STARDOJO: BENCHMARKING OPEN-ENDED BEHAVIORS OF AGENTIC MULTIMODAL LLMS IN PRODUCTION—LIVING SIMULATIONS WITH STARDEW VALLEY

Anonymous authors

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

034

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Autonomous agents navigating human society must master both production activities and social interactions, yet existing benchmarks rarely evaluate these skills simultaneously. To bridge this gap, we introduce StarDojo, a novel benchmark based on Stardew Valley, designed to assess AI agents in open-ended production-living simulations. In StarDojo, agents are tasked to perform essential livelihood activities such as farming and crafting, while simultaneously engaging in social interactions to establish relationships within a vibrant community. StarDojo features 1,000 meticulously curated tasks across five key domains: farming, crafting, exploration, combat, and social interactions. Additionally, we provide a compact subset of 100 representative tasks for efficient model evaluation. The benchmark offers a unified, user-friendly interface that eliminates the need for keyboard and mouse control, supports all major operating systems, and enables the parallel execution of multiple environment instances, making it particularly well-suited for evaluating the most capable foundation agents, powered by multimodal large language models (MLLMs). Extensive evaluations of state-of-the-art MLLMs agents demonstrate substantial limitations, with the best-performing model, GPT-4.1, achieving only a 12.7% success rate, primarily due to challenges in visual understanding, multimodal reasoning and low-level manipulation. As a user-friendly environment and benchmark, StarDojo aims to facilitate further research towards robust, open-ended agents in complex production-living environments. Code and demos can be found at https://stardojo2025.github.io/stardojo.

1 Introduction

The complexity of human society is embodied in "**production-living systems**" (Fu et al., 2023), where individuals simultaneously engage in resource-generating activities (like farming, manufacturing, and crafting) while participating in social exchanges, cultural practices, and adaptive responses to environmental changes. These systems represent the fundamental ways humans navigate the world, requiring constant balancing between productive output, resource management, social integration, and environmental adaptation. Developing human-level agentic multimodal LLMs (MLLMs) that can navigate production-living systems represents a crucial frontier in AI research, moving beyond narrow task completions such as question answering (Achiam et al., 2023), code generation (Li et al., 2022) and math reasoning (Trinh et al., 2024), towards systems that can handle the interplay between productive activities and social dynamics. This comprehensive capability is essential for Artificial General Intelligence (**AGI**), as it encompasses the holistic integration of cognitive, practical, and social intelligence that has enabled human thriving across diverse environments. However, the capability is seldom evaluated by existing interactive benchmarks such as Atari (Bellemare et al., 2013), PySC2 (Vinyals et al., 2019), MineDojo (Fan et al., 2022) and OSWorld (Xie et al., 2025).

To bridge the gap, we introduce StarDojo, a novel environment and benchmark consisting of 1,000 comprehensive tasks based on the production-living dynamics of the popular simulation game, Stardew Valley. As illustrated in Figure 1, agents undertake tasks such as clearing farmland, tilling soil, planting and watering crops, harvesting produce, raising animals, mining and foraging for resources, and crafting essential tools and equipment. Additionally, agents must explore diverse maps, collect various items, combat monsters, trade with merchants to generate income, upgrade and



Figure 1: StarDojo leverages the open-world richness of Stardew Valley to provide a diverse array of scenarios and activities. Agents are required not only to engage in a wide range of production activities, such as farming, crafting, mining, logging, and animal husbandry, but also to participate in various social events, including trading, conversing with NPCs, and even starting families. In addition, agents need to adapt to dynamic changes in time and weather, thereby reflecting the multifaceted nature of real-world living and social interaction.

expand their farms, and complete numerous in-game quests. Social progress involves participating in festivals and community events, building interpersonal relationships, and may eventually lead to marriage and raising children. The environment realistically simulates time, stamina, weather patterns, and seasonal cycles, all of which significantly impact both production tasks and social interactions, presenting further challenges for agents. To facilitate development and evaluation for researchers, we also present StarDojo-Lite, a curated subset comprising 100 core tasks that focus on the essential skills typically encountered during the game's early stages.

To facilitate the development and evaluation of agent behaviors, StarDojo offers four key features: 1) **Unified User-friendly Interface**: Provides an intuitive Python interface to interact seamlessly with Stardew Valley's game engine implemented in C#, simplifying interaction and internal state capture, eliminating the need for manual screenshots and keyboard/mouse inputs to obtain observation and action. 2) **Automated Evaluation**. Includes comprehensive evaluation scripts for all tasks, ensuring reliable and reproducible agent performance assessments. 3) **System Compatibility**. Supports major operating systems (Ubuntu, macOS, and Windows) to ensure wide accessibility. 4) **Parallelized Environments**. Allows multiple headless environment instances to run concurrently, significantly enhancing efficiency for evaluation and data collection.

Through extensive evaluations, we demonstrate that tasks within StarDojo present significant challenges even for agents with state-of-the-art MLLMs. Our assessments cover cutting-edge models, including GPT-4.1 (& mini) (OpenAI, 2025), Claude-3.7 Sonnet (Anthropic, 2025), Gemini 2.5 Pro (Google, 2025), and the open-source Llama 4 Maverick (Meta, 2025), Qwen2.5-VL-72B (Bai et al., 2025) and Gemma 3 27B (Team et al., 2025). These agents achieve performance ranging from 4% to 12.7% on the StarDojo-Lite task suite. While agents successfully complete some easy-level tasks requiring fewer than 30 steps, they exhibit near-zero success rates on medium and hard tasks demanding extended action sequences. Additionally, tasks involving long-term navigation and combat are particularly challenging. We observed that even advanced models struggle significantly with accurately identifying crucial visual elements, such as character locations, entrances, trees, NPCs, and crops, which makes them difficult to determine whether the object is already within the reach of interaction. This makes it hard for them to tell whether something is close enough to interact with, which in turn prevents them from finishing tasks that require moving across different areas or going

in and out of buildings, impeding their ability to engage in comprehensive long-term planning. We release StarDojo as an open-source environment and benchmark, providing setup instructions, robust evaluation scripts, comprehensive documentation, and baseline implementations. We hope it can facilitate research into robust, open-ended decision-making agents capable of lifelong learning and long-term planning in production-living systems.

Table 1: Comparison of StarDojo with representative existing environments. The columns indicate whether the environment supports open-ended interaction and continuous learning, evaluates the ability to perform long-term planning, reflects scenarios relevant to real-world daily life, includes production activities (e.g., farming, hunting, and crafting), simulates human society with social interactions, implements a realistic economy supporting production and social activities, and provides language-based APIs for interaction with LLMs.

Environment	Open- Endness	Long-Term Planning		Production Activities		,	Language APIs
Atari (Bellemare et al., 2013)	×	×	X	×	×	X	×
VirtualHome (Puig et al., 2018)	×	×	\checkmark	×	×	X	×
SMAC (Samvelyan et al., 2019)	×	×	X	×	×	X	×
TextWorld (Côté et al., 2019)	×	×	\checkmark	×	×	X	\checkmark
Smallville (Park et al., 2023)	√	\checkmark	\checkmark	×	\checkmark	X	\checkmark
OSWorld (Xie et al., 2025)	×	\checkmark	\checkmark	×	×	X	\checkmark
Crafter (Hafner, 2021)	√	×	X	\checkmark	×	X	×
MineDojo (Fan et al., 2022)	√	\checkmark	X	\checkmark	×	X	\checkmark
CivRealm (Qi et al., 2024)	√	\checkmark	X	\checkmark	\checkmark	\checkmark	\checkmark
StarDojo	√	√	√	√	√	√	

2 Related Work

Benchmarks for MLLMs Agents. The development of robust benchmarks for evaluating decisionmaking agents across various scenarios has been a critical focus in AI research. Traditional reinforcement learning (RL) benchmarks (Bellemare et al., 2013; Todorov et al., 2012; Guss et al., 2019; Samvelyan et al., 2019) predominantly emphasize low-level control tasks and non-realistic environments. More realistic simulators (Chang et al., 2017; Kolve et al., 2017; Li et al., 2021; Puig et al., 2023; 2018), usually focus on embodied tasks in household scenarios, lacking significant environmental dynamics and diverse social activities. Recent advancements in generative agents have enabled large language models (LLMs) to simulate human social behaviors in interactive environments (Park et al., 2023; Albrecht et al., 2022; Côté et al., 2019). Their limited action spaces, primarily restricted to dialogue, hinder engagement in broader, real-world-inspired tasks that involve production, consumption, and resource management. On the other side, while GUI benchmarks (Shi et al., 2017; Zhou et al., 2023; Koh et al., 2024; Xie et al., 2025) also provide interactive environments, they focus on short-term software manipulations. The most relevant work is MineDojo (Fan et al., 2022), which offers a diverse range of tasks within the open-ended environment of Minecraft. However, its gameplay primarily focuses on interactions with nature, with limited opportunities for human-like social interactions. Additionally, its complex 3D navigation controls pose significant challenges, particularly for LLM-based agents to complete even simple tasks, further limiting its suitability for more complex activities. Another relevant benchmark, CivRealm (Qi et al., 2024), evaluates agents' strategic decision-making at a country level within a turn-based, Civilization-like game. While CivRealm presents a variety of tasks including managing population, production, and economy, its scope remains at a macro-strategic level, distinct from StarDojo's focus on granular, individual-level decision-making. Additionally, recent works (Paglieri et al., 2024; Zheng et al., 2025) tend to evaluate MLLMs on traditionally RL environments in simple game scenarios, lacking semantic richness and social interaction.

As shown in Table 1, StarDojo addresses the limitations of prior benchmarks by providing a comprehensive, open-ended evaluation platform for decision-making agents in a dynamic environment. StarDojo allows for rich, multimodal interactions that encompass farming, trading, crafting, exploration, and social relationships. Its unique combination of complexity, realism, and diversity makes it a valuable testbed for advancing agents in real-world-like environments.

MLLMs Agents. Traditional reinforcement learning (RL) agents (Mnih et al., 2015; Lillicrap, 2015; Schulman et al., 2017; Haarnoja et al., 2018) primarily focus on low-level control and fail to leverage natural language understanding, making them unsuitable for tasks requiring complex reasoning, long-term planning, and social interactions. Recent advancements in LLM-based agents have significantly

expanded AI capabilities by integrating reasoning mechanisms such as chain-of-thought (CoT) prompting (Yao et al., 2023) and reflection (Shinn et al., 2023). Modular frameworks and multi-agent architectures have enabled LLM-based agents to achieve remarkable performance in tasks like code generation (Hong et al., 2023; Wu et al., 2023; Wang et al., 2024b) and GUI manipulation (Zheng et al., 2024; Zhang et al., 2024; Wu et al., 2024; Wang et al., 2024a). Additionally, Voyager (Wang et al., 2023) has demonstrated strong in-context lifelong learning abilities, showcasing exceptional proficiency in the open-ended world of Minecraft. However, Voyager's strong reliance on built-in APIs makes it challenging to adapt to other games. Cradle (Tan et al., 2024) successfully completes meaningful tasks across multiple commercial video games and software applications with a unified interface without the need to access APIs. Its preliminary experiments in Stardew Valley highlight the limitations of current state-of-the-art agents, particularly in handling multi-modal understanding, long-term planning, and resource management, which reveals the importance of extending Stardew Valley as a well-developed benchmark for decision-making agents.

3 STARDOJO

3.1 Introduction to Stardew Valley

Stardew Valley is a globally popular open-ended simulation RPG game where players inherit a run-down farm. Players must thoughtfully manage their farming strategies, explore the surrounding village, build meaningful relationships with villagers, and gather diverse resources to revitalize the farm. More details can be found in Appendix B.

Realistic Dynamics. The game incorporates a realistic cycle, featuring days that run from 6 AM to 2 AM, a daily energy system, and four 28-day seasons with changing weather, all of which influence both production and social activities. Success depends on effective time and energy management, as well as the ability to quickly adapt to dynamics.

Rich Production Activities. As the core gameplay of Stardew Valley, players can engage in a variety of production activities such as clearing land, growing crops, raising animals, mining, and foraging. These activities improve character skills and unlock over 100 crafting recipes for useful tools, machinery, and decorative items that significantly enhance productivity and efficiency.

Diverse Social Interaction. The game also features social interaction with 45 unique non-player characters (NPCs), each with their own personality and routines. Players can build friendships and may even date, marry, and raise children with villagers. Improving friendships not only enriches the narrative experience but also facilitates production, as villagers may send useful gifts, share recipes, or offer assistance in various activities. Rich town festivals and quests provide further opportunities for community engagement and resource gathering.

Comprehensive Economic System. Agents must engage in strategic resource management, investment, and efficient planning to generate income from both production and social activities while adapting to seasonal demands and market conditions to ensure long-term financial success.

Overall, Stardew Valley serves as an ideal environment for decision-making agents in the production—living simulation. Its well-integrated systems of time management, resource allocation, economic planning, and social interaction provide a dynamic and complex environment that requires strategic thinking and adaptability. The game's structured yet open-ended nature makes it an excellent testbed for evaluating decision-making capabilities in simulated real-world conditions.

3.2 ARCHITECTURE

As a typical commercial video game, Stardew Valley only supports human-like interaction, e.g., observing gameplay through screenshots and using keyboard and mouse to control. Additionally, the game window must remain active and in the foreground, significantly restricting automated gameplay and preventing simultaneous execution of multiple instances. As shown in Figure 2, to overcome these limitations, we introduce the carefully designed StarDojo environment, enabling efficient interaction and comprehensive evaluation of agents.

Unified User-friendly Interface. We present StarDojoMod, a novel extension built upon the Stardew Modding API (SMAPI)(Plamondon-Willard), which is a widely adopted, open-source modding framework designed specifically for Stardew Valley. SMAPI offers developers extensive APIs that expose key game events and internal states, facilitating the creation of interactive and sophisticated mods. Based on SMAPI, StarDojoMod provides structured and efficient interactions between agents and the game environment. It communicates in real-time with the Stardew Valley game engine through a socket server, granting agents direct access to rendered gameplay images, saving the

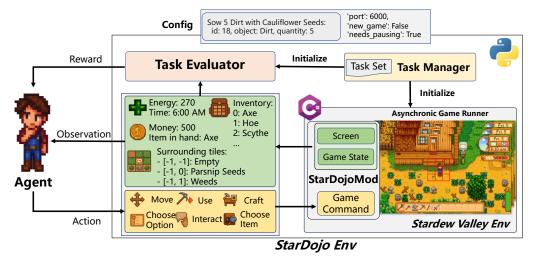


Figure 2: StarDojo environment is initiated by configurable task files. It communicates with parallel game engines through StarDojoMod to obtain internal game states and execute commands, which will be encapsulated as observations and actions by the Python Wrapper.

time-consuming screen captures, internal game states (such as character positions, statuses, and environmental information), and enabling diverse callable functions as action skills beyond traditional keyboard and mouse inputs. Moreover, we implemented a configurable pause-and-resume mechanism by directly modifying the inner states of the game, allowing the game to pause during model inference and agent planning, and resume before action execution. Inherited from SMAPI, StarDojoMod is implemented in C# to be consistent with the game engine. To enhance ease-of-use and accessibility of the environment, we provide a user-friendly Python Wrapper based on the StarDojoMod for observation retrieval, action execution, and task customization, empowering users to engage with the StarDojo environment effortlessly.

System Compatibility. Stardew Valley is one of the few games that can be played on all mainstream operating systems (Linux, macOS and Windows). We also ensured the compatibility of StarDojoMod and the Python Wrapper, enabling the entire environment to run seamlessly across different systems.

Parallel Execution. Our architecture is designed for scalability and parallel execution. Each instance of Stardew Valley is independently managed through unique ports, enabling simultaneous control of multiple game instances without interference. Communication efficiency is further enhanced through the use of shared memory, reducing observation retrieval time to as little as 30 ms. Furthermore, StarDojoMod supports headless operation through the X Virtual Framebuffer (Xvfb), enabling compatibility with Linux systems without graphical interfaces, thus broadening accessibility across diverse hardware and system configurations.

3.3 OBSERVATION AND ACTION SPACES

Observation Space. StarDojo offers a comprehensive observation space that integrates both visual and textual modalities to accommodate a wide range of agent architectures. Each observation includes a gameplay screenshot alongside detailed textual information describing the game state. This textual state contains character status (such as health and energy), local tile information $(n \times n)$ tiles surrounding the agent), and global information (such as time, weather, and the positions of NPCs and buildings). By combining visual and textual observations, StarDojo ensures that agents can leverage both high-level context and fine-grained environmental cues for more robust and informed decision-making. For our experiments, we selected a subset of this information to fairly evaluate agent behaviors. Full details of the observation space are provided in Appendix C.1.

Action Space. The action space in StarDojo is designed to encompass the full range of activities that can be performed in the original Stardew Valley using a keyboard and mouse. Actions are carefully abstracted to eliminate redundant operations while retaining the core decision-making challenges inherent to the game. We define ten fundamental actions, which together are sufficient to cover most gameplay scenarios: move(x, y), use(direction), interact(direction), $choose_item(slot_index)$, $attach_item(slot_index)$, $detach_item()$, craft(item), $choose_option(option_index$, quantity, direction) and $menu(option, menu_name)$. Detailed descriptions of the action space can be found in Appendix C.2.

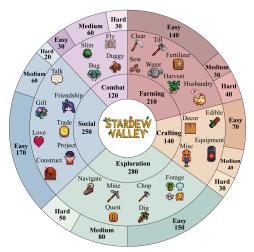


Figure 3: Distribution of 1000 tasks across five categories: Farming, Crafting, Exploration, Combat and Social in StarDojo, each with Easy, Medium, and Hard difficulties.

Table 2: Task statistics of StarDojo-Lite. The task suite is made up of the most representative early-stage tasks from each category.

Category	Easy	Medium	Hard	Total
Farming	14	3	4	21
Crafting	7	4	3	14
Exploration	15	8	5	28
Combat	3	6	3	12
Social	17	6	2	25
Total	56	23	21	100

3.4 TASKS

As shown in Figure 3, we carefully curate 1000 tasks to benchmark agents' various behaviors in StarDojo. These tasks are divided into five distinct categories: Farming, Crafting, Exploration, Combat and Social, which cover most of the production-living activities in the early and middle stages of the game. Each task is classified into three difficulties: easy, medium, and hard. For easy-level tasks, agents are provided with all necessary items or tools from the start (e.g., mature crops ready for harvest, ingredients for crafting). Agents are initialized near the target location (e.g., inside a shop with enough budget for trading). These tasks primarily test the agent's basic ability to complete atomic activities. For medium-level tasks, agents need to fulfill the prerequisites of the tasks on their own, such as planting and harvesting crops, gathering crafting ingredients, traveling from the farm to the target area, and earning sufficient funds to purchase specific items. Many of these resources can be acquired through multiple sources and methods. This flexibility gives agents considerable freedom in choosing their strategies. Hard-level tasks typically require several in-game days to complete and are often composed of multiple medium-level tasks. Agents have greater freedom to allocate their time and energy each day, choosing how to prioritize activities and sequence actions. Success often demands strategic long-term planning, adaptive decision-making, and efficient resource management, as there are multiple viable paths to achieving the goal.

As shown in Table 2, to facilitate efficient agent evaluation, we curate a representative smaller task suite, called StarDojo-Lite, comprising 100 core tasks from the full task collection, balancing coverage and practicality. This lite task set covers most of the representative activities in the early stage of the game. Detailed descriptions of the task set can be found in Appendix C.3.

Initial Config and Setup. Some tasks require the agent to possess specific equipment (e.g., a sword), have certain items (e.g., sufficient crop seeds in inventory), or have completed particular game progress (e.g., unlocking the mines). To establish the initial state for each task, we provide multiple saved game files reflecting various stages of progression, along with specialized task-specific functions utilizing StarDojoMod commands. At the start of each task, StarDojo automatically loads the corresponding saved game and executes these task-specific commands, ensuring all necessary prerequisites, such as items, equipment, skill levels, date and farm status, are appropriately configured.

Automatic Evaluation. Given the extensive number of existing tasks and the vast potential for future additions, it is essential to establish a reusable, scalable, adaptive, and efficient evaluation mechanism. StarDojo addresses this need by implementing a general evaluation system that continuously monitors task progression and provides immediate reward feedback to agents. Specifically, after agents finish executing actions at each step, the evaluator is invoked to determine whether tasks have been successfully completed or have reached the predefined maximum number of steps. To achieve this, the evaluator maintains the previous game state information, compares it with the current state obtained via StarDojoMod, identifies incremental changes indicating progress, and accurately assesses task completion criteria based on these observed differences. By leveraging a comprehensive observation

Table 3: Success rates(%) and standard deviation of agents with different base models on StarDojo-Lite task set, ranging over five categories (Farming, Crafting, Exploration, Combat and Social) and three levels of difficulty. Each task is evaluated over three runs.

Model	Task	Farming	Crafting	Exploration	Combat	Social	Total
	Easy	31.0 ±4.1	52.4 ±8.3	15.6 ±3.9	11.1 ±19.3	7.8±6.8	21.4 ±1.8
GPT-4.1	Medium	$0.0 {\pm} 0.0$	$0.0 {\pm} 0.0$	8.3 ± 7.2	$0.0 {\pm} 0.0$	$0.0 {\pm} 0.0$	2.7 ± 2.3
Gr 1-4.1	Hard	$0.0 {\pm} 0.0$	0.0 ± 0.0	$0.0 {\pm} 0.0$	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
	Total	20.6 ±2.8	26.2 ±4.1	10.7 ± 0.0	2.8 ± 4.8	5.3 ± 4.6	12.7 ± 0.6
	Easy	26.2±4.1	47.6±8.3	6.7±6.7	0.0 ± 0.0	11.8 ±5.9	17.9±1.8
Gemini 2.5 Pro	Medium	$0.0 {\pm} 0.0$	$0.0 {\pm} 0.0$	8.3 ± 7.2	$0.0 {\pm} 0.0$	$0.0 {\pm} 0.0$	2.7 ± 2.3
Geiiiiii 2.3 1 10	Hard	$0.0 {\pm} 0.0$	$0.0 {\pm} 0.0$	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
	Total	17.5 ± 2.8	23.8 ± 4.1	6.0 ± 2.1	0.0 ± 0.0	8.0 ± 4.0	10.7 ± 0.6
	Easy	31.0 ±4.1	28.6±0.0	8.9 ± 10.2	0.0 ± 0.0	7.8±3.4	16.1±4.7
Claude 3.7 Sonnet	Medium	0.0 ± 0.0	$0.0 {\pm} 0.0$	16.7 ±7.2	0.0 ± 0.0	0.0 ± 0.0	5.3 ±2.3
Claude 3.7 Somet	Hard	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	$0.0 {\pm} 0.0$	$0.0 {\pm} 0.0$	0.0 ± 0.0
	Total	20.6 ±2.8	14.3 ± 0.0	9.5 ± 5.5	$0.0 {\pm} 0.0$	5.3 ± 2.3	10.3 ± 2.5
	Easy	26.2±4.1	4.8±8.3	4.4±7.7	11.1 ±19.3	7.8±3.4	11.3±2.7
GPT-4.1 mini	Medium	0.0 ± 0.0	0.0 ± 0.0	8.3 ± 7.2	$0.0 {\pm} 0.0$	$0.0 {\pm} 0.0$	2.7 ± 2.3
O1 1-4.1 IIIIII	Hard	$0.0 {\pm} 0.0$	$0.0 {\pm} 0.0$	$0.0 {\pm} 0.0$	$0.0 {\pm} 0.0$	0.0 ± 0.0	0.0 ± 0.0
	Total	17.5 ± 2.8	2.4 ± 4.1	4.8 ± 5.5	2.8 ±4.8	5.3 ± 2.3	7.0 ± 1.7
	Easy	28.6 ± 0.0	$28.6{\pm}0.0$	$0.0 {\pm} 0.0$	0.0 ± 0.0	7.8 ± 3.4	13.1 ± 1.0
Llama 4 Maverick	Medium	0.0 ± 0.0	0.0 ± 0.0	4.2 ± 7.2	$0.0 {\pm} 0.0$	$0.0 {\pm} 0.0$	1.3 ± 2.3
Liama + Maverier	Hard	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	$0.0 {\pm} 0.0$	$0.0 {\pm} 0.0$	0.0 ± 0.0
	Total	19.1 ± 0.0	14.3 ± 0.0	1.2 ± 2.1	0.0 ± 0.0	5.3 ± 2.3	7.7 ± 0.6
	Easy	16.7±8.3	9.5 ± 8.3	13.3 ± 13.3	$0.0 {\pm} 0.0$	9.8 ± 3.4	11.9±5.5
Owen2.5 VL 72B	Medium	0.0 ± 0.0	$0.0 {\pm} 0.0$	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
QWCII2.5 VL 72B	Hard	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
	Total	11.1 ± 5.5	4.8 ± 4.1	7.1 ± 7.1	$0.0 {\pm} 0.0$	6.7 ± 2.3	6.7 ± 3.1
	Easy	9.5±4.1	0.0 ± 0.0	2.2±3.9	0.0 ± 0.0	11.8 ±0.0	6.6±2.1
Gemma 3 27B	Medium	0.0 ± 0.0	0.0 ± 0.0	4.2 ± 7.2	0.0 ± 0.0	$0.0{\pm}0.0$	1.3 ± 2.3
Gennia 3 2/D	Hard	0.0 ± 0.0	0.0 ± 0.0	$0.0 {\pm} 0.0$	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
	Total	6.4 ± 2.8	0.0 ± 0.0	2.4 ± 2.1	$0.0 {\pm} 0.0$	8.0 ± 0.0	4.0 ± 1.0

space for incremental progress tracking, our evaluation approach effectively mitigates discrepancies caused by varying initial conditions, ensuring robust adaptability across diverse tasks.

4 EMPIRICAL STUDIES

In this section, we benchmark the current state-of-the-art MLLM-based agents on StarDojo and provide a comprehensive study and analysis of agents' behaviors and limitations.

Agentic MLLM Baselines. We evaluated seven cutting-edge agentic MLLMs on StarDojo-Lite task set: GPT-4.1 series (OpenAI, 2025)(gpt-4.1-2025-04-14 & gpt-4.1-mini-2025-04-14), Claude 3.7 Sonnet (Anthropic, 2025) (claude-3-7-sonnet-20250219), Gemini 2.5 Pro (Google, 2025) (gemini-2.5-pro-preview-03-25), from the closed-source community and Llama 4 Maverick (Meta, 2025) (Llama-4-Maverick-17B-128E-Instruct), Qwen2.5 VL (Bai et al., 2025) (qwen2.5-vl-72b-instruct) and Gemma 3 (Team et al., 2025) (gemma-3-27b-it) from the open-source community.

Settings. If not mentioned explicitly, all experiments are conducted under the following settings: Agents have access to both visual and textual observations for their decision-making process. Visual observations are provided at a resolution of 720P (1280×720). Textual observations include 7×7 agent-centered surrounding information and other global information like time, date and budget. In addition to receiving observations from the current timestep, agents are also provided with history information from the previous timestep, enabling agents to reflect on past states and facilitating consistency in decision-making. Agents can output at most two skills as an action to be executed sequentially. After executing all the actions, the environment is paused until the agent outputs the next action. Each task is evaluated over three runs with a heuristic maximum of 30, 50 and 150 steps for easy, medium and hard level tasks. More details of settings are provided in Appendix D

4.1 QUALITATIVE ANALYSIS

Overall Results. As shown in Table 3, a clear gap between flagship commercial models and open-source models can be observed. GPT-4.1, Gemini 2.5 Pro, and Claude 3.7 Sonnet all exceed a 10% overall success rate, with GPT-4.1 achieving the highest at 12.7%. In contrast, all open-source models remain below 8%, with Llama 4 Maverick performing best, largely due to its larger model size. Most successful completions are limited to easy tasks, whereas all models struggle significantly with medium and hard tasks, achieving near-zero success rates due to increased task complexity and longer sequences of required actions. Models show some proficiency in farming and crafting tasks but exhibit considerable difficulty in exploration, combat, and social interactions.

Low-Level Control as Tool-Use. All models demonstrate reasonable ability to control agents with non-trivial movements and interactions with actions clearly specified in the prompts. Larger models such as GPT-4.1, Gemini 2.5 Pro, Claude 3.7 Sonnet and Llama 4 Maverick perform significantly better than others on Crafting tasks, highlighting their superior tool-use capabilities and in-context understanding when handling more complex function calls. This contrast explains the weaker performance of smaller models like GPT-4.1 Mini and Gemma 3 27B, which often struggle to provide valid and accurate parameters when calling actions like *choose_option* and *craft*. However, all models continue to struggle in fast-paced tasks such as Combat, which demand dynamic control.

Navigation is the Key Bottleneck. We found that all models struggled most severely with navigation, which is the core skill required in the game. Due to limited image understanding, they consistently failed to accurately identify and locate target objects, building entrances or exits, and map transitions, which usually only appear in the provided image rather than the textual information. These short-comings heavily impaired tasks that involve moving across scenarios or exploring areas beyond the immediate field of view. As a result, success rates were particularly low in Exploration and Social tasks, both of which often require cross-map navigation to locate targets. In contrast, agents are less influenced and perform relatively better in easy-level Farming tasks, since these are usually confined to a small farmland where target products remain within visual proximity and are directly accessible.

Long-Horizon Tasks Remain Far From Solved. Medium and hard-level tasks can typically be decomposed into multiple simpler subtasks involving different category combinations. However, given the low success rates even on easy-level tasks, which serve as the atomic activities in the benchmark, all state-of-the-art models remain far from being able to complete complex long-horizon tasks in StarDojo. This gap reveals substantial potential for more capable MLLM models.

4.2 ERROR ANALYSIS

To further investigate models' behaviors, we conducted the error analysis on GPT-4.1-based agent to identify and categorize failure modes in StarDojo-Lite task set. For tasks with multiple errors, we report primary errors that block further progress.

The error analysis reveals several key failure modes, with the most significant being limited **Visual Understanding**, which accounts for 42% of all errors. The model often struggles to reliably recognize target objects, many of which are only around 10x10 pixels shown in the image. And even when detected, it frequently fails to accurately determine their position relative to the character, hindering effective navigation and interaction. Additionally, models often struggle to interpret the status of tiles, even those immediately adjacent to the character, such

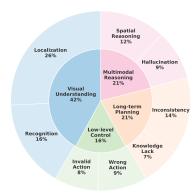


Figure 4: Error Analysis.

as whether a tile is tilled, seeded, watered, or obstructed. The second major source of failure, at 21%, is **Multimodal Reasoning**, where the model overly relies on its limited visual perception instead of integrating both visual and textual inputs, leading to confusion even when clear positional or status information is available in the text. Additionally, the model sometimes hallucinates task progress by incorrectly assuming the success of previous actions. **Long-Term Planning** issues also contribute to 21% of failures, as the model, while capable of proposing reasonable subtasks, often abandons plans prematurely and switches strategies inconsistently, which undermines progress on longer-horizon tasks; this is sometimes compounded by insufficient domain knowledge, such as not knowing a crafting recipe or which NPC to interact with. Lastly, **Low-Level Control** accounts for a smaller but persistent 16% of errors, where the model, despite generally choosing appropriate actions, occasionally selects incorrect or suboptimal ones, demonstrating instability in fine-grained execution, with occasional formatting or parameter errors.

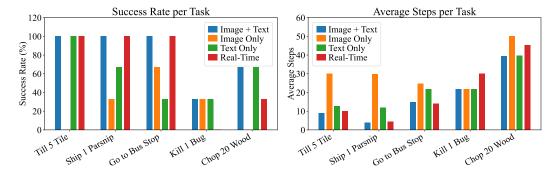


Figure 5: Ablation results of GPT-4.1-based agents on five representative tasks under four different settings: *Image* + *Text* (default setting), *Image Only*, *Text Only*, and *Real-Time* (without game pausing). Except for *Chop 20 Wood*, which is a medium-level task with a maximum of 50 steps, the remaining four are easy-level tasks capped at 30 steps. Each task is evaluated over three runs.

4.3 ABLATION STUDIES

To demonstrate the flexible customization capabilities of StarDojo and enable more comprehensive evaluations of agent performance, we benchmark GPT-4.1 under the following three additional settings using five representative basic tasks. 1) **Image Only**: Agents receive only visual observations without textual state information, evaluating their capability to rely solely on visual cues. 2) **Text Only**: Agents receive textual state observations without visual input, restricting the agent's perception to local 7×7 grid-based information. This setting is particularly relevant for LLM agents. 3) **Real-time**: The environment progresses continuously during action generation, simulating real gameplay conditions where agents must plan and respond without game pausing.

As illustrated in Figure 5, removing textual input (*Image Only*) significantly affects agents' performance across all tasks, reflecting the defect of base models' poor visual-based control, emphasizing textual information's importance in grounding detailed action decisions for the current stage of agents. On the other side, eliminating visual input (*Text Only*) substantially reduces success in tasks that require navigation like *Ship 1 Parsnip* and *Go to Bus Stop*, demonstrating the essential contribution of visual cues to spatial reasoning and movement. Disabling the feature to pause the environment (*Real-time*) remarkably affects performance across tasks demanding timely reactions or prolonged action sequences. For example, in combat scenarios like *Kill 1 Bug*, the target (bug) continues moving during model inference, which can take over 10 seconds per request for GPT-4.1 via API. By the time the action is executed, the bug has often moved far from its previous position, rendering the action ineffective. Similarly, in long-horizon tasks like *Chop 20 Wood*, the in-game clock advances continuously during inference. With pausing enabled, the task may be completed in just 2 in-game hours; without pausing, it can take over 12 in-game hours, frequently pushing completion into the night or spanning multiple in-game days. These findings highlight the practical importance of real-time evaluation, an aspect often overlooked in prior benchmarks.

5 LIMITATIONS AND CONCLUSION

Limitations. This work has several potential limitations. 1) Although StarDojo is an open-source environment and benchmark, users need an official copy of Stardew Valley to run it. 2) Fishing, an optional gameplay, is currently not included in StarDojo, due to its nature as a complex, real-time mini-game. The operations are independent of other game controls, which is not applicable to evaluate MLLMs. 3) The benchmark primarily focuses on early- and mid-game content within the main valley map. Other advanced areas, such as the Desert and Ginger Island, are not yet supported. 4) Due to budget limitations, our evaluations were conducted mainly on StarDojo-Lite with 7 models. A broader assessment across diverse tasks and models remains an important direction for future work.

Conclusion. Overall, we introduce StarDojo, a novel environment and benchmark designed to evaluate the open-ended behaviors of MLLM agents in Stardew Valley. StarDojo bridges the gap in existing environments by enabling comprehensive assessment of agents across both production and daily living activities within a simulated nature and society. Featuring a set of diverse tasks, StarDojo exposes significant challenges in current agents' visual understanding, multimodal reasoning, long-term planning, and real-time inference, highlighting key areas for future research and development.

REPRODUCIBILITY STATEMENT

StarDojo environment and benchmark are already open-sourced and can be found at https://stardojo2025.github.io/stardojo. We also provide the code in the supplementary material.

LARGE LANGUAGE MODELS USAGE STATEMENT

492 493 494

486

487 488

489

490 491

> In this work, we utilized LLMs exclusively for language polishing and grammar refinement. LLMs were not employed for tasks such as information retrieval, discovery and research ideation.

495 496

499

500

501

502

504

505

506

507

509

510

511

512

513

514 515

516

517 518

519

521

522

523

524

525

526

527

528

529

530 531

532

534

535

536

538

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- Joshua Albrecht, Abraham Fetterman, Bryden Fogelman, Ellie Kitanidis, Bartosz Wróblewski, Nicole Seo, Michael Rosenthal, Maksis Knutins, Zack Polizzi, James Simon, et al. Avalon: A benchmark for rl generalization using procedurally generated worlds. Advances in Neural Information Processing Systems, 35:12813–12825, 2022.
- Anthropic. Claude 3.7 sonnet and claude code. https://www.anthropic.com/news/ claude-3-7-sonnet, 2025. Accessed: 2025-05-10.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. arXiv preprint arXiv:2502.13923, 2025.
- Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. Journal of Artificial Intelligence Research, 47: 253–279, 2013.
- Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. arXiv preprint arXiv:1709.06158, 2017.
- Marc-Alexandre Côté, Akos Kádár, Xingdi Yuan, Ben Kybartas, Tavian Barnes, Emery Fine, James Moore, Matthew Hausknecht, Layla El Asri, Mahmoud Adada, et al. Textworld: A learning environment for text-based games. In Computer Games: 7th Workshop, CGW 2018, Held in Conjunction with the 27th International Conference on Artificial Intelligence, IJCAI 2018, Stockholm, Sweden, July 13, 2018, Revised Selected Papers 7, pp. 41–75. Springer, 2019.
- 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.
- Jingying Fu, Qiang Gao, Dong Jiang, Xiang Li, and Gang Lin. Spatial-temporal distribution of global production–living–ecological space during the period 2000–2020. Scientific Data, 10(1): 589, 2023.
- Google. Gemini 2.5: Our most intelligent ai model. https://blog.google/technology/ google-deepmind/gemini-model-thinking-updates-march-2025/ #gemini-2-5-thinking, 2025. Accessed: 2025-05-10.
- 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.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In International conference on machine learning, pp. 1861–1870. PMLR, 2018.

542

543

544

546

547

548

549

550

551

552 553

554

555

558

559

561

562

563

564

565 566

567

568 569

570

571 572

573

574

575

576

577

578

579

580 581

582

583

584

585

586

588

589

590

- 540 Danijar Hafner. Benchmarking the spectrum of agent capabilities. arXiv preprint arXiv:2109.06780, 2021.
 - Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multi-agent collaborative framework. arXiv preprint arXiv:2308.00352, 2023.
 - Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. arXiv preprint arXiv:2401.13649, 2024.
 - Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Matt Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, et al. Ai2-thor: An interactive 3d environment for visual ai. arXiv preprint arXiv:1712.05474, 2017.
 - Chengshu Li, Fei Xia, Roberto Martín-Martín, Michael Lingelbach, Sanjana Srivastava, Bokui Shen, Kent Vainio, Cem Gokmen, Gokul Dharan, Tanish Jain, et al. igibson 2.0: Object-centric simulation for robot learning of everyday household tasks. arXiv preprint arXiv:2108.03272, 2021.
 - Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode. Science, 378(6624):1092–1097, 2022.
 - TP Lillicrap. Continuous control with deep reinforcement learning. arXiv preprint arXiv:1509.02971, 2015.
 - Meta. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation. https://ai.meta.com/blog/llama-4-multimodal-intelligence/, 2025. Accessed: 2025-05-10.
 - Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. nature, 518(7540):529–533, 2015.
 - OpenAI. Introducing gpt-4.1 in the api. https://openai.com/research/gpt-4-1, 2025. Accessed: 2025-05-10.
 - Davide Paglieri, Bartłomiej Cupiał, Samuel Coward, Ulyana Piterbarg, Maciej Wolczyk, Akbir Khan, Eduardo Pignatelli, Łukasz Kuciński, Lerrel Pinto, Rob Fergus, et al. Balrog: Benchmarking agentic llm and vlm reasoning on games. arXiv preprint arXiv:2411.13543, 2024.
 - Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th* annual acm symposium on user interface software and technology, pp. 1-22, 2023.
 - Jesse Plamondon-Willard. SMAPI: Stardew Modding API. https://smapi.io/. Accessed: 2025-07-09.
 - Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba. Virtualhome: Simulating household activities via programs. In *Proceedings of the IEEE conference* on computer vision and pattern recognition, pp. 8494–8502, 2018.
 - Xavier Puig, Eric Undersander, Andrew Szot, Mikael Dallaire Cote, Tsung-Yen Yang, Ruslan Partsey, Ruta Desai, Alexander William Clegg, Michal Hlavac, So Yeon Min, et al. Habitat 3.0: A co-habitat for humans, avatars and robots. arXiv preprint arXiv:2310.13724, 2023.
 - Siyuan Qi, Shuo Chen, Yexin Li, Xiangyu Kong, Junqi Wang, Bangcheng Yang, Pring Wong, Yifan Zhong, Xiaoyuan Zhang, Zhaowei Zhang, et al. Civrealm: A learning and reasoning odyssey in civilization for decision-making agents. arXiv preprint arXiv:2401.10568, 2024.
 - Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. arXiv preprint arXiv:1902.04043, 2019.

- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
 - Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. World of bits: An open-domain platform for web-based agents. In *International Conference on Machine Learning*, pp. 3135–3144. PMLR, 2017.
 - Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. Reflexion: language agents with verbal reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=vAElhFcKW6.
 - Weihao Tan, Wentao Zhang, Xinrun Xu, Haochong Xia, Gang Ding, Boyu Li, Bohan Zhou, Junpeng Yue, Jiechuan Jiang, Yewen Li, et al. Cradle: Empowering foundation agents towards general computer control. In *NeurIPS 2024 Workshop on Open-World Agents*, 2024.
 - Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*, 2025.
 - Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In 2012 IEEE/RSJ international conference on intelligent robots and systems, pp. 5026–5033. IEEE, 2012.
 - Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. Solving olympiad geometry without human demonstrations. *Nature*, 625(7995):476–482, 2024.
 - Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *nature*, 575(7782):350–354, 2019.
 - 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.
 - Junyang Wang, Haiyang Xu, Haitao Jia, Xi Zhang, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. Mobile-agent-v2: Mobile device operation assistant with effective navigation via multi-agent collaboration. *arXiv preprint arXiv:2406.01014*, 2024a.
 - Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024b.
 - Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023.
 - Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao Yu, and Lingpeng Kong. OS-copilot: Towards generalist computer agents with self-improvement. *arXiv* preprint arXiv:2402.07456, 2024.
 - Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Jing Hua Toh, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *Advances in Neural Information Processing Systems*, 37:52040–52094, 2025.
 - Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. ReAct: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.
 - Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang, Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qingwei Lin, Saravan Rajmohan, et al. UFO: A UI-focused agent for Windows OS interaction. *arXiv* preprint arXiv:2402.07939, 2024.

Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-exemplar prompting with memory for computer control. In *ICLR*, 2024.

Xiangxi Zheng, Linjie Li, Zhengyuan Yang, Ping Yu, Alex Jinpeng Wang, Rui Yan, Yuan Yao, and Lijuan Wang. V-mage: A game evaluation framework for assessing visual-centric capabilities in multimodal large language models. *arXiv preprint arXiv:2504.06148*, 2025.

Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.

A APPENDIX

B INTRODUCTION TO STARDEW VALLEY

Stardew Valley is an open-ended simulation RPG game where the player inherits a run-down farm and works to restore it. The game is presented in a 2D top-down perspective, offering a wider field of view compared to 3D games, which helps reduce operational difficulty and makes navigation smoother. The game does not have mandatory main storyline tasks and players can freely explore and live in the open-ended world.

B.1 REALISTIC GAMEPLAY MECHANISM

Stardew Valley is an open-ended simulation RPG game where the player inherits a run-down farm and works to restore it. The game is presented in a 2D top-down perspective, offering a wider field of view compared to 3D games, which helps reduce operational difficulty and makes navigation smoother. The game does not have mandatory main storyline tasks and players can freely explore and live in the open-ended world.

There are many mechanisms in Stardew Valley to simulate the real world, which raises additional challenges for players to plan reasonably according to these rules during the gameplay.

Time and Daily Routine. Each in-game day lasts about 14 minutes in real-time, running from 6 AM to 2 AM. Players must balance production, and social activities before exhaustion sets in. Staying up past 12 AM reduces energy the next day, and if players don't sleep by 2 AM, they will pass out and lose gold or wake up at the clinic. Nightfall typically occurs between 6 and 8 PM, depending on the season, as the outdoors gradually darkens over time.

Energy Management. At the beginning, players have 270 energy points per day. Every action, from farming to mining, consumes energy. Energy can be restored by sleeping or eating food. Strategic time and energy management are essential to maximize productivity and also the main challenge in Stardew Valley.

Seasons and Weather Effects. The game also has four seasons each lasting 28 days and affecting crop growth, fish availability, and town events. Weather varies daily, with rain saving time on watering crops, storms potentially damaging them, and snow limiting farming options. Some crops and activities are only available in specific seasons, requiring careful planning.

B.2 PRODUCTION

Simulating real-world society, production activities are the main gameplay within the game, where players improve life quality through labor output.

Labor Activities. The game world is vast and diverse, featuring locations like Pelican Town, the Mines, the Beach and the forest, each offering unique activities such as mining, foraging, and fishing. Farming is central to gameplay, requiring players to plant, water, harvest crops and raise animals like cows and chickens for valuable products. Fishing and foraging provide additional resources, with seasonal variations and rare finds. Mining is crucial for gathering ores and materials, with progressively challenging levels and combat against monsters.

Skill Progression and Crafting. As players engage in these activities, they gradually improve their skills in Farming, Mining, Foraging, Fishing, and Combat. Gaining experience in each skill unlocks new crafting recipes, efficiency boosts, and profession choices that provide specialized benefits. Crafting is an essential part of progression, allowing players to create tools, machines, and decorations that enhance farm efficiency and exploration. In addition to upgrading tools and structures, they can fully customize their farm and home, arranging decorations and personalizing interiors to create their ideal living space.

B.3 SOCIETY

Besides production, engaging with the community is a core aspect of the game.

Festivals and Quests. Every season features unique festivals and events, such as the Egg Festival, Stardew Valley Fair, and Winter Star Festival. These events provide mini-games, rare items, and opportunities to strengthen relationships with villagers, adding depth to the community experience. Alongside festivals, quests play a crucial role in guiding players through different aspects of the game. NPCs also post daily requests on the town board, offering gold and friendship points for completing specific tasks. By participating in festivals and completing quests, players engage more deeply with the world, fostering a sense of purpose and progression throughout their journey.

Relationship and Friendship System. The game also provides various NPCs with unique personalities, schedules, and heart events. Players can be friend them, give gifts, and even date or marry eligible characters. Higher friendship levels unlock new interactions, cutscenes, and benefits, such as helpful spouses assisting with farm chores.

Economy System. Stardew Valley features a comprehensive economic system where players generate income through various means, including farming, fishing, mining, and crafting. Players must manage their resources, reinvest in production, and make strategic decisions to ensure financial growth. The economy fluctuates based on seasonal demand, production choices, and market conditions. Strategic planning, investment in high-value goods, and efficient time management are essential to achieving long-term prosperity. The economic system provides depth and challenges players to optimize their approach to wealth generation and sustainability.

Stardew Valley serves as an ideal benchmark for decision-making agents in Production–Living Simulation. Its well-integrated systems of time management, resource allocation, economic planning, and social interaction provide a dynamic and complex environment that requires strategic thinking and adaptability. The game's structured yet open-ended nature makes it an excellent testbed for evaluating decision-making capabilities in simulated real-world conditions.

C STARDOJO ENVIRONMENT

Our environment is developed based on the game Stardew Valley, which offers comprehensive official mod development documentation. A large community of game enthusiasts has created and open-sourced their own mods. This provides significant convenience, as we can build a completely new mod on top of existing open-source mods to achieve the interaction between the LLM agent, RL agent, and the game.

Stardew Valley is available for purchase on Steam, and users must first buy and install the game to use our environment. Since our mod heavily relies on the official mod framework SMAPI, users will need to install both SMAPI and the mod we developed. This process can be easily carried out on Windows, macOS, and Linux systems with a graphical interface. Additionally, we have ensured compatibility for Linux systems without a graphical interface. After purchasing the game, users can install Stardew Valley through Steamcmd, the command-line tool of the Steam client to install and update dedicated servers for Steam games. We also provide installation commands for SMAPI, allowing users to install it directly. To use our developed mod, users simply need to copy it to the designated mod folder. Since Stardew Valley is a graphical application, it cannot be directly opened on systems without a graphical interface. To address this limitation, we use the X virtual framebuffer (Xvfb), which supports all graphical operations in virtual memory without displaying any screen output. This directly ensures that our environment is compatible with various system and hardware configurations.

The interaction between our algorithm and Stardew Valley is achieved through the mod we developed, rather than simulating keyboard and mouse inputs. This approach allows us to open multiple game instances and control each one independently. Specifically, when launching multiple games, we assign a unique port number to each one. Actions provided by the algorithms being trained or tested are transmitted through this port to the corresponding game, and the returned information is received via the same port. Moreover, we use shared memory for each port to improve communication efficiency when necessary. By ensuring that the actions executed and information transmitted in all game processes do not interfere with one another, we enable parallel training and testing. To help researchers efficiently establish the initial state for each task, we provide the simulator APIs to configure the game environment.

C.1 OBSERVATION SPACE

StarDojo provides a rich and flexible observation space to accommodate the diverse needs of various agents. The observation space includes both visual and textual observations, which can be customized and extended by developers to suit their specific goals.

Visual Observation. StarDojo leverages the game engine to directly retrieve rendered gameplay images, eliminating the need for inefficient screen-capture methods. This approach ensures high-fidelity visual observations that are consistent with the gameplay environment. The resolution of these images can be configured dynamically, ranging from 360P to 4K. Notably, the resolution scaling is not merely a resizing operation; as the image size increases, the field of view also expands, providing agents with a broader perspective of the game world. This feature is particularly useful for tasks requiring detailed spatial awareness or long-term planning.

Textual Observation. While visual observations are essential, recent work (Tan et al., 2024) highlights the challenges faced by state-of-the-art MLLMs in accurately interpreting the unique art style and precise manipulation requirements of Stardew Valley. To address this limitation, StarDojo provides structured textual observations through built-in APIs. These observations are designed to complement visual data, offering agents a more comprehensive understanding of the game state. The textual observations are organized into four main categories:

- Character Information: Includes health, energy, gold, position, inventory, location, facing direction, professions, skills (e.g., Farming, Mining, Combat, Fishing, Foraging), dating partners and spouse.
- Surrounding Information: Includes tile information within an N×N grid centered on the player's position. Each tile in the grid contains detailed information about its contents and properties, such as debris, crops, NPCs, exits, buildings, furniture, terrain features, and other properties. The grid size can be adjusted to balance granularity and computational efficiency.
- Map-level Global Information: Provides global information within the current location map. Details such as crops, exits, NPCs, buildings, shop counters, furniture in rooms, and animals or pets on the farm are included as part of the map-level global information. By providing a structured overview of the entire map, this information enables the model to locate and reason about specific targets more efficiently, while reducing the need for exhaustive exploration of the environment.
- Game-level Global Information: Provides global information across the entire game. Game states such as time, day of month, season, year and weather are included in this part, along with other necessary information like the current menu information if any menu is showing.

As shown in our experiment, we selected a subset of the available information as the general configuration. Specifically, we adopted character information, surrounding information, and gamelevel global information to construct the observation space for Stardojo. Although map-level global information can significantly boost agent performance in certain tasks, our goal is to evaluate the agent's exploration ability based on visual input. Since map-level global information offers shortcuts for locating specific targets, it is deliberately excluded from our experimental setup. This combination of visual and textual observations ensures that agents have access to both high-level contextual information and fine-grained environmental details, enabling more robust and informed decision-making.

The observation space captures a comprehensive snapshot of the game's state in a JSON-like structure, with an additional screenshot RGB map. All relevant details are organized within the following nested objects. Figure 6 shows an example of the screenshot included in the observation.

Observation Space Format

The observation space consists of the following structured fields:

- Health: Integer representing the agent's current health.
- Energy: Float indicating the agent's current energy level.
- Money: Integer showing the amount of money the agent holds.
- Current Time: String formatted as hh:mm AM/PM.
- Day: Integer indicating the current day in the season.
- Season: String, one of spring, summer, fall, or winter.
- Item in Your Hand: A dictionary with fields:
 - index: Integer slot index.
 - currentitem: String name of the item.
- **Toolbar**: A list of 36 item slot descriptions in the format:

```
"slot_index N: [Item Name] (quantity: Q)" or "slot_index N: No item"
```

- Current Menu: A dictionary with keys such as: type, message, shopmenudata, animalsmenudata, etc.
- **Surrounding Blocks**: A list of nearby tiles, each with:
 - position: A 2D integer coordinate offset relative to the agent.
 - object: A list of string attributes (e.g., Type: Dirt, Diggable: True).
 - (Optional) npc on this tile: Information about an NPC, if present.



Figure 6: Example of game screenshot as part of the observation space.

C.2 ACTION SPACE

To better align with the actual action space of human players while simplifying redundant operations that do not contribute to the model's decision-making capabilities, we designed a simplified minimal action space.

The action space defines the set of skills (or actions) that an agent can perform. Each action is implemented as a function with a specific call template and a thorough comment. The full list of actions is provided in **Table 4**, which details each available action along with its parameters and intended behavior. In our experiments, we excluded the *navigate* action from the available action space. While *navigate* provides a high-level shortcut for moving between maps, our primary objective is to evaluate the agent's realistic exploration abilities and its performance using an action space that more closely mirrors human interactions.

864 865

Table 4: Complete Action Space with Call Templates and Parameter Descriptions

866 Action Description 867 Call Template: move (x = ..., y = ...)move(x, y)868 **Parameters:** x, y - X and Y coordinates of the destination. Move to the position (x, y). 870 craft (item) Call Template: craft (item = ...) 871 Parameters: item - The name of the item to craft. 872 Craft an item based on its name. 873 874 use (direction) **Call Template:** use (direction = ...) Parameters: direction - A string: up, right, down, or left. 875 Use an item in a specified direction. Requires proper positioning. 876 877 choose_item(slot_index) Call Template: choose_item(slot_index = ...) **Parameters:** slot_index - Inventory index (0-35). 878 Choose the item in the specified inventory slot. 879 880 interact (direction) Call Template: interact (direction = ...) 881 **Parameters:** direction – A string: up, right, down, or left. Interact with an object or NPC in a specific direction. 882 883 choose_option(Call **Template:** choose_option(option_index = ..., 884 option_index, quantity = ..., direction = ...) quantity, direction) **Parameters:** 885 option_index - Index of the option to choose. 886 quantity (optional) - Quantity of items to buy/sell. 887 direction (optional) - "in" for buy/take, "out" for sell/put. Choose from a list of options, with optional quantity and direction. attach_item(slot_index) Call Template: attach_item(slot_index = ...) 290 **Parameters:** slot_index - Inventory index (0-35). 891 Attach the item in the given inventory slot. 892 Call Template: unattach_item() unattach_item() 893 Parameters: None. Unattach the currently attached item. 895 Call Template: menu (option = ..., menu_name = ...) menu (option, 896 menu_name) **Parameters:** 897 option - "open" or "close" 898 menu_name - Menu name (e.g., "map") 899 Open or close a specific menu. 900 navigate (name) Call Template: navigate (name = ...) 901 Parameters: name - Name of the location to navigate to. 902 Navigate to a specified location. 903 904

C.3 TASK

905

906 907

908 909

910

911

912

913

914

915

916

917

Our benchmark provides 1000 tasks organized into five categories:

• Farming: Farming tasks can be broadly categorized into two types: cultivation (growing crops) and husbandry (raising animals). Easy tasks involve routine agricultural work such as clearing and tilling tiles, fertilizing, sowing seeds, watering plants, harvesting mature crops and animal products, as well as animal care—including feeding, grazing, and interaction. These simple, discrete operations test whether agents possess the most fundamental production capabilities. Medium tasks require agents to independently procure farming resources such as seeds, fertilizer, water, and hay. These resources can be obtained through foraging, crafting, or purchasing. Hard tasks typically span multiple days, requiring agents to perform daily routine farming activities. These activities form a cohesive production chain, following a specific sequence where each step is interdependent. For example, if a task involves growing a plant from seed to harvest, the agent must plant the seed in tilled dirt, water it daily over several days, and finally reap the crop. Any missed step, such as failing to water the plant on a given day, could delay or even prevent maturation. Thus, hard farming

tasks demand that agents autonomously allocate time and resources efficiently, presenting significant tests of multi-step reasoning and long-term planning capabilities.

- Crafting: Crafting tasks encompass both fabrication and cooking. Most crafting activities simply require adequate raw materials, though some may need equipment like furnace or cookout kit. Typically, easy tasks provide all necessary materials and tools directly in the agent's inventory, requiring only proper execution of crafting procedures. Occasionally, agents might need to gather readily available resources like a few woods or stones. Medium tasks demand greater autonomy, which agents must identify, locate and acquire appropriate materials and equipment through careful planning. Hard tasks involve procuring diverse materials through demanding methods (like deep mining for ores), followed by complex, multi-stage processes like crafting intermediary components, testing an agent's comprehensive crafting capabilities.
- Exploration: Exploration tasks can be categorized into three types: map navigation/pathfinding, resource gathering in wilderness areas, and completing challenging in-game quests. Easy tasks may involve traveling to an accessible location or collecting specified items within a small nearby area. Medium tasks present greater complexity, some require venturing into more distant and hazardous environments (like the second floor of the mines), while others demand searching expansive zones (such as an entire forest) for randomly spawning resources. Hard tasks challenge agents to locate extremely rare resources with highly randomized spawn locations (like amethysts in mines). Additionally, built-in game quests that chain multiple sub-tasks of varying types, challenging enough to occasionally defeat even experienced human players also qualify as high-difficulty exploration objectives.
- Combat: Agents are tasked with eliminating a specified number of monsters in mines. As difficulty escalates, targets become both more formidable and numerous. Basic adversaries like slimes, bugs, and grubs present minimal threat, their predictable movement patterns, low health pools, and weak attacks make them easy to dispatch or evade. However, advanced creatures employ deadly specialties: flies attack with erratic, lightning-fast strikes; duggies ambush from subterranean positions; rock crabs retreat into impregnable armored stances. Defeating these requires dynamic positioning, tactical strike timing, and adaptive combat strategies.
- Social: Social tasks encompass two primary objectives: cultivating relationships with NPCs, and conducting transactional interactions with specialized NPCs (such as carpenter, blacksmith, etc.). Easy tasks require only basic interactions like conversations or gift-giving, and agents are teleported directly to designated transaction locations (e.g., at the counter of Pierre's General Store). Medium tasks demand agents autonomously select and navigate to appropriate venues. Hard tasks challenge agents to build high friendship levels with NPCs. This requires strategic planning to increase rapport through daily greetings, thoughtful gift-giving, and fulfilling requests. Each NPC possesses unique behavioral patterns and preferences. Agents must develop customized engagement strategies, as actions that delight one NPC (e.g., a favored gift) may offend another (e.g., a disliked item). This nuanced system tests agents' adaptive social intelligence.

Each category of tasks is encoded in a YAML file. This dictionary-like file format is both highly readable and convenient for processing by Python programs. The structure of a task is as follows:

Task Format

sow_5_dirt_with_cauliflower_seeds: The key serves as both the name and the description of the task, and will be used as the prompt input to LLM.

- id: A unique identifier for the task within its category.
- object: The target object of the task—such as item, location, character, or quest—that the player
 tries to acquire, reach, interact with, or complete in a specific quantity.
- quantity: The required number of the target object. For non-quantifiable objects, such as location and NPC, quantity is simply set to 1.
- tool: The tool required to complete the task.
- save: The initial game save for the task, which has pre-configured some common environmental settings to reduce frequent calls to simulator APIs.
- init_commands: This is a list of commands to invoke the simulator APIs. After loading the task, StarDojo will automatically execute these commands one by one, invoking the corresponding API to fine-tune the initial conditions. Combined with the save file, StarDojo achieves efficient standardization for each task.
- evaluator: The type of evaluator to assess the task.
- difficulty: The difficulty level of the task.

The 100 tasks in StarDojo-Lite are listed in Table 5.

Table 5: Complete list of tasks in StarDojo-Lite.

Task	Category ID		Quantity	Tool	Save	Init Commands	Evaluator	Difficult
clear_10_weeds_with_scythe	Farming 0	Weeds	10	Scythe	save_new		clear	easy
clear_5_stone_with_pickaxe clear_30_debris_with_scythe_and_pickaxe_and_axe	Farming 1 Farming 2	Stone Debris	5 30	Pickaxe Scythe, Pickaxe, Ax	save_new e save_new		clear clear	easy medium
till_5_tile_with_hoe	Farming 3	Tile	5	Hoe	save_new		till	easy
fertilize_5_dirt_with_basic_retaining_soil fertilize_1_dirt_with_speed_gro	Farming 4 Farming 5	Dirt Dirt	5 1	Basic Retaining Soi Speed-Gro	I save_farming save_farming	add_item_by_name("Basic Retaining Soil", 5)	fertilize fertilize	easy
sow_5_dirt_with_cauliflower_seeds	Farming 6	Dirt	5	Cauliflower Seeds	save_farming	add_item_by_name("Cauliflower Seeds", 5)	sow	easy
sow_1_dirt_with_potato_seeds	Farming 7 Farming 8	Dirt	1 5	Potato Seeds	save_new	set_time(time=900)	sow water	medium
water_5_crop_with_watering_can harvest_5_parsnip	Farming 8 Farming 9	Crop Parsnip	5	Watering Can	save_farming save_farming		harvest	easy
cultivate_and_harvest_1_garlic	Farming 10	Garlic	1		save_new	add_item_by_name("Garlic Seeds")	harvest	hard
pet_3_animal pet_8_animal	Farming 11 Farming 12		3 8		save_farming save_farming		pet pet	easy medium
open_1_deluxe_coop	Farming 12		1		save_farming		open	easy
fill_1_pet_bowl_with_watering_can	Farming 14	Pet Bowl	1	Watering Can	save_new		fill	easy
fill_1_feeding_bench_with_hay	Farming 15		1	Hay	save_farming save_farming	add_item_by_name("Hay", 12)	fill harvest	easy
harvest_1_egg harvest_1_milk_with_milk_pail	Farming 16 Farming 17		1	Milk Pail	save_farming	add_item_by_name("Milk Pail")	harvest	easy
harvest_3_milk_with_milk_pail	Farming 18	Milk	3	Milk Pail	save_farming		harvest	hard
incubate_1_chicken_with_incubator earn_50_friendship_with_1_cat	Farming 19 Farming 20		1 50	Incubator	save_farming save_new		incubate friendship	hard hard
craft_1_cherry_bomb	Crafting 0	Cherry Bomb	1		save_new	set_time(time=900)	craft	medium
craft_1_wood_fence	Crafting 1	Wood Fence	1		save_new		craft	easy
craft_1_sprinkler	Crafting 2	Sprinkler	1		save_new	add_item_by_name("Furnace"), add_item_by_name("Copper Ore", 5), add_item_by_name("Iron Ore", 5), add_item_by_name("Coal", 2)	craft	medium
craft_1_basic_retaining_soil		Basic Retaining Soi	1 1		save_new		craft	easy
craft_1_spring_seeds	Crafting 4	Spring Seeds	1		save_new	add item by second Account add item by second ("Monle Cond")	craft	hard
craft_1_field_snack	Crafting 5	Field Snack	1		save_new	add_item_by_name("Acorn"), add_item_by_name("Maple Seed"), add_item_by_name("Pine Cone")	craft	easy
craft_1_torch	Crafting 6	Torch	1		save_new		craft	easy
craft_1_furnace craft_1_chest	Crafting 7 Crafting 8	Furnace Chest	1		save_new save_new	set_time(time=900) set_time(time=900)	craft craft	medium medium
craft_1_criest craft_1_scarecrow	Crafting 9	Scarecrow	1		save_new	add_item_by_name("Wood", 50), add_item_by_name("Coal"),	craft	easy
						add_item_by_name("Fiber", 20) add_item_by_name("Furnace"), add_item_by_name("Copper Ore", 5),		
produce_1_copper_bar_with_furnace	Crafting 10		1	Furnace	save_new	add_item_by_name("Coal")	craft	easy
produce_1_refined_quartz_with_furnace cook_1_fried_egg_with_stove	Crafting 11 Crafting 12		1	Furnace Stove	save_new save_new	upgrade_house(1), add_item_by_name("Egg")	craft craft	hard easy
cook_1_nied_egg_with_stove cook_1_salad_with_cookout_kit	Crafting 12		1	Cookout Kit	save_new	upgrade_nouse(1), add_nem_by_name(Egg)	craft	hard
go_to_bed	Exploration 0	Bed	1		save_new		sleep	easy
go_to_coop	Exploration 1	Coop	1		save_new		location	easy
	Exploration 2	BusStop	1		save_new		location	easy
go_to_backwoods go_to_pierre's_general_store	Exploration 3 Exploration 4	Backwoods SeedShop	1		save_new save_new	set_time(time=900)	location location	easy
	Exploration 5	AnimalShop	1		save_new	set_time(time=900)	location	easy
	Exploration 6	FishShop	1		save_new	set_time(time=900)	location	easy
go_to_carpenter's_shop go_to_the_mines_2nd_floor	Exploration 7 Exploration 8	ScienceHouse UndergroundMine2	1		save_new save_new	set_time(time=900)	location location	easy medium
		UndergroundMine5		Elevator	save_new		location	medium
go_to_the_mines_10th_floor		UndergroundMine1			save_new		location	hard
	Exploration 11 Exploration 12		10 20	Axe Axe	save_new save_new		harvest	easy medium
		Wild Horseradish	1		save_new	warp("forest")	harvest	medium
	Exploration 14		1		save_new	warp("town")	harvest	medium
	Exploration 15 Exploration 16		1		save_new save_new	warp("mountain") warp("beach")	harvest	medium easy
	Exploration 17		10	Scythe	save_new	waip(beach)	silo	easy
forage_1_quartz	Exploration 18		1		save_new	warp_mine(1)	harvest	medium
	Exploration 19 Exploration 20		1	Hoe Pickaxe	save_new save_new	warp_mine(13) warp_mine(1)	harvest harvest	easy hard
mine_1_copper_ore_with_pickaxe	Exploration 21		1	Pickaxe	save_new	warp_mine(2)	harvest	easy
	Exploration 22		1	Pickaxe	save_new	warp_mine(1)	harvest	medium
quit_1_quest take_1_quest_reward	Exploration 23 Exploration 24		1		save_new save_quests		quit reward	easy
complete_1_help_wanted_quest		Help Wanted Quest			save_new		complete_help	hard
	Exploration 26		1		save_new		complete_story	hard
1 7-1 -0 0-	Exploration 27		1		save_quests		complete_story	hard
kill_1_green_slime_with_rusty_sword kill_5_green_slime_with_rusty_sword	Combat 0 Combat 1	Green Slime Green Slime	1 5	Rusty Sword Rusty Sword	save_new save_new	warp_mine(2) warp_mine(2)	kill kill	easy medium
kill_10_green_slime_with_rusty_sword kill_10_green_slime_with_rusty_sword	Combat 1	Green Slime	10	Rusty Sword Rusty Sword	save_new save_new	warp_mine(2) warp_mine(2)	kill	hard
kill_1_bug_with_rusty_sword	Combat 3	Bug	1	Rusty Sword	save_new	warp_mine(2)	kill	easy
kill_5_bug_with_rusty_sword kill_10_bug_with_rusty_sword	Combat 4 Combat 5	Bug Bug	5 10	Rusty Sword Rusty Sword	save_new save_new	warp_mine(2) warp_mine(2)	kill kill	medium hard
kill_1_fly_with_rusty_sword	Combat 6	Fly	1	Rusty Sword	save_new	warp_mine(2)	kill	medium
kill_1_duggy_with_rusty_sword	Combat 7	Duggy	1	Rusty Sword	save_new	warp_mine(6)	kill	medium
kill_1_grub_with_rusty_sword kill_5_grub_with_rusty_sword	Combat 8 Combat 9	Grub Grub	1 5	Rusty Sword Rusty Sword	save_new save_new	warp_mine(15) warp_mine(15)	kill kill	easy medium
			10	Rusty Sword	save_new	warp_mine(15)	kill	hard
kill_10_grub_with_rusty_sword	Combat 10	Grub			save_new	warp_mine(2)	kill	medium
kill_1_rock_crab_with_rusty_sword	Combat 10 Combat 11		1	Rusty Sword				easy
kill_1_rock_crab_with_rusty_sword ship_1_parsnip_with_shipping_bin	Combat 11 Social 0	Rock Crab Parsnip	1	Shipping Bin	save_new	add_item_by_name("Parsnip")	sell	
kill_1_rock_crab_with_rusty_sword ship_1_parsnip_with_shipping_bin purchase_5_beer	Combat 11 Social 0 Social 1	Rock Crab Parsnip Beer	1 1 5		save_new save_new	warp_shop("gus")	purchase	easy
kill_1_rock_crab_with_rusty_sword ship_1_parsnip_with_shipping_bin	Combat 11 Social 0	Rock Crab Parsnip	1		save_new			easy medium easy
kill_l_rock_emb_with_rusty_sword ship_l_parsnip_with_shipping_bin purchase_5_beer purchase_l_muscle_remedy sell_5_parsnip_to_pierre sell_l_parsnip_to_pierre	Combat11Social0Social1Social2Social3Social4	Parsnip Beer Muscle Remedy Parsnip Parsnip	1 5 1 5		save_new save_new save_new save_new save_new	<pre>warp_shop("gus") set_time(time=900) warp_shop("pierre"), add_item_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_time(time=900)</pre>	purchase purchase sell sell	medium easy medium
kill_l_rock_crab_with_rusty_sword ship_l_parsnip_with_shipping_bin purchase_l_person purchase_l_nuscle_remedy sell_5_parsnip_to_pierre sell_l_parsnip_to_pierre upgrade_to_copper_pickaxe	Combat 11 Social 0 Social 1 Social 2 Social 3 Social 4 Social 5	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Copper Pickaxe	1 5 1 5 1 1		save_new save_new save_new save_new save_new save_new	warp_shop("gus") set_time(time=900) warp_shop("pierre"), add_tiem_by_name("Parsnip", 5) add_tiem_by_name("Parsnip"), set_time(time=900) ad_tiem_by_name("Copper Bar", 5)	purchase purchase sell sell upgrade_tool	medium easy medium hard
kill_l_rock_crab_with_rusty_sword ship_l_parsinp_with_shipping_bin purchase_5_beer purchase_1_muscle_remedy sell_5_pansinp_to_pierre sell_1_parsinp_to_pierre sell_1_parsinp_to_pierre upgrade_to_coopper_pickaxe break_5_goodd	Combat11Social0Social1Social2Social3Social4	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Copper Pickaxe Geode	1 5 1 5		save_new save_new save_new save_new save_new save_new save_new	warp_shop("gus") set_ime(time=900) warp_shop("pierre"), add_item_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_time(time=900) add_item_by_name("Copper Bar", 5) warp_shop("clint"), add_item_by_name("Goode", 5)	purchase purchase sell sell upgrade_tool break	medium easy medium
kill_l_rock_crab_with_rusty_sword ship_l_parsinp_with_shipping_bin purchase_5_beer purchase_1_muscle_remedy sell_5_parsing_bo_pierre sell_1_parsing_bo_pierre sell_1_parsing_bo_pierre upgrade_to_copper_pickaxe break_5_gende_purchase_joig_membership purchase_miscarts_development_project	Combat 11 Social 0 Social 1 Social 2 Social 3 Social 4 Social 5 Social 6 Social 7 Social 8	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Copper Pickaxe Geode Joja Membership Minecarts Repaired	1 5 1 5 1 1 5 1 1 5		save_new save_new save_new save_new save_new save_new save_new save_new save_new	warp_shop("gus") set_ime(time=900) warp_shop("pierre"), add_item_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_ime(time=900) add_item_by_name("Copper Bar", 5) warp_shop("clint"), add_item_by_name("Geode", 5) warp_shop("clint"), add_item_by_name("Geode", 5) warp_shop("clint"), indictim_by_name("Geode", 5) joja_membership(), set_ime(time=900)	purchase purchase sell sell upgrade_tool break jojamart jojamart	medium easy medium hard easy easy medium
kill_l_rock_crab_with_rusty_sword ship_l_parsnip_with_shipping_bin purchase_5_beer purchase_1_muscle_remedy sell_5_pansnip_to_pierre sell_1_pansnip_to_pierre sell_1_pansnip_to_pierre upgrade_to_copper_pickaxe break_5_goode purchase_joja_membership purchase_inja_membership purchase_inja_membership purchase_inja_membership	Combat 11	Rock Crab Parsnip Beer Muscle Remedy Parsnip Copper Pickaxe Geode Joja Membership Minecarts Repaired Large Pack	1 5 1 5 1 1 5 1 1 5		save_new	warp_shop("gus") set_time(time=900) warp_shop("piers"), add_item_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_time(time=900) add_item_by_name("Copper Bar", 5) warp_shop("clint"), add_item_by_name("Goods", 5) warp_shop("clint"), add_item_by_name("Goods", 5) warp("oja", 21, 26) joja_nembership(), set_time(time=900) warp_shop("piers")	purchase purchase sell sell upgrade_tool break jojamart jojamart backpack	medium easy medium hard easy easy medium easy
kill_Irock_crab_with_nusty_sword ship_I_pursnip_with_shipping_bin purchase_5_beer purchase_1_muscle_remedy sell_5_pursnip_to_pierre sell_1_pursnip_to_pierre sell_1_pursnip_to_pierre upgrade_to_copper_pickaxe break_5_geode purchase_joi_membership purchase_minecarts_development_project upgrade_to_lurge_puck	Combat 11	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Copper Pickaxe Geode Joja Membership Minecarts Repaired Large Pack Chicken	1 5 1 5 1 1 5 1 1 5		save_new	warp_shop("gus") set_time(time=900) warp_shop("piere"), add_item_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_time(time=900) add_item_by_name("Copper Bar", 5) warp_shop("clint"), add_item_by_name("Geode", 5) warp_shop("clint"), add_item_by_name("Geode", 5) warp("sigs", 21, 26) joja_nembership(), set_time(time=900) warp_shop("pierre") warp_shop("marnie")	purchase purchase sell sell upgrade_tool break jojamart jojamart backpack purchase_animal	medium easy medium hard easy easy medium easy easy
kill_l_nock_crab_with_nusty_sword ship_l_parsnip_with_shipping_bin purchase_5_beer purchase_1_muscle_remedy sell_5_parsnip_to_pierre sell_1_parsnip_to_pierre upgrade_to_copper_pickaxe break_5_geode purchase_joja_membership purchase_minearts_development_project upgrade_to_large_pack purchase_lichicken sell_1_chicken	Social O Social 1 Social 2 Social 3 Social 4 Social 5 Social 6 Social 7 Social 8 Social 9 Social 10 Social 11 Social 10 Social 11 Social 1	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Copper Pickaxe Geode Joja Membership Minecarts Repaired Large Pack Chicken Chicken	1 5 1 5 1 1 5 1 1 5 1		save_new	warp_shop("gus") set_time(time=900) warp_shop("pierre"), add_tiem_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_time(time=900) add_tiem_by_name("Copper Bar", 5) warp_shop("clint"), add_item_by_name("Geode", 5) warp_shop("clint"), add_item_by_name("Geode", 5) warp_shop("pierre") warp_shop("pierre") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie")	purchase purchase sell sell upgrade_tool break jojamart jojamart backpack purchase_animal sell_animal	medium easy medium hard easy easy medium easy easy easy
kill_l_rock_erab_with_nusty_sword ship_l_parsnip_with_shipping_bin purchase_5_beer purchase_1_muscle_remedy sell_5_parsnip_to_pierre sell_1_parsnip_to_pierre sell_1_parsnip_to_pierre upgrade_to_copper_pickaxe break_5_geode purchase_joja_membership purchase_minecarts_development_project upgrade_to_large_pack purchase_l_chicken sell_t_chicken build_l_big_coop	Social O Social 1 Social 2 Social 3 Social 4 Social 5 Social 6 Social 7 Social 8 Social 9 Social 10 Social 11 Social 12 Social 13 Social 14 Social 14 Social 15 Social 1	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Copper Pickaxe Geode Joja Membership Minecarts Repaired Large Pack Chicken Chicken Big Coop	1 5 1 5 1 1 5 1 1 1 1 1 1 1 1 1		save_new	warp_shop("gus") set_time(time=900) warp_shop("piers"), add_item_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_time(time=900) add_item_by_name("Copper Bar", 5) warp_shop("clint"), add_item_by_name("Goods", 5) warp_shop("clint"), add_item_by_name("Goods", 5) warp_shop("piers") warp_shop("piers") warp_shop("robin"), add_item_by_name("Woods", 400), add_item_by_name("Woods", 400), add_item_by_name("Stones", 150)	purchase purchase sell sell upgrade_tool break jojamart jojamart backpack purchase_animal sell_animal build	medium easy medium hard easy easy medium casy easy easy easy easy
kil_l_rock_erab_with_rusty_sword ship_l_parsnip_with_shipping_bin purchase_5_bee purchase_l_mesele_remedy sell_5_pansnip_to_pierre sell_1_pansnip_to_pierre sell_1_pansnip_to_pierre sell_1_pansnip_to_pierre upgrade_to_copper_pickaxe break_5_geode purchase_joia_membership purchase_incearts_development_project upgrade_to_large_pack purchase_l_chicken sell_1_chicken build_l_big_coop move_l_coop upgrade_fammouse	Social O Social 1 Social 2 Social 3 Social 4 Social 5 Social 6 Social 7 Social 8 Social 9 Social 10 Social 11 Social 10 Social 11 Social 1	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Copper Pickave Geode Joja Membership Minecarts Repaired Large Pack Chicken Chicken Big Coop Coop	1 1 5 1 5 1 1 5 1 1 1 1 1 1 1 1 1		save_new	warp_shop("gus") set_time(time=900) warp_shop("pierre"), add_tiem_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_time(time=900) add_item_by_name("Copper Bar", 5) warp_shop("clint"), add_item_by_name("Goode", 5) warp_shop("clint"), add_item_by_name("Goode", 5) warp_shop("pierre") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") sad_item_by_name("Wood", 400), add_item_by_name("Stone", 150) set_time(time=900)	purchase purchase sell sell upgrade_tool break jojamart jojamart backpack purchase_animal sell_animal	medium easy medium hard easy easy medium easy easy easy easy easy easy medium
killrock_crab_with_nusty_sword ship_l_parsnip_with_shipping_bin purchase_5_beer purchase_1_muscle_remedy sell_5_parsnip_to_pierre sell_l_parsnip_to_pierre sell_l_parsnip_to_pierre upgrade_to_coppet_pickase break_5_geode purchase_ipia_membership purchase_ininecarts_development_project upgrade_to_large_pack purchase_l_chicken sell_l_chicken sell_l_chicken build_l_big_coop move_l_coop upgrade_farmhouse demoish_l_shipping_bin	Social Oscial Social ISOcial ISOCIA ISOCIA	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Copper Pickaxe Geode Joja Membership Minecarts Repaired Large Pack Chicken Chicken Big Coop Coop Farmhouse Shipping Bin	1		save_new	warp_shop("gus") set_ime(time=900) warp_shop("pierre"), add_item_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_ime(time=900) add_item_by_name("Copere Bar", 5)) warp_shop("clint"), add_item_by_name("Goode", 5) warp_shop("clint"), add_item_by_name("Goode", 5) warp_shop("robin"), set_time(time=900) warp_shop("marnie") warp_shop("marnie") warp_shop("robin"), add_item_by_name("Wood", 400), add_item_by_name("Stone", 150) set_time(time=900) warp_shop("robin"), add_item_by_name("Wood", 450) set_time(time=900)	purchase purchase sell sell upgrade_tool break jojamart jojamart backpack purchase_animal sell_animal build move upgrade_farmhous demolish	medium easy medium hard easy easy medium easy easy easy easy easy medium e casy medium
killnockcrabwithnusty_sword ship L_parsnipwith_shipping_bin purchase Sheer _ purchase Imuscleremedy sell S_parsniptopierre sell I_parsniptopierre sell I_parsniptopierre upgradetocopperpickaxe break Sgoode purchaseiojamembership purchaseiojamembership purchasel_picken sell Ichicken sell Ichicken build Ibigcoop move Icoop upgradefarmhouse demolishshippingbin talk toletx	Combat	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Parsnip Copper Pickaxe Geode Joja Membership Mincearts Repaired Large Pack Chicken Chicken Chicken Big Coop Farmhouse Shipping Bin Alex	1		save_new	warp_shop("gas") set_time(time=900) warp_shop("piere"), add_tiem_by_name("Parsnip", 5) add_tiem_by_name("Parsnip"), set_time(time=900) add_tiem_by_name("Suspir"), set_time(time=900) add_tiem_by_name("Copper Bas", 5) warp_shop(clint"), add_tiem_by_name("Goode", 5) warp_shop(clint"), add_tiem_by_name("Goode", 5) warp_shop("piere") warp_shop("robin"), add_tiem_by_name("Wood", 400), add_tiem_by_name("Stone", 150) set_time(time=900) warp_shop("robin"), add_tiem_by_name("Wood", 450) set_time(time=900) set_time(time=900) set_time(time=900) set_time(time=900)	purchase purchase sell sell upgrade_tool break jojamart backpack purchase_animal sell_animal build move upgrade_farmhous demolish talk	medium easy medium hard easy easy medium easy easy easy easy easy medium e easy easy easy easy medium exity medium exity medium exity medium exity
killrock_crab_with_nusty_sword ship_l_parsnip_with_shipping_bin purchase_5_beer purchase_1_muscle_remedy sell_5_parsnip_to_pierre sell_l_parsnip_to_pierre sell_l_parsnip_to_pierre upgrade_to_coppet_pickase break_5_geode purchase_ipia_membership purchase_ininecarts_development_project upgrade_to_large_pack purchase_l_chicken sell_l_chicken sell_l_chicken build_l_big_coop move_l_coop upgrade_farmhouse demoish_l_shipping_bin	Social Oscial Social ISOcial ISOCIA ISOCIA	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Parsnip Lopper Pickase Geode Joja Membership Minecarts Repaired Large Pack Chicken Chicken Chicken Chicken Sig Coop Coop Farmhouse Shipping Bin Alex Sebastian	1		save_new	warp_shop("gus") set_ime(time=900) warp_shop("pierre"), add_item_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_ime(time=900) add_item_by_name("Copere Bar", 5)) warp_shop("clint"), add_item_by_name("Goode", 5) warp_shop("clint"), add_item_by_name("Goode", 5) warp_shop("robin"), set_time(time=900) warp_shop("marnie") warp_shop("marnie") warp_shop("robin"), add_item_by_name("Wood", 400), add_item_by_name("Stone", 150) set_time(time=900) warp_shop("robin"), add_item_by_name("Wood", 450) set_time(time=900)	purchase purchase sell sell upgrade_tool break jojamart jojamart backpack purchase_animal sell_animal build move upgrade_farmhous demolish	medium easy medium hard easy easy medium easy easy easy easy easy medium e casy medium
kill_l_rock_crab_with_musty_sword ship_l_parsnip_with_shipping_bin purchase_5_beer purchase_1_muscle_remedy sell_5_parsnip_to_pierre sell_l_parsnip_to_pierre sell_l_parsnip_to_pierre upgrade_1o_copper_pickaxe break_5_geode purchase_pioj_membership purchase_minecarts_development_project upgrade_1o_lapse_pack purchase_ling_membership purchase_ling_membership purchase_ling_membership purchase_ling_membership purchase_ling_membership purchase_l_chicken sell_l_chicken sell_l_chicken sell_l_chicken build_l_big_coop move_l_coop upgrade_farmhouse demoish_l_shipping_bin talk_to_alex talk_to_sebastian talk_to_vincent give_abigail_l_amethyst	Combat	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Parsnip Copper Pickase Geode Joja Membership Minecarts Repaired Large Pack Chicken Chicken Chicken Chicken Slig Coop Coop Farmhouse Shipping Bin Alex Sebastian Vincent Abigail	1	Shipping Bin	Save_new	warp_shop("gus") set_time(time=900) warp_shop("pierre"), add_item_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_time(time=900) add_item_by_name("Parsnip"), set_time(time=900) add_item_by_name("parsnip", set_time(time=900) warp_shop("clint"), add_item_by_name("Goode", 5) warp_shop("pierre") warp_shop("pierre") warp_shop("pierre") warp_shop("robin"), add_item_by_name("Wood", 400), add_item_by_name("Stone", 150) set_time(time=900) warp_shop("robin"), add_item_by_name("Wood", 450) set_time(time=900) set_time(time=500) set_time(time=500) set_time(time=500) add_item_by_name("None", 150) add_item_by_name("None", 150) set_time(time=500)	purchase purchase sell sell upgrade_tool break jojamart jojamart backpack purchase_animal sell_animal build move pupgrade_farmhous demolish talk talk talk tgift	medium easy medium hard easy easy medium easy easy easy easy easy medium easy easy medium esy easy
killrock_crab_with_nusty_sword ship_l_parsnip_with_shipping_bin purchase_5_beer purchase_1_muscle_remedy sell_5_parsnip_to_pierre sell_1_parsnip_to_pierre sell_1_parsnip_to_pierre upgrade_to_copper_pickaxe break_5_geode purchase_joig_membership purchase_minecarts_development_project upgrade_to_large_pack purchase_l_ichicken sell_l_chicken sell_l_chicken build_l_big_coop move_l_coop upgrade_farmhouse demoisth_l_shipping_bin talk_to_slex talk_to_sebastain talk_to_vincent give_abigail_l_amethyst give_halve_l_coconut	Combat	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Parsnip Copper Pickave Geode Joja Membership Minecarts Repaired Large Pack Chicken Chicken Chicken Chicken Goop Coop Farmhouse Shipping Bin Alex Sebastian Vincent Abigail Haley	1 5 1 5 1 5 1 1 5 1 1 1 1 1 1 1 1 1 1 1	Shipping Bin Amethyst Coconut	Save_new	warp_shop("gus") set_ime(time=900) warp_shop("pierre"), add_item_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_ime(time=900) add_item_by_name("Copere Bar", 5) warp_shop("clint"), add_item_by_name("Goode", 5) warp_shop("clint"), add_item_by_name("Goode", 5) warp_shop("pierre") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") set_time(time=900) set_time(time=900) set_time(time=1500) set_time(time=1500) set_time(time=1500) add_item_by_name("Cooraut", 1,4), set_time(time=900) add_item_by_name("Cooraut", 1,4), set_time(time=1100) add_item_by_name("Cooraut", 1,4), set_time(time=1100)	purchase purchase sell sell upgrade_tool break jojamart jojamart jojamart backpack purchase_animal build move demolish talk talk talk gift gift	medium easy medium hard easy easy easy easy easy easy medium easy easy easy medium e e casy medium easy easy medium easy easy easy easy easy easy easy easy
kill_I_rock_crab_with_nasty_swerd ship_I_parsnip_with_shipping_bin purchase_5_beer purchase_1_muscle_remedy sell_5_parsnip_to_pierre sell_1_parsnip_to_pierre sell_1_parsnip_to_pierre upgrade_to_copper_pickaxe break_5_geode purchase_join_membership purchase_join_membership purchase_to_to_come_to	Combat	Rock Crab Parsnip Beer Muscle Remedy Parsnip Parsnip Parsnip Copper Pickase Geode Joja Membership Minecarts Repaired Large Pack Chicken Chicken Chicken Chicken Chicken Slig Coop Coop Farmhouse Shipping Bin Alex Sebastian Vincent Abigail Haley Jas	1	Shipping Bin	Save_new	warp_shop("gus") set_time(time=900) warp_shop("piers"), add_item_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_time(time=900) add_item_by_name("Copper Bar", 5) warp_shop("clint"), add_item_by_name("Goods", 5) warp_shop("clint"), add_item_by_name("Goods", 5) warp_shop("robin"), add_item_by_name("Woods", 400), add_item_by_name("Sions", 150) set_time(time=900) warp_shop("robin"), add_item_by_name("Woods", 450) set_time(time=900) warp_shop("robin"), add_item_by_name("Woods", 450) set_time(time=900) set_time(time=900) add_item_by_name("Amethysis", 1, 4), set_time(time=900) add_item_by_name("Coronut", 1, 4), set_time(time=900) add_item_by_name("Coronut", 1, 4), set_time(time=900) add_item_by_name("Coronut", 1, 4), set_time(time=900) add_item_by_name("Croconut", 1, 4), set_time(time=900)	purchase purchase sell sell sell upgrade_tool break jojamart jojamart backpack purchase_animal sell_animal build move upgrade_farmhous demolish talk talk talk talk gift gift gift	medium easy medium hard casy easy medium casy casy casy easy easy medium casy medium casy medium casy easy medium casy easy medium casy easy easy casy casy casy
kill_Irok_crab_with_nusty_sword ship_I_pursinj_with_shipping_bin purchase_5_beer purchase_1_muscle_remedy sell_5_pursinj_to_pierre sell_1_pursinj_to_pierre sell_1_pursinj_to_pierre upgrade_to_copper_pickaxe break_5_geode purchase_joi_membership purchase_micentre_development_project upgrade_to_lurge_puck purchase_I_ichicken sell_1_chicken sell_1_chicken build_1_big_coop move_I_coop upgrade_farmhouse demolish_1_shipping_bin nuk_to_alex tulk_to_sebastian uuk_to_vincent give_laley_I_coconut	Combat	Rock Craib Parsnip Beer Muscle Remedy Parsnip Parsnip Parsnip Copper Pickase Geode Joja Membership Minecarts Repaired Large Pack Chicken Chicken Chicken Chicken Geop Farmhouse Coop Farmhouse Shipping Bin Alex Schastian Vincent Abigail Haley Jas Jodi	1	Amethyst Coconut Fairy Rose	Save_new	warp_shop("gus") set_ime(time=900) warp_shop("pierre"), add_item_by_name("Parsnip", 5) add_item_by_name("Parsnip"), set_ime(time=900) add_item_by_name("Copere Bar", 5) warp_shop("clint"), add_item_by_name("Goode", 5) warp_shop("clint"), add_item_by_name("Goode", 5) warp_shop("pierre") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") warp_shop("marnie") set_time(time=900) set_time(time=900) set_time(time=1500) set_time(time=1500) set_time(time=1500) add_item_by_name("Cooraut", 1,4), set_time(time=900) add_item_by_name("Cooraut", 1,4), set_time(time=1100) add_item_by_name("Cooraut", 1,4), set_time(time=1100)	purchase purchase sell sell upgrade_tool break jojamart jojamart jojamart backpack purchase_animal build move demolish talk talk talk gift gift	medium easy medium hard easy easy easy easy easy easy medium easy easy easy medium e e casy medium easy easy medium easy easy easy easy easy easy easy easy

C.4 SIMULATOR APIS

1026

1027 1028

1029

1032

1033

1034

1035

1036

1037 1038

1039

1040 1041

1042

1043

1045

1047

1049

1051

1052

1053

1054

1055

1056

1057

1058

1059

1063

1066

1068

1069

1070

1071

1072

1073

1074 1075

1077

1078

To establish the initial state for each task, we develop a set of simulator APIs that allow fine-grained customization of the game environment. Many tasks require strict initial conditions and world settings, such as resources (e.g., sufficient crop seeds in inventory), weather conditions (e.g., a rainy day), and game progress (e.g., unlocking the mines). The simulator APIs can configure these task-specific conditions, endowing each task with a unique environment setup. This mechanism significantly expands the benchmark's diversity and flexibility, allowing varied and nuanced task designs that closely mirror the dynamic production-living settings in human society.

To efficiently standardize our tasks, we also create a series of tailored game save files, which have pre-configured some common environmental settings using simulator APIs. These files reduce the excessively frequent calls to simulator APIs in real time, improving the efficiency of task execution. At the beginning of each task, StarDojo automatically loads the corresponding saved game and invokes specific simulator APIs, ensuring all necessary prerequisites are appropriately configured.

The complete APIs are as follows:

Simulator APIs

· Player Settings

- set base health(amount: int): Set the health capacity of the player.
- set_health(amount: int): Set the current health of the player.
- set_base_energy(amount: int): Set the energy capacity of the player.
- set_energy(amount: int): Set the current energy of the player.
- set_inventory_size(size: int): Set the inventory size.
- clear_inventory(): Clear the inventory.
- set_money(amount: int): Set the amount of money that the player possesses.
- add_item_by_id(id: str, count: int, quality: int): Add the specific item of given count and quality (e.g., 0, 1) to inventory by ID.
- add_item_by_name(name: str, count: int, quality: int): Add the specific item of given count
 and quality to inventory by name.
- lookup(name: str): Look up and print the item ID by name.
- **current_position()**: Print the current position of the player.
- add_recipe(type: str, recipe: str): Teach the player the specific crafting / cooking recipe.
- set_max_luck(): Set player's luck to the maximum.
- print_luck(): Print player's luck.

Surrounding Settings

- world_clear(entity: str, location: str): Remove all entities of the given type (e.g., "crops", "trees") from a location.
- set_terrain(terrain: str, id: str, x: int, y: int): Set the terrain feature of the given tile.
- place_item(item: str, type: str, x: int, y: int): Place the specific item on the given tile.
- remove_item(x: int, y: int): Remove the item on the given tile.
- place_crop(crop: str, x: int, y: int): Place the specific crop on the given tile.
- grow_crop(day: int, x: int, y: int): Grow the crop on the given tile for a specific number of days.
- grow_tree(day: int, x: int, y: int): Grow the tree on the given tile for a specific number of days.
- build(type: str, force: bool, x: int, y: int): Build the specific building at the given coordinate.
- build_stable(x: int, y: int): Build a stable at the given coordinate.
- move_building(x_source: int, y_source: int, x_dest: int, y_dest: int): Move the building
 from the source coordinate to the destination coordinate.
- remove_building(x: int, y: int): Remove the building at the given coordinate.
- upgrade_house(level: int): Upgrade the farmhouse to the given level.

· Character Settings

- spawn_pet(type: str, breed: str, name: str, x: int, y: int): Spawn a pet of given type (e.g., "cat", "dog"), breed (e.g., "0", "1"), and name on a tile.
- spawn_animal(type: str, name: str): Spawn an animal of given type and name in the animal house.

1	080	
1	081	

- 1082
- 1083 1084
- 1086
- 1088 1090
- 1093
- 1094 1096
- 1097
- 1099 1100

- 1102 1103
- 1104 1105 1106
- 1107 1108 1109
- 1110 1111
- 1112 1113 1114
- 1116 1117

1115

1118 1119

1120 1121 1122

1124 1125

1123

1126 1127

1128 1129

1130

1131

1132 1133

- grow_animal(name: str): Set the specific animal in current location to day 1 of adulthood, unless already adult.
- animal friendship(name: str, friendship: int): Set the friendship with the specific animal.
- npc friendship(npc: str, friendship: int): Set the friendship with the specific NPC.
- all_npc_friendship(friendship: int): Set the friendship with all NPCs.
- dating(npc: str): Make the specific NPC be the player's boyfriend / girlfriend.

Location Settings

- warp(location: str, x: int, y: int): Warp the player to given location and position.
- warp_mine(level: int): Warp the player to given mine level.
- warp_volcano(level: int): Warp the player to given volcano level.
- warp home(): Warp the player back home.
- warp_shop(npc: str): Warp the player to the shop run by the given NPC.
- warp_character(npc: str, location: str, x: int, y: int): Warp the specific NPC to given location and position.

· World Settings

- set_date(year: int, season: str, day: int): Set the date.
- set_time(time: int): Set the current time.
- rain(): Set the weather to rainy.

Progression Settings

- set_deepest_mine_level(level: int): Set the deepest mine level reached by the player.
- set_monster_stats(monster: str, kills: int: Set the kill stats for the specific monster to the given value.
- print_monster_stats(monster: str): Print the kill stats for the specific monster.
- start_quest(id: str): Start the quest of given ID.
- start_help_quest(type: str): Start a random help wanted quest of given type.
- complete_quest(id: str): Complete the quest of given ID.
- joja_membership(): Give the player JojaMart membership.
- spawn_junimo_note(id: str): Spawn the junimo note of given ID in the Community Center.
- mark_bundle(id: str): Mark the completion of the specific bundle.
- **complete_room_bundles(id: str)**: Complete all bundles in the specific room of the Community Center.
- community_development(id: str): Complete the community development project of given ID.
- receive_mail(mail: str): Add the specific mail to mailbox.
- trigger_event(id: str): Trigger the event of given ID.
- seen_event(id: str, see_or_forget: bool): Mark the viewed / unviewed status of the specific event.
- **load_save(save: str)**: Load the game save of given name.

C.5 EVALUATION

Given the scale of our extensive task set, it is imperative to design an efficient and reusable evaluation mechanism that not only accurately monitors task progression but also delivers immediate reward feedback to agents. Therefore, we implement a unified evaluation system based on textual observation comparison. The evaluation workflow follows a consistent pattern across all tasks, as outlined in the following steps and Algorithm 1:

- Maintain the previous observation: Store the agent's prior textual observation to enable temporal compari-
- Capture the current observation: Acquire the latest textual observation.
- Compare two observations: Based on the task's evaluator type (e.g., harvest, sell) and target object (e.g., item, NPC), the system detects related game state changes, such as items in the inventory and surrounding tiles, to quantify task progress.
- Accumulate incremental progress: Accumulate incremental changes captured per step into a sum.
- Check completion criterion: Validate whether predefined success conditions (e.g., quantity thresholds, event triggers) are met.

Algorithm 1 Task Evaluation

1134

11501151

1152

1153

1154

1155

1156 1157

1158 1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

11691170

1171 1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

11821183

1184

1185

1186

1187

```
1135
          1: function EVALUATE(obs)
1136
                 if last\_obs is null then
          2:
1137
          3:
                     last\_obs \leftarrow obs
1138
          4:
                     return {completed : False, current_quantity : 0}
1139
          5:
1140
          6:
                 quantity\_change \leftarrow Compare(evaluator, object, obs, last\_obs)
1141
          7:
                 last \ obs \leftarrow obs
1142
          8:
                 current\_quantity \leftarrow current\_quantity + quantity\_change
          9:
1143
                 if current\_quantity \ge quantity then
         10:
                      completed \leftarrow True
1144
                 else
         11:
1145
         12:
                     completed \leftarrow False
1146
         13:
                 end if
1147
         14:
                 return {completed : completed, current_quantity : current_quantity}
1148
         15: end function
1149
```

The evaluation system relies on the comprehensive observation space, allowing for step-by-step tracking of progress toward the task goal. This incremental inspection approach avoids differences caused by varying initial conditions, making the proposed evaluation mechanism more generalizable. The evaluation output provides standardized metrics, including completion status and current progress. These metrics also serve as reward feedback and, along with observations, are passed to the agents to prompt them to adjust their behavior.

D EXPERIMENT SETTINGS

If not mentioned explicitly, all experiments are conducted under the following settings: Agents have access to both visual and textual observations for their decision-making process. Visual observations are provided at a resolution of 720p (1280×720 pixels). Textual observations include detailed information about current state. Specifically, there are information about the player itself, including health, energy, gold, chosen item, and inventory information. There are also global information, which are current time, current day, season, and current open menu. Finally, we provide 7×7 agent-centered surrounding information, containing details about each tile, ranging from the terrain, debris, buildings, object, exits, NPCs, and furniture information, to other tile properties encoded in the game. In addition to receiving observations from the current timestep, agents are also provided with action and visual information from the previous timestep. Incorporating previous timestep information enables agents to reflect on past states and facilitate consistency in decision-making. All the agents can output at most two skills as an action to be executed sequentially. After executing all the actions, the environment remains paused until the agent outputs the next action. All experiments are repeated three times to ensure reliability.

Prompt Used in Main Experiment

You are a helpful AI assistant integrated with 'Stardew Valley' on the PC, equipped to handle various tasks in the game. Your advanced capabilities enable you to process and interpret gameplay screenshots and other relevant information. By analyzing these inputs, you gain a comprehensive understanding of the current context and situation within the game. Utilizing this insight, you are tasked with identifying the most suitable in-game action to take next, given the current task. You control the game character and can execute actions from the available action set. Upon evaluating the provided information, your role is to articulate the precise action you would deploy, considering the game's present circumstances, and specify any necessary parameters for implementing that action.

Here is some helpful information to help you make the decision.

Your Current task is: <\$task_description\$>

Basic knowledge about the environment:

- 1. Hoe is used to till the soil, Watering Can is used to water the soil, Pickaxe is used to break rocks, Axe is used to chop trees, Scythe is used to harvest crops.
- 2. When you want to go through a door, move in front of it by 1 tile, and interact towards it.
- 3. Please go to bed at night (after 18:00) even if your task is not yet complete!

1241

4. Call interact(direction) with a box, a shipping bin or anything else. Call use(direction) to use 1189 an item or tool in your inventory. 1190 Health: <\$health\$> 1191 1192 Energy: <\$energy\$> 1193 1194 Money: <\$money\$> 1195 1196 Current Time: <\$time\$> 1197 Day: <\$day\$> 1198 1199 Season: <\$season\$> 1200 Item in your hand: <\$chosen_item\$> 1201 1202 Toolbar of items which you can choose from: <\$toolbar_information\$> 1203 1204 Current menu: <\$current_menu\$> 1205 Surrounding blocks (Objects surrounds you): <\$surroundings\$> 1207 Valid action set in Python format to select the next action: <\$skill_library\$> 1208 1209 Last executed action: <\$pre action\$> 1210 1211 <\$image_introduction\$> 1212 Based on the above information, analyze the current situation and provide the reasoning for what you 1213 should do for the next step to complete the task. Then, you should output the exact action you want to 1214 execute in the game. You should respond to me with: 1215 Reasoning: You should think step by step and provide detailed reasoning to determine the next action 1216 executed on the current state of the task. You need to answer the following questions step by step. You 1217 cannot miss the last question: 1218 1. Is there an open menu? What is the current menu saying? What are the options? Which one 1219 should you choose or should you exit the menu? Use choose_option to make a choice. 1220 2. What is the current map? Do you need to move to another map? 1221 3. Refer to surroundings, what are the important tiles? What are their positions? 1222 4. You are always at (0, 0). You can only affect points adjacent to you, such as (0,1), (0,-1), (1,0), (-1,0). Is your target at these positions? If not, move next to it first. 1225 5. Analyze the information in the toolbar. Does it contain all the necessary items for completing the task? What is the current item? 1226 6. When calling use or interact, you need to decide the direction. For example, if you are at 1228 - call interact("up") or use("up") to interact with or use against (0,-1) - call interact("right") or use("right") to interact with or use against (1,0) 1230 - call interact("down") or use("down") to interact with or use against (0,1) 1231 - call interact("left") or use("left") to interact with or use against (-1,0) 1232 7. What is the current image showing? What additional information can be learned from the 1233 image? 1234 8. Do all the selected actions exist in the valid action set? If no, regenerate the actions and give 1235 the reasons. 1236 Actions: The requirements that the generated action needs to follow. The best action, or short sequence 1237 of actions without gaps, to execute next to progress in achieving the goal. Pay attention to the names of the available skills and to the previous skills already executed, if any. You should also pay more attention to the following action rules: 1239 1. You should output actions in Python code format and specify any necessary parameters to 1240

execute that action. If the function has parameters, you should also include their names and

decide their values, like move(x=0, y=1). If it does not have a parameter, just output the action, like unattch_item().
 You can only output at most 2 actions in the output.
 If you want to interact with the objects in the toolbar, you need to make sure that the target object is already selected. You need to use choose_item() to select them before executing use().
 If you want to plant a seed or harvest a mature crop, please use interact(). If you want to use tools, like axe, hoe, watering can, pickaxe and scythe, please use use().
 Your action should strictly follow the analysis in the reasoning. Do not output any additional action not mentioned in the reasoning.

You should only respond in the format described below, and you should not output comments or other information.

```
Reasoning:
```

```
1. ...
```

 2. ...

3. ...

Actions:

```
python
    action(args1=x, args2=y)
```

Conclusion

The above cases collectively reveal key limitations of current large language models when deployed in grounded, spatially structured environments like Stardew Valley. Despite having access to both visual and textual information, the model consistently exhibits hallucinations in spatial understanding — including confusion about direction, misjudgment of proximity, and failure to plan indirect paths around obstacles.

These behaviors suggest that LLMs, while capable in language-based reasoning, still lack robust internal representations of space and geometry. Moreover, the persistence of such errors even with explicit instructions points to fundamental weaknesses in grounding language to actionable physical reasoning. Bridging this gap will require future work in multimodal integration, spatial memory, and instruction-following mechanisms tailored to embodied agents. Our findings underscore the importance of evaluating LLMs not just on language tasks, but within interactive environments where spatial and physical reasoning are essential.

E FINANCIAL COST

The financial costs of evaluating different closed-source MLLMs on StarDojo-Lite for one round are summarized in Table 6.

Table 6: Financial costs of running one round on StarDojo-Lite with different closed-source MLLMs.

MLLM	Model Name	Cost
GPT-4.1	gpt-4.1-2025-04-14	\$20.6
Gemini 2.5 Pro	gemini-2.5-pro-preview-03-25	\$25.8
Claude 3.7 Sonnet	claude-3-7-sonnet-20250219	\$38.7
GPT-4.1 mini	gpt-4.1-mini-2025-04-14	\$4.3