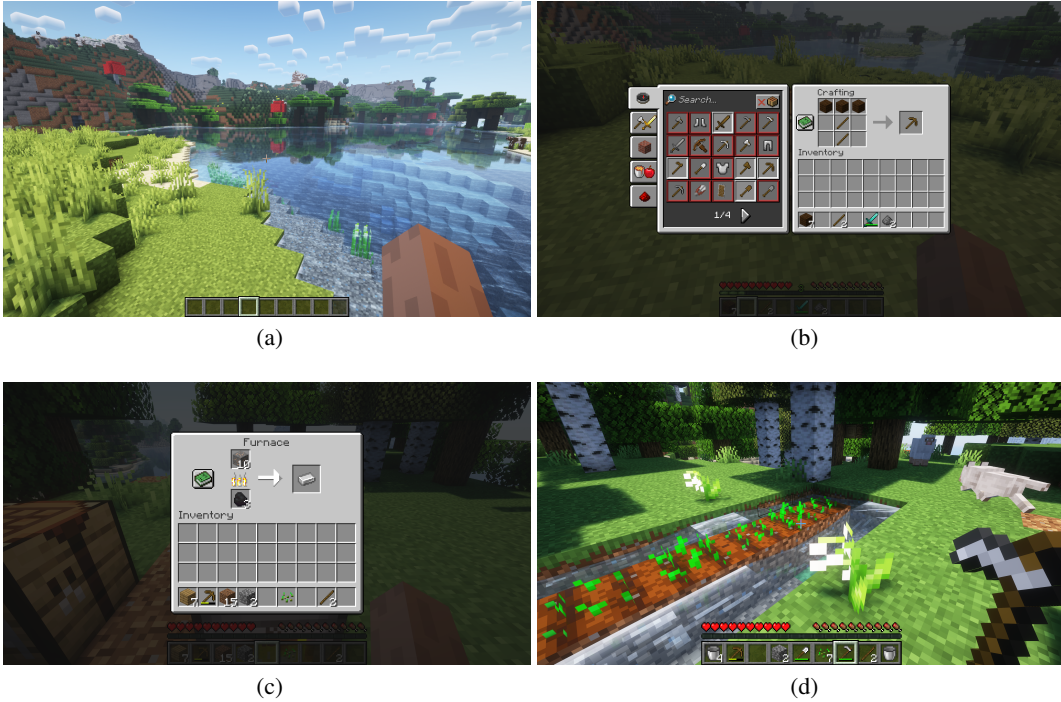




Figure 6: The screenshots in Minecraft. **(a)**. The world has different complex terrains, including plains, river, forest and mine. **(b)**. The agent can use crafting table to craft tools and items with recipes. **(c)**. The agent can use the furnace to smelt ore to obtain precious ingot. **(d)**. The agent can grow wheat near the river.

considerations such as secure bases and farms, the creative aspect of building personalized structures constitutes a significant part of the Minecraft experience.

Freedom. In Minecraft, player can do anything they can imagine. Player can craft tools, smelt ore, brew potions, trade with villagers and wandering traders, attack mobs, grow crops, raise animals in captivity, etc. Player even can use redstone  to build a computer. This is a world of freedom and infinite possibilities.

More Challenge than Diamond . Progression beyond the Overworld is fairly limited: Eventually, you can build a nether portal to reach the Nether, where you can get materials for more complex crafting, the resources to brew potions, and the top tier of tools and armor. The Nether materials also let you reach the End dimension, where you must defeat the Ender Dragon to unlock the outer End Islands, where you can get an elytra that lets you fly, and shulker boxes for more storage.

B.2 Observation and Action Spaces

Observation. Our observation space is completely consistent with human players. The agent only receives an RGB image with dimensions of 640×360 during the gameplay process, including the hotbar, health indicators, food saturation, and animations of the player’s hands. It is worth helping the agent see more clearly in extremely dark environments, we have added a night vision effect for the agent, which increases the brightness of the environment during the night.

Action Spaces. Our action space is almost similar to human players, except for craft and smelt actions. It consists of two parts: the mouse and the keyboard. The keypresses are responsible for controlling the movement of agents, such as jumping, forward, back, etc. The mouse movements are responsible for controlling the perspective of agents and the cursor movements when the GUI is opened. The left and right buttons of the mouse are responsible for attacking and using or placing items. In Minecraft, precise mouse movements are important when completing complex tasks that need open inventory or crafting table. In order to achieve both the same action space with MineDojo

Table 4: Our action space.

Index	Action	Human Action	Description
1	Forward	key W	Move forward.
2	Back	key S	Move back.
3	Left	key A	Strafe left.
4	Right	key D	Strafe right.
5	Jump	key Space	Jump. When swimming, keeps the player afloat.
6	Sneak	key left Shift	Slowly move in the current direction of movement.
7	Sprint	key left Ctrl	Move quickly in the direction of current movement.
8	Attack	left Button	Destroy blocks (hold down); Attack entity (click once).
9	Use	right Button	Place blocks, entity, open items or other interact actions defined by game.
10	hotbar [1-9]	keys 1-9	
11	Open/Close Inventory	key E	Opens the Inventory. Close any open GUI.
12	Yaw	move Mouse X	Turning; aiming; camera movement.Ranging from -180 to +180.
13	Pitch	move Mouse Y	Turning; aiming; camera movement.Ranging from -180 to +180.
14	Craft	-	Execute a crafting recipe to obtain new item
15	Smelt	-	Execute a smelting recipe to obtain new item.

[7], we abstract the craft and the smelt action into action space. The detailed action space is described in Table 4.

B.3 Long-horizon Tasks

Long-horizon Tasks are complex tasks that require world knowledge to solve and consist of multiple indispensable subtask sequences. In Minecraft, technology has six levels, including wood 🪵, stone 🪨, iron ⚒️, golden 🛠️, diamond 💎, and netherite ⚔️. Wooden tools can mine stone-level blocks, but can’t mine iron-level and upper-level blocks. Stone tools can mine iron-level blocks, but can’t mine diamond-level and upper-level blocks. Iron-level tools can mine diamond-level blocks, but can’t mine netherite-level blocks. Diamond-level tools can mine any level blocks.

For example, the agent now wants to complete the task “Craft iron sword ⚔️”. The agent needs to craft wood-level tools to mine stone 🪵, and craft stone-level tools to mine iron ore 🪨. In order to craft tools, the agent needs a crafting table 🪵. To smelt iron ore 🪨 into iron ingot ⚒️, the agent needs a furnace 🪪. Moreover, craft crafting table needs 4 planks, and craft furnace needs 8 cobblestone. In summary, the agent needs to obtain many raw materials, wood-level and stone-level tools, 1 crafting table, 1 furnace, and most importantly, 2 iron ingots. The process of this task is shown in Figure 7.

C Theory

In this section, we briefly introduce the relevant theory of cognitive science. For more details, please refer to the original articles.

Our ability to understand and predict the world around us depends on our long-term memory stores, which have historically been divided into two distinct systems [27, 41, 45]. The semantic memory system provides a conceptual framework for describing the similar meanings of words and objects as they are encountered in different contexts (e.g., a bee is a flying insect with yellow and black stripes that produces honey), whereas the episodic memory system records our personal experiences characterized by the co-occurrence of words and objects at different times and places (e.g., being stung by a bee while eating honey at a picnic last weekend). These information stores and the interactions between them play a crucial role in guiding our behaviour and giving us the flexibility to adapt to the various demands of our environment.

In this paper, inspired by the above theory, we divide the agent memory module into two parts: knowledge and experience. Based on this, we propose Hierarchical Directed Knowledge Graph and Abstracted Multimodal Experience Pool to enable the agent to acquire, store, and utilize knowledge and experience during the execution of tasks. Extensive experimental results demonstrate the effectiveness of the proposed methodology

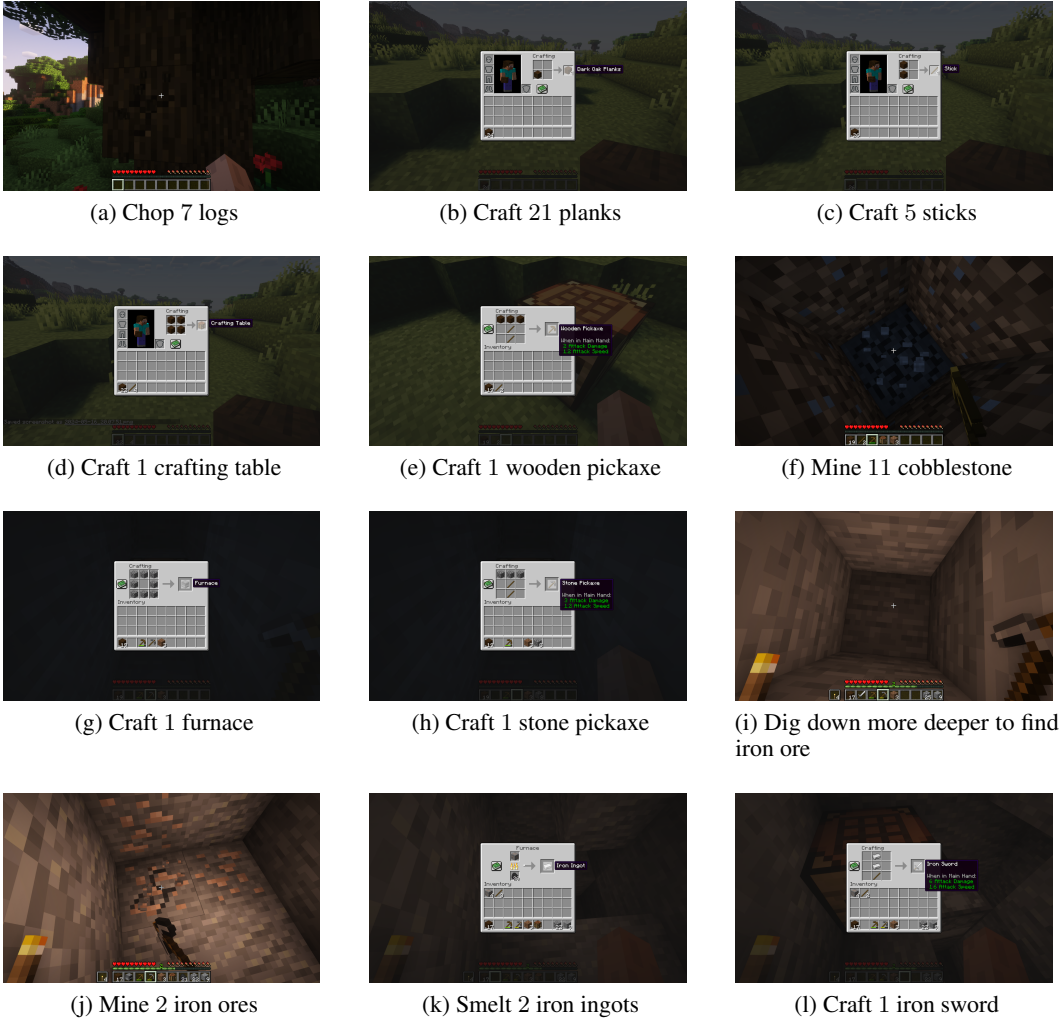


Figure 7: Processing of task "Craft 1 iron sword". Optimus-1 needs thousands of steps to complete this task. To craft and smelt precisely, the mouse movements action can't have any error.

D Benchmark Suite

D.1 Benchmark

We constructed a benchmark of 67 tasks to evaluate Optimus-1's ability to complete long-horizon tasks in Minecraft. According to recommended categories in Minecraft, we have classified these tasks into 7 groups: Wood 🌳, Stone 🪨, Iron ⚒️, Gold 🏆, Diamond 💎, Redstone 🔴, Armor 🛡️. The statistics for benchmark are shown in Table 5. Due to the varying complexity of these tasks, we adopt different maximum gameplay steps (Max. Steps) for each task. The maximum steps are determined by the average steps that human players need to complete the task. Due to the randomness of Minecraft, the world and initial spawn point of the agent could vary a lot. In our benchmark setting, We initialize the agent with an empty inventory, which makes it necessary for the agent to complete a series of sub-goals (mining materials, crafting tools) in order to perform any tasks. This makes every task challenging, even for human players.

Note that Diamonds are a very rare item that only spawns in levels 2 to 16 and have a 0.0846% chance of spawning in Minecraft 1.16.5. Diamonds are usually found near level 9, or in man-made or natural mines no higher than level 16. In order to reduce the huge impact that diamond generation

Table 5: Setting of 7 groups encompassing 67 Minecraft long-horizon tasks.

Group	Task Num.	Example Task	Max. Steps	Initial Inventory	Avg. Sub-goal Num.
Wooden	10	Craft a wooden axe	3600	Empty	4.60
Stone	9	Craft one stone pickaxe	7200	Empty	8.78
Iron	16	Craft a iron pickaxe	12000	Empty	12.75
Golden	6	Mine gold and smelt into golden ingot	36000	Empty	15.83
Redstone	6	Craft a piston	36000	Empty	17.50
Diamond	7	Dig down and mine a diamond	36000	Empty	15.14
Armor	13	Craft one iron helmet	36000	Empty	15.85

Table 6: We evaluate Optimus-1 on these tasks in ablation study which are the subset of our benchmark.

Group	Task	Sub-Goal Num.	Max. Step	Initial Inventory
Wooden	Craft a wooden axe	5	3600	Empty
	Craft a crafting table	3	3600	Empty
Stone	Craft a stone pickaxe	10	7200	Empty
	Craft a stone axe	10	7200	Empty
	Craft a furnace	9	7200	Empty
Iron	Craft a iron pickaxe	13	12000	Empty
	Craft a bucket	13	12000	Empty
	Craft a rail	13	12000	Empty
	Craft a iron sword	12	12000	Empty
	Craft a shears	12	12000	Empty
Golden	Craft a golden pickaxe	16	36000	Empty
	Craft a golden axe	16	36000	Empty
	Smelt a golden ingot	15	36000	Empty
Diamond	Craft a diamond pickaxe	15	36000	Empty
	Craft a diamond axe	16	36000	Empty
	Craft a diamond hoe	15	36000	Empty
	Craft a diamond sword	15	36000	Empty
	Dig down and mine a diamond	15	36000	Empty

probability has on agent’s likelihood of completing a task, we have adjusted the diamond generation probability to 20%, spawns in levels 2 to 16. This setting applies to human players as well.

In the ablation study, we select the subset of our benchmark as the test set (shown in Table 6). The environment setting is the same as the benchmark.

D.2 Baselines

Existing Baseline. On the one hand, we employ GPT-3.5 and GPT-4V as baseline, which are evaluated without integrating hybrid multimodal memory modules. During the planning phase, they generate a plan for the action controller based on task prompt (and observation). During the reflection phase, they generate reflection results in a zero-shot manner. On the other hand, we compare existing SOTA Agents [50, 51] in Minecraft.

Human-level Baseline. To better demonstrate agent’s performance level in Minecraft, we hired 10 volunteers to play the game as a human-level baseline. The volunteers played the game with the same environment and settings, and every volunteer asked to perform the each task on the benchmark 10 times. Ultimately, we used the average scores of 10 volunteers as the human-level baseline. The results of the human-level baseline are shown in Table 1. To ensure the validity of the experiment, we ensured that each volunteer had at least 20 hours of Minecraft gameplay before conducting the experiment. For each volunteer, we pay \$25 as reward.

Table 7: Statistics for various Minecraft agents.

Agent	Pub.	Env.	Input	Output	Planning	Reflection	Knowledge	Experience
VPT [1]	NeurIPS' 22	MineRL	V	low-level action				
MineDOJO [7]	NeurIPS' 22	MineDOJO	T+V	low-level action				
STEVE-1 [25]	NeurIPS' 23	MineRL	T+V	low-level action				
Voyager [46]	NeurIPS' 23	Mineflayer	T+V	code	✓	✓		✓
DEPS [50]	NeurIPS' 23	MineDOJO	T+V	code	✓	✓		
GROOT [3]	ICLR' 24	MineRL	T+V	low-level action				
MP5 [32]	CVPR' 24	MineDOJO	T+V	code	✓	✓		
Jarvis-1 [51]	-	MineRL	T+V	low-level action	✓	✓		✓
Optimus-1	-	MineRL	T+V	low-level action	✓	✓	✓	✓

D.3 Minecraft Agents

In this section, we summarise the differences between existing Minecraft agents. As shown in the Table 7, earlier work [1, 7, 25, 3] constructed Transformer-based policy network as agent. Recent work [46, 50, 51, 32] introduces the Multimodal Large Language Model, which empowers the agent to complete long-horizon tasks by exploiting the powerful language comprehension and planning capabilities of LLM.

In the Mineflayer and Minedojo environments, agents [7, 46, 50, 32] can accomplish sub-goals by calling APIs (in the form of codes), which is a different behavioral pattern from humans. In MineRL [11], agents [1, 25, 3, 51] must generate low-level actions to perform tasks, which is more challenging to accomplish long-horizon tasks.

Moreover, existing agents lack knowledge and experience, and their performance in Minecraft is still vastly gapped from the human level. In this paper, we introduce Hybrid Multimodal Memory, which empowers Optimus-1 with hierarchical knowledge and multimodal experience. This makes Optimus-1 significantly outperform all existing agents on challenging long-horizon tasks benchmark, and exhibits near human-level performance on many tasks.

E Implementation Details

E.1 Hybrid Multimodal Memory

E.1.1 Abstracted Multimodal Experience Pool

Relevant studies [5, 28, 17, 15] have demonstrated the importance of memory for agents to complete long-horizon tasks. To implement the memory mechanism, Minedojo [7] and Voyager [46] only considered unimodal storage of historical information. Jarvis-1 [51] considered a multimodal memory mechanism to store task planning and visual information as experience, but it stores all historical information without summarisation. This approach stores all visual images, which poses a huge challenge in storage size and retrieval efficiency. To solve the problem, we propose the Abstracted Multimodal Experience Pool structure, which summarizes all historical information during the agent’s execution of the task, which maintains the integrity of long sequential information and greatly improves the storage and retrieval efficiency of the experience.

As shown in Figure 2, we first input the visual image stream to the video buffer, which filters the image stream at a fixed frequency. It makes the length of the image stream substantially shorter. Empirically, we set the frequency of filtering to 1 second/frame, meaning that the video buffer takes one frame per second from the original image stream to compose the filtered image stream. We found that above this frequency makes the visual information redundant (too much similarity between images), and below this frequency does not preserve enough complete visual information.





Then, we feed the filtered frames into an image buffer with a window size of 16. We dynamically compute the similarity between images in the image buffer, when a new image comes in, we compute the similarity between the new image and the most recent image, and then we remove the image with the highest similarity in order to keep the image buffer’s window size to 16.

Subsequently, we introduce MineCLIP [7], a pre-trained model of video-text alignment with a structure similar to CLIP [34], as our visual summariser. For a given sub-goal, it calculates the

correlation between the visual content within the current memory bank and the sub-goal, and when this correlation exceeds a pre-set threshold, the frames within the memory bank are saved as the visual memories corresponding to that sub-goal. Finally, we store the visual memories with the sub goal’s textual description into the Abstracted Multimodal Experience Pool. In addition, we incorporate the environment information, agent initial state, plan from Knowledge-Guided Planner, etc. into the experience memory of the given task. In this way, we consider the history information of each sub-goal and summaries and summarise it to finally compose the multimodal experience of the given task.

Note that we also store these visual memories as failure cases when the feedback from the reflection phase is REPLAN. Therefore, when Optimus-1 executes a long-horizon task, it can retrieve past successes and failures as references and update memory after the task is finished. In the reflection phase, Optimus-1 retrieve the most relevant cases from Abstracted Multimodal Experience Pool, which contains the three scenarios COMPLETE, CONTINUE, and REPLAN, to help the agent better assess which state the current situation belongs to. This approach of considering both successful and failed cases for in-context learning is inspired by related research [8, 31], and its effectiveness is validated in Section 3.3.

E.1.2 Hierarchical Directed Knowledge Graph

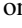


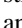

As shown in the Figure 2, crafting a diamond sword  requires two diamonds  and a wooden stick , while mining diamonds requires an iron pickaxe , which in turn requires additional raw materials and crafting steps. We transform this mine and craft knowledge into a graph structure, where the nodes of the graph are objects, and the nodes point to objects that can be crafted or completed by that object. With directed graph, we show that connections between objects are established, and that this knowledge can be stored and updated efficiently. For a given object, we only need to retrieve the corresponding node to extract the corresponding subgraph from the knowledge graph. Then by topological sorting, we can get the antecedents and required materials for the object, and this information is provided to the Knowledge-Guided Planner as a way to generate a more reasonable sequence of sub-goals. With Hierarchical Directed Knowledge Graph, we can significantly enhance the world knowledge of the agent in a train-free manner, as shown in the experimental results in Section 3.3.

Our HDKG can be efficiently updated and expanded. When adding new nodes, the HDKG can be updated by simply merging the nodes and relationships into the graph. This method involves local linear modifications to the graph rather than altering the entire graph, making the process efficient and time-saving. For example, when M new nodes and N edges are added, the HDKG can be updated with M+N times of operations. Moreover, an HDKG containing 851 objects (nodes) requires less than 1 MB of memory. Thus, the HDKG can be efficiently updated and maintained.

E.2 Hybrid Multimodal Memory Driven Optimus-1

In order to implement the proposed Hybrid Multimodal Memory and to progressively increase the capacity of Optimus-1 in a self-evolution manner, we propose a non-parametric learning method named “free exploration-teacher guidance”.

In the free exploration phase, we randomly initialize the environment, materials, and tasks. For the task “craft a wooden pickaxe”, we provide initial materials (three planks, two sticks), and then Optimus-1 (only the action controller activated) attempts to complete the task. If the environment feedback indicates the task is successful, the knowledge {3 planks, 2 sticks → wooden pickaxe} is added to the HDKG. Note that we randomly initialize materials and their quantities, which means that the task may not always succeed. As a result, each free exploration may not acquire the corresponding knowledge, but it can record the relevant experience (whether successful or fail). In the free exploration phase, Optimus-1 learns simple atomic operations, such as crafting sticks in the Wooden Group and mining diamonds in the Diamond Group.

In the teacher guidance phase, Optimus-1 need to learn a small number of long-horizon tasks based on extra knowledge. For example, during the free exploration phase, Optimus-1 mastered crafting stick  and mining diamond , but did not know that “a diamond sword  is obtained by a stick  and two diamonds ”. So we provide some task plans, which will serve as extra knowledge to guide

Optimus-1 to complete the task of “craft diamond sword”. We built the following automated process to get the task plan needed for “free exploration”:

- We randomly select 5 tasks for each Group (7 groups in total) that are not included in the benchmark.
- For each selected task, we use a script to automatically obtain the crafting relationships from the Minecraft Wiki ⁴. Taking the task “craft a wooden sword” as an example, we use the script to automatically obtain the crafting relationships: 1 wooden stick, 2 planks, 1 crafting table → 1 wooden sword, 1 log → 4 planks, 2 planks → 4 sticks, 4 planks → 1 crafting table.
- These relationships are converted into a directed acyclic graph through an automated script. By performing a topological sort, the graph can be converted into tuples of materials and their quantities: (wooden sword, 1), (crafting table, 1), (wooden stick, 1) (planks, 8), (log, 2).
- We prompt GPT-4 to construct a plan in order from basic materials to advanced materials.
- Finally, we get the plan: 1. Get two logs 2. Craft eight planks 3. Craft a crafting table 4. Craft a wooden stick 5. Craft a wooden sword

During the teacher guidance phase, Optimus-1’s memory is further expanded and it gains the experience of executing complete long-horizon tasks. Teacher guidance phase allows Optimus-1 to acquire advanced knowledge and learn multimodal experiences through complete long-horizon tasks.

E.3 Backbone of Optimus-1

Optimus-1 consists of Knowledge-Guided Planner, Experience-Driven Reflector, and Action Controller. In this paper, we employ OpenAI’s GPT-4V (gpt-4-turbo) ⁵ as Knowledge-Guided Planner and Experience-Driven Reflector, and STEVE-1 [25] as Action Controller. We also employ open-source models like Deepseek-VL [26] and InternLM-XComposer2-VL [6] as Knowledge-Guided Planner and Experience-Driven Reflector.

All experiments were implemented on 4x NVIDIA A100 GPUs. We employ multiple Optimus-1 to perform different tasks at the same time, and this parallelized inference greatly improves our experimental efficiency. In the free exploration and teacher guidance phases, there is no need to access OpenAI’s API, and the learning process takes approximately 16 hours on 4x A100 80G GPUs. During the inference phase, it takes about 20 hours on 4x A100 80G GPUs.

Throughout the experiment, we spent about \$5,000 to access the GPT-4V API. However, we also offer more cost-effective solutions. As shown in Figure 5, if we employ Deepseek-VL [26] or InternLM-XComposer2-VL [6] as Optimus-1’s backbone, we can get comparable performance with low-cost!

E.4 Prompt for Optimus-1

We show the prompt templates for Experience-Driven Reflector and Action Controller as follows.

```
System: You are a MineCraft game expert and you can guide agents to complete complex tasks.
User: For a given game screen and task, you need to complete <goal inference> and <visual
inference>.
<goal inference>: According to the task, you need to infer the weapons, equipment, or
materials required to complete the task.
<visual inference>: According to the game screen, you need to infer the following aspects:
health bar, food bar, hotbar, environment.
I will give you an example as follow:
[Example]
<task>: craft a stone sword.
<goal inference>: stone sword
<visual inference>
health bar: full
food bar: full
```

⁴<https://minecraft.wiki/>

⁵<https://openai.com/index/gpt-4v-system-card/>

```

hotbar: empty
environment: forest
Here is a game screen and task, you MUST output in example format.
<task>: {task}.
<game screen>: {image}

Assistant:

=====
User: Now you need to make a plan with the help of <visual info> and <craft graph>.
<visual info>: Consists of the following aspects: health bar, food bar, hotbar, environment.
Based on the current visual information, you need to consider whether prequel steps
needed to ensure that agent can complete the task.
<craft graph>: a top-down list of all the tools and materials needed to complete the task.
I will give you an example of planning under specific visual conditions as follow:
[Example]
{example}
Here is a game screen and task, you MUST output in example format. Remember <task planning>
MUST output in example format.
<task>: {task}
<game screen>: {image}
<craft graph>: {graph}
Assistant:

```

Listing 1: Prompt for Knowledge-Guided Planner.

```

System: You are a Minecraft game expert and you can guide agents to complete complex tasks.
Agent is executing the task: {task}.
Given two images about agent's state before executing the task and its current state, you
should first detection the environment (forest, cave, ocean, etc.,) in which the agent
is located, then determine whether the agent's current situation is done, continue, or
replan.
<done>: Comparing the image before the task was performed, the current image reveals that the
task is complete.
<continue>: Current image reveals that the task is NOT complete, but agent is in good state (
good health, not hungry) with high likelihood to complete task.
<replan>: Current image reveals that the task is NOT complete, and agent is in bad state (bad
health, or hungry) or situation (in danger, or in trouble), need for replanning. For
replan, you need to further determine whether the agent's predicament is "drop_down" or
"in_water". "drop_down" means that the agent has fallen into a cave or is trapped in a
mountain or river, while "in_water" means that the agent is in the ocean and needs to
return to land immediately.

User: I'll give you some examples to illustrate the different situations. Each example
consists of two images, where the first image is the state of the agent before
performing the task and the second image is the current state of the agent.

[Examples]
<done>: {image1},{image2}
<continue>: {image1},{image2}
<replan>: {image1},{image2}

Now given two images about agent's state before executing the task and its current state, you
MUST and ONLY output in following format:
Environment: <environment>
Situation: <situation>
(if situation is replan) Predicament: <predicament>

```

Listing 2: Prompt for Experience-Driven Reflector.

F Additional Experimental Results

F.1 Full Results on Our Benchmark

We list the results of each task on the benchmark below, with details including task name, sub-goal numbers, success rate (SR), average number of steps (AS), average time (AT), and eval times. All tasks are evaluated in Minecraft 1.16.5 Survival Mode. Note that each time Optimus-1 performs a task, we initial it with an empty initial inventory and a random start point. This makes it challenging for Optimus-1 to perform each task.

Table 8: The results of Optimus-1 on various tasks in the Wood group. SR, AS, AT denote success rate, average number of steps, and average time (seconds), respectively.

Task	Sub-Goal Num.	SR	AS	AT(s)	Eval Times
Craft a wooden shovel	6	95.00	995.58	49.78	40
Craft a wooden pickaxe	5	100.00	1153.91	57.70	30
Craft a wooden axe	5	96.67	1010.28	50.51	30
Craft a wooden hoe	5	100.00	1042.80	52.14	30
Craft a stick	4	97.14	372.97	18.65	70
Craft a crafting table	3	98.55	448.63	22.43	69
Craft a wooden sword	5	100.00	1214.90	60.74	30
Craft a chest	4	100.00	573.80	28.69	30
Craft a bowl	4	100.00	744.30	37.21	30

Table 9: The results of Optimus-1 on various tasks in the Stone group. SR, AS, AT denote success rate, average number of steps, and average time (seconds), respectively.

Task	Sub-Goal Num.	SR	AS	AT(s)	Eval Times
Craft a stone shovel	8	90.32	2221.00	111.05	31
Craft a stone pickaxe	10	96.77	2310.09	115.50	31
Craft a stone axe	10	96.88	2112.59	105.63	32
Craft a stone hoe	8	94.64	2684.60	134.23	56
Craft a charcoal	9	88.57	3083.35	154.17	35
Craft a smoker	9	90.24	3118.89	155.94	41
Craft a stone sword	8	94.29	2067.92	103.40	35
Craft a furnace	9	93.75	2842.71	142.14	32
Craft a torch	8	85.71	2109.00	105.45	95

Moreover, in MineRL [11] environment, ‘steps’ refers to the number of interactions between the agent and the environment, occurring at a frequency of 20 times per second. For example, if an agent takes 2 seconds to complete the task “chop a tree”, it interacts with the environment 40 times, resulting in a recorded steps number of 40. Experimental results show that Optimus-1’s average task completion step (AS) is significantly lower than other baselines.

F.2 Results on Other Benchmark

For a more comprehensive comparison with current Minecraft Agents, we also report Optimus-1’s performances on the benchmark used by Voyager [46], MP5 [32], and DEPS [50] below. Due to the different environments and settings, agents perform tasks with varying degrees of difficulty. For example, Optimus-1 requires low-level action to perform any task in MineRL [11], and we initialize its inventory to be empty. While Voyager [46] performs tasks in Mineflayer⁶ environment only through encapsulated code, MP5 [32] performs tasks in MineDOJO [7] environment only needs a specific control signal to craft tools, no low-level actions (mouse movement and click) are needed.

Optimus-1’s success rate in completing tasks with these baselines is shown in the Table 15 and Table 16, and Optimus-1’s efficiency in unlocking the tech tree in Minecraft is shown in the Figure 8. These results reveal that Optimus-1 outperforms a variety of powerful baseline agents, even in challenging environmental settings!

⁶<https://github.com/PrismarineJS/mineflayer>

Table 10: The results of Optimus-1 on various tasks in the Iron group. SR, AS, AT denote success rate, average number of steps, and average time (seconds), respectively.

Task	Sub Goal Num.	SR	AS	AT(s)	Eval Times
Craft an iron shovel	13	54.79	5677.35	637.81	73
Craft an iron pickaxe	13	59.42	6157.39	591.81	69
Craft an iron axe	13	54.29	6026.26	676.97	70
Craft a stone_hoe	13	52.70	6650.97	743.82	74
Craft a bucket	13	54.29	6124.61	591.35	70
Craft a hopper	14	46.67	7242.14	710.17	60
Craft a rail	13	42.19	6713.07	754.48	64
Craft an iron sword	12	57.14	5625.49	633.91	70
Craft a shears	12	53.62	5058.00	570.35	69
Craft a smithing table	12	44.93	5317.39	594.81	69
Craft a tripwire hook	13	48.57	4968.74	562.66	70
Craft a chain	13	44.93	5764.42	645.33	69
Craft an iron bars	12	42.00	6508.43	723.13	50
Craft an iron nugget	12	30.99	4697.23	525.29	71
Craft a blast furnace	14	25.71	7760.67	711.05	35
Craft a stonecutter	13	34.78	5993.38	675.52	46

Table 11: The results of Optimus-1 on various tasks in the Gold group. SR, AS, AT denote success rate, average number of steps, and average time (seconds), respectively.

Task	Sub Goal Num.	SR	AS	AT(s)	Eval Times
Craft a golden shovel	16	9.80	13734.75	686.74	51
Craft a golden pickaxe	16	13.75	9672.00	783.60	80
Craft a golden axe	16	4.44	10158.75	707.94	45
Craft a golden hoe	16	3.33	13120.50	756.03	27
Craft a golden sword	16	3.33	9792.00	789.60	26
Smelt and craft a gold_ingot	15	16.42	9630.27	681.51	67

Table 12: The results of Optimus-1 on various tasks in the Diamond group. SR, AS, AT denote success rate, average number of steps, and average time (seconds), respectively.

Task	Sub Goal Num.	SR	AS	AT(s)	Eval Times
Craft a diamond shovel	15	18.75	23696.75	1184.84	64
Craft a diamond pickaxe	15	15.71	32189.50	1609.46	70
Craft a diamond axe	16	4.00	21920.50	1096.03	75
Craft a diamond hoe	15	4.61	24031.00	1201.55	65
Craft a diamond sword	15	14.52	27555.50	1377.78	62
Dig down and mine a diamond	15	9.09	20782.13	1039.11	64
Craft a jukebox	15	14.58	25056.00	1252.80	48

Table 13: The results of Optimus-1 on various tasks in the Redstone group. SR, AS, AT denote success rate, average number of steps, and average time (seconds), respectively.

Language Instruction	Sub-Goal Num.	SR	AS	AT(s)	Eval Times
Craft a piston	16	28.57	6457.10	822.85	35
Craft a redstone torch	16	29.63	6787.87	939.39	27
Craft an activator rail	18	15.68	8685.62	934.28	51
Craft a compass	23	15.00	14908.67	845.43	40
Craft a dropper	16	37.50	7272.80	1063.64	40
Craft a note block	16	24.32	6727.89	936.39	37

Table 14: The results of Optimus-1 on various tasks in the Armor group. SR, AS, AT denote success rate, average number of steps, and average time (seconds), respectively.

Task	Sub Goal Num.	SR	AS	AT(s)	Eval Times
Craft shield	14	43.33	7229.00	861.45	30
Craft iron chestplate	14	47.22	7230.24	851.51	36
Craft iron boots	14	23.81	6597.33	729.87	42
Craft iron leggings	14	6.67	9279.00	763.95	30
Craft iron helmet	14	58.14	6287.11	814.36	43
Craft diamond helmet	17	2.08	7342.00	867.10	48
Craft diamond chestplate	17	2.70	7552.00	777.60	37
Craft diamond leggings	17	9.68	7664.67	883.23	31
Craft diamond boots	17	16.67	10065.60	803.28	30
Craft golden helmet	17	12.50	11563.25	778.16	32
Craft golden leggings	17	14.60	10107.33	805.37	41
Craft golden boots	17	6.06	10311.00	915.55	33
Craft golden chestplate	17	9.67	10407.58	820.38	31

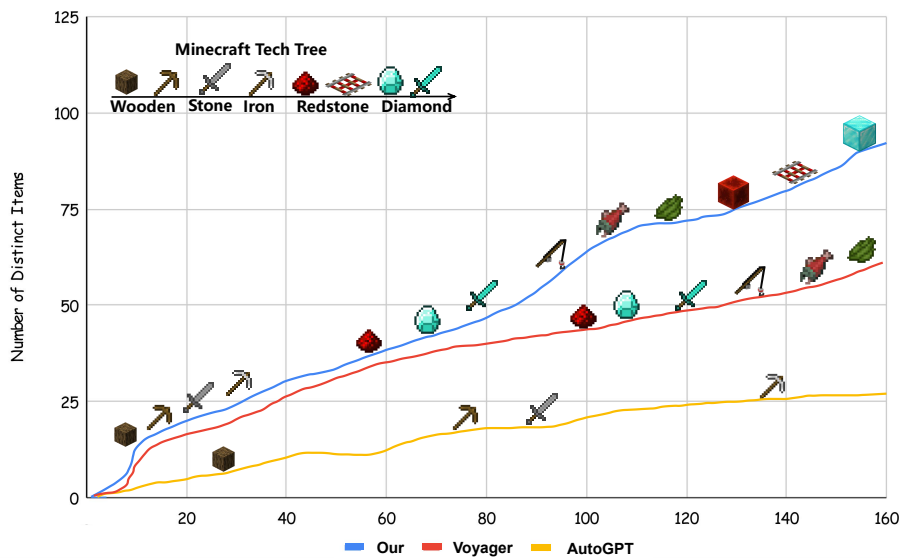


Figure 8: An illustration of Optimus-1 unlocking the tech tree in Minecraft.

Table 15: Result on Process-Dependent Tasks compared with MP5 [32]. SR, AS, AT denote success rate, average number of steps, and average time (seconds), respectively.

Task Level		MP5 [32]	Optimus-1		
		SR	SR	AS	AT(s)
Basic Level	log	96.67	100.00	586.58	29.33
	sand	96.67	94.32	1540.33	77.02
	planks	96.67	100.00	571.06	28.55
	stick	96.67	97.14	372.97	18.65
	crafting table	93.33	98.55	448.63	22.43
	Average	96.00	98.00	703.91	35.20
Wooden Level	bowl	93.33	100.00	744.30	37.21
	boat	93.33	92.86	1170.00	58.50
	chest	90.00	100.00	573.80	28.69
	wooden sword	86.67	100.00	1214.90	60.74
	wooden pickaxe	80.00	100.00	1153.91	57.70
	Average	88.67	98.57	971.38	48.56
Stone Level	cobblestone	80.00	95.29	1492.00	74.60
	furnace	80.00	93.75	2842.71	142.14
	stone pickaxe	80.00	96.77	2310.09	115.50
	iron ore	60.00	50.00	3017.00	150.85
	glass	80.00	81.11	3870.75	193.54
	Average	76.00	83.38	2706.51	135.32
Iron Level	iron ingot	56.67	59.42	4697.23	634.86
	shield *	56.67	43.33	7229.00	661.45
	bucket	53.33	54.29	6124.61	606.23
	iron pickaxe	50.00	59.42	6157.39	607.87
	iron door	43.33	48.28	5528.00	676.40
	Average	52.00	52.94	5947.25	637.36
Diamond Level	diamond ore *	30.00	9.09	20782.13	1039.10
	mind redstone	20.00	25.12	6787.87	739.39
	compass	16.67	15.00	14908.67	745.43
	diamond pickaxe	23.33	15.71	32189.50	1609.48
	piston	20.00	28.57	6457.10	622.85
	Average	22.00	18.70	14963.83	948.19

Table 16: Results (success rate) on 8 META TASK groups compared with DEPS [49].

Meta-Task	Task Object	InnerMonologue	Code-as-Policy	DEPS [49]	Ours
Basic MT1	planks	83.3	83.3	83.3	100.0
	stick	83.3	83.3	86.7	97.1
	chest	0.0	50.0	76.7	100.0
	sign	0.0	43.3	86.7	94.3
	boat	26.7	56.7	73.3	92.9
	trapdoor	56.7	56.7	76.7	96.2
	bowl	23.3	46.7	80.0	100.0
Tool(Simple) MT2	crafting_table	70.0	70.0	90.0	98.5
	wooden_pickaxe	80.0	80.0	80.0	100.0
	wooden_sword	83.3	83.3	86.7	100.0
	wooden_shovel	76.7	76.7	90.0	95.0
	furnace	40.0	40.0	66.7	93.7
	stone_pickaxe	36.7	53.3	73.3	96.7
	stone_axe	30.0	30.0	70.0	96.8
	stone_hoe	36.7	56.7	66.7	94.6
	stone_shovel	36.7	36.7	66.7	90.3
stone_sword	53.3	36.7	80.0	94.2	
Hunt and Food MT3	bed	6.7	6.7	43.3	90.0
	painting	16.7	16.7	86.7	92.2
	carpet	0.0	13.3	43.3	91.3
	cooked_porkchop	0.0	0.0	50.0	90.0
	cooked_beef	0.0	0.0	63.3	90.0
	cooked_mutton	0.0	0.0	66.7	90.0
Dig-down MT4	stone_stairs	36.7	16.7	66.7	90.3
	stone_slab	16.7	33.3	73.3	91.2
	lever	46.7	46.7	83.3	91.0
	coal	6.7	0.0	20.0	86.5
	torch	6.7	0.0	13.3	85.7
Equipment MT5	leather_boots	13.3	13.3	60.0	68.2
	leather_chestplate	0.0	6.7	36.7	64.2
	leather_helmet	6.7	0.0	70.0	65.9
	leather_leggings	20.0	0.0	56.7	65.5
	iron_chestplate	0.0	0.0	0.0	47.2
	iron_leggings	0.0	0.0	3.3	6.6
	iron_helmet	0.0	0.0	3.3	58.1
	iron_boots	0.0	0.0	20.0	23.8
	shield	0.0	6.7	13.3	43.3
Tool Complex MT6	bucket	0.0	3.3	6.7	54.3
	shears	0.0	0.0	30.0	53.6
	iron_pickaxe	6.7	0.0	10.0	59.4
	iron_axe	0.0	0.0	16.7	54.3
	iron_hoe	0.0	0.0	13.3	52.7
	iron_shovel	0.0	0.0	13.3	57.8
	iron_sword	0.0	3.3	6.7	54.7
Iron-Stage MT7	iron_bars	0.0	0.0	6.7	42.0
	hopper	0.0	0.0	6.7	46.7
	iron_door	0.0	0.0	3.3	48.3
	tripwire_hook	6.7	0.0	30.0	48.6
	rail	0.0	0.0	6.7	42.2
Challenge MT8	diamond	0.0	0.0	0.6	9.1

G Case Study

This section introduces several cases to comprehensively demonstrate Optimus-1’s capabilities.

Figures 9, 10, and 11 demonstrate the superiority of our reflection mechanism, which dynamically adjusts the plan based on the current game progress.

- Figure 9 illustrates Optimus-1’s replanning ability. When Optimus-1 realizes it cannot complete a task (such as a craft failure shown in the figure), it will replan the current task and continue execution.
- Figures 10 and 11 showcase Optimus-1’s ability to make judgments based on visual signals. When Optimus-1 determines that it has completed a task (such as “kill a cow 🐮” in Figure 10), it will finish the current task and move on to the next one. If Optimus-1 discovers that it has not yet completed the task and the task has not failed(as shown in Figure 11), it will continue executing the task.

Figures 12 and 13 illustrate the advantages of planning with knowledge. With the Hierarchical Directed Knowledge Graph, we can generate a high-quality plan in one step and dynamically adjust the plan based on current visual signals.

- Figure 12 demonstrates the importance of knowledge. For a long-horizon task such as “Mine 1 diamond 📍,” Optimus-1 first generates a plan based on the Hierarchical Directed Knowledge Graph. However, this plan needs to be adjusted based on the current visual signals. For example, in this figure, Optimus-1 appears in a cave, so the primary task is not to “chop a tree” but to “leave the cave” first. Only after exiting the cave can Optimus-1 proceed with the initial plan.
- Figure 13 demonstrates the high efficiency of our method. Agents like MP5 [32] and Voyager [46] use an iterative planning approach, which is very time-consuming, generating the final plan step by step. During this process, agent does not take any action. As shown in Figure 13, a zombie is gradually approaching the agent, but the agent is still iterating on its plan. Optimus-1, however, generates the plan in one step based on the Hierarchical Directed Knowledge Graph and makes reasonable plans based on the current visual signals.



Figure 9: The process of completing the task "Craft 1 wooden pickaxe". Optimus-1 gives wrong planning. When Optimus-1 realizes it cannot complete the task, it will replan the current task.



Figure 10: The process of completing the task "Find a cow and kill it". Hierarchical Directed Knowledge Graph indicates that having a wooden sword will make the task easier to complete. Therefore, Optimus-1 first crafts a wooden sword and then proceeds to find and kill a cow.



Figure 11: The process of completing the task "Chop tree to obtain 10 logs". Hierarchical Directed Knowledge graph indicates that no tools are needed to complete this goal. After finding a tree, Optimus-1 starts chopping it down. The task requires a substantial amount of wood, so midway through, Optimus-1 performs a reflection. The task is not yet complete but is progressing smoothly, and the result of the reflection is to continue.



Figure 12: The process of completing the task "Mine 1 diamond". Mining diamonds is a highly complex task. Diamonds can only be mined with an iron pickaxe, so an iron pickaxe must be crafted first. Crafting an iron pickaxe requires iron ingots, which are smelted from iron ore. Mining iron ore requires a stone pickaxe. Crafting a stone pickaxe requires stone, which in turn must be mined with a wooden pickaxe. Crafting a wooden pickaxe requires wooden planks and sticks. All these crafting processes require a crafting table, and smelting requires a furnace. In this case, the agent spawns at a cave, so Optimus-1 must leave the cave to chop logs.

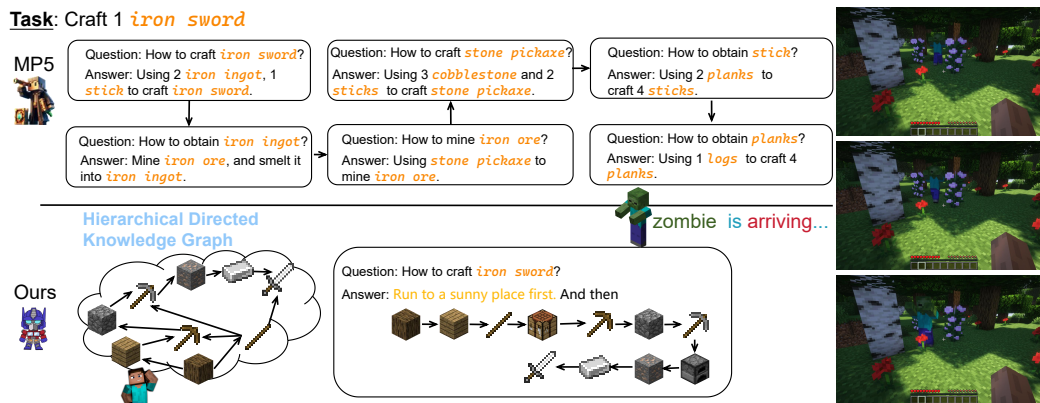


Figure 13: In this example, a zombie is slowly approaching the agent. Agents like MP5 [32] and Voyager [46] use an iterative planning strategy to generate the plan, which consumes a great deal of time and puts the agent in danger. While Optimus-1 directly generates a plan in one step based on the knowledge graph. Using the current visual information, it makes a plan to "run to a sunny place," allowing the agent to avoid danger then begin to achieve sub-goals.