

TIMEARENA: Shaping Efficient Multitasking Language Agents in a Time-Aware Simulation

Yikai Zhang[♣], Siyu Yuan[♣], Caiyu Hu[♣],
 Kyle Richardson[♡], Yanghua Xiao^{♣*}, Jiangjie Chen^{♣*}
[♣]Fudan University [♡]Allen Institute for AI
 {ykzhang22, syuan21, cyhu24}@m.fudan.edu.cn,
 kyler@allenai.org, {shawyh, jjchen19}@fudan.edu.cn

Abstract

Despite remarkable advancements in emulating human-like behavior through Large Language Models (LLMs), current textual simulations do not adequately address the notion of time. To this end, we introduce TIMEARENA, a novel textual simulated environment that incorporates complex temporal dynamics and constraints that better reflect real-life planning scenarios. In TIMEARENA, agents are asked to complete multiple tasks as soon as possible, allowing for parallel processing to save time. We implement the dependency between actions, the time duration for each action, and the occupancy of the agent and the objects in the environment. TIMEARENA grounds to 30 real-world tasks in cooking, household activities, and laboratory work. We conduct extensive experiments with various state-of-the-art LLMs using TIMEARENA. Our findings reveal that even the most powerful models, e.g., GPT-4, still lag behind humans in effective multitasking, underscoring the need for enhanced temporal awareness in the development of language agents.¹

1 Introduction

Large language models (LLMs) (OpenAI, 2022, 2023; Team and Google, 2023; Touvron et al., 2023a,b) have enabled the development of language agents (*a.k.a.* LLM-based agents), which aim to simulate human behaviors in real-world scenarios through their reasoning and planning capabilities (Wu et al., 2023a; Liu et al., 2023; Gong et al., 2023; Akata et al., 2023). However, planning in the real world involves temporal and resource constraints (Russell and Norvig, 2010), which are rarely implemented in most textual simulations for LLMs and language agents (Wang et al., 2022; Park et al., 2023).

*Corresponding authors.

¹Project page: <https://time-arena.github.io>.

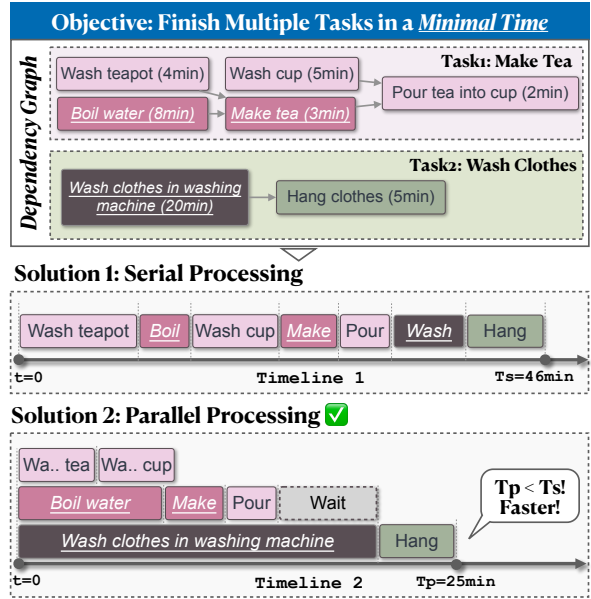


Figure 1: An example illustrating multitasking with temporal constraints in TIMEARENA. The completion of tasks requires actions in a predetermined dependency and order. Underlined actions do not occupy the agent, allowing other actions to be processed by the agent simultaneously. The Wait action skips the current time step, meaning the agent is idle.

The integration of time in simulated environments challenges agents to navigate and align with human-like efficient multitasking skills. Such a simulation requires the agent to consider the following three factors: 1) **Time Duration and Dependency**: Actions will have durations upon dependencies, requiring agents to strategize and prioritize based on time constraints and task completion progress. 2) **Agent Occupancy**: Agents will be occupied by certain actions thus they might be unable to perform other actions at the same time. 3) **Object Occupancy**: Some objects might be occupied for some time, and agents must use available objects in the environment for the tasks. These factors are common in real-life but are seldom addressed by current textual simulations.

To help illustrate, Figure 1 shows an example of completing the task of “make tea” (Task 1) and “wash clothes” (Task 2). The actions of each task might depend on previous actions (*e.g.*, agents must **boil water** before **make tea**.), and each action takes a specific duration in time (*e.g.*, **wash cup** takes 5 minutes). In particular, some actions let agents be idle, allowing agents to carry out other actions. For example, **wash clothes in washing machine** allows agents to perform other actions at the same time. Moreover, actions temporarily occupy objects, making them unavailable for other actions and hindering parallel processing. For example, **boil water** occupies the **pot**, delaying other actions like **cook soup** until it is available. When no action is currently available for the agent, the only option is to wait. For example, in Solution 2, the agent must wait for the completion of **wash clothes in washing machine**, before **hang clothes**.

In this work, we introduce TIMEARENA, a textual simulated environment featuring 30 real-world involving cooking, household activities, and laboratory work. TIMEARENA is the first textual simulation to evaluate language agents on multitasking efficiency. Specifically, we incorporate the time duration of each action and set two types of actions based on agent occupancy. One type occupies agents (*e.g.*, **wash cup**) and another lets agents be idle (*e.g.*, **boil water**). Additionally, we simulate resource competition by implementing object occupancy, *i.e.*, an object used for one task cannot be simultaneously used for another, which is common in parallel processing. Therefore, agents must focus on parallel processing, taking into account the occupancy of agents and objects, to minimize time consumption. We design four metrics in TIMEARENA to evaluate the average progress, completion speed, task completion rate and average completion time. These metrics help to assess and analyze the efficient multitasking capabilities of language agents. Our comprehensive evaluation of 7 LLMs on TIMEARENA shows that current language agents struggle in efficient multitasking. Even the most powerful LLM, GPT-4, still faces challenges in parallel processing.

In summary, our contributions are as follows:

- To the best of our knowledge, we are the first to explore the notion of time of language agents in a textual environment, which is important for more realistic simulation.

- We create TIMEARENA, a novel text-based simulated environment consisting of 30 tasks, where LLMs can complete multiple tasks in parallel.
- Using TIMEARENA, we conduct rich experiments to evaluate the efficient multitasking capabilities of language agents. Our results demonstrate that efficient multitasking in TIMEARENA poses a significant challenge for current LLMs and language agents.

2 Related Work

Simulation-based Evaluation For language Agents With the great success of LLMs (Ouyang et al., 2022; OpenAI, 2022, 2023; Team and Google, 2023; Touvron et al., 2023a,b), recent works have shifted the focus from traditional NLP tasks to explore language agents in simulation environments that mimic real-world scenarios (Wu et al., 2023a; Liu et al., 2023; Gong et al., 2023; Akata et al., 2023). These simulation environments can be divided into two categories: 1) Social Simulations (Park et al., 2023; Mukobi et al., 2023; Zhou et al., 2023), which aim to evaluate the behaviors of language agents in some social scenarios; 2) Problem-solving simulations, which are created based on competitive games (Wang et al., 2023c; Xu et al., 2023; Chen et al., 2023a) or cooperation games (Chen et al., 2023b; Zhang et al., 2023; Agashe et al., 2023) and scientific scenarios (Wang et al., 2022). These simulations mainly test the planning and reasoning abilities of language agents and need agents to solve specific problems within dynamic and evolving environments. In this paper, we focus on problem-solving simulations to investigate the efficient multitasking capabilities of language agents.

Language Planning Language planning aims to decompose a complex task into steps (Schank and Abelson, 1975, 2013). Early studies mainly endow the planning capabilities of LMs through training them on specific planning datasets (Peng et al., 2018; Hua et al., 2019; Kong et al., 2021), which exhibits poor generalization. Recent studies have identified that LLMs can effectively decompose tasks into procedural steps (Olmo et al., 2021; Lu et al., 2022; Ruan et al., 2023; Wu et al., 2023b; Song et al., 2023; Wang et al., 2023d; Yuan et al., 2023; Shen et al., 2023). However, existing work mostly focuses on planning the logical structure

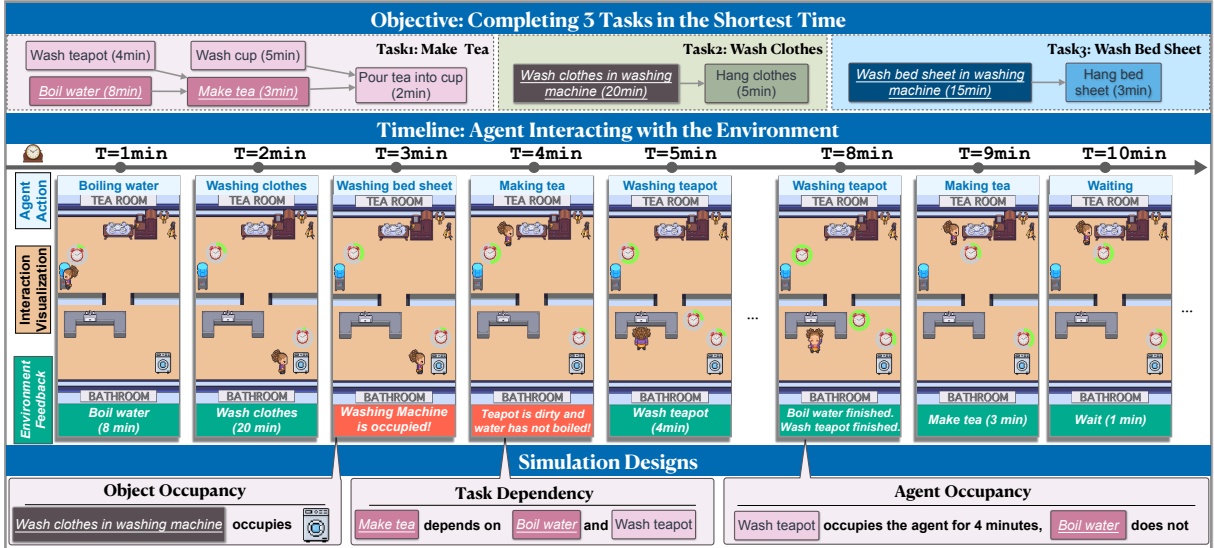


Figure 2: An overview of TIMEARENA, with a multitasking example that shows our designs of the simulation. TIMEARENA first sets an objective for the agent, and then the agent interacts with TIMEARENA over time, with the design of task dependency, object occupancy, and agent occupancy.

of actions, where these actions lack specified time durations and do not accommodate the possibility of agents performing multiple actions concurrently. Multitask planning with parallel processing in dynamic environments still remains under-studied.

Temporal Reasoning Numerous research efforts address diverse challenges in temporal reasoning. Temporal reasoning involves comprehending, structuring, and interpreting events, actions, and states through the lens of time (Allen, 1991; Vila, 1994; Stock, 1998). Previous studies in temporal reasoning focus on temporal relation extraction (Miller et al., 2015; Vashishtha et al., 2019; Mathur et al., 2021; Wang et al., 2023b), event temporal reasoning (Zhou et al., 2021; Qin et al., 2021; Dhingra et al., 2022; Mathur et al., 2022; Yang et al., 2023; Wang and Zhao, 2023) and explore the temporal reasoning capability of LLMs with several contemporary time-sensitive QA datasets (Zhang and Choi, 2021; Shang et al., 2022; Tan et al., 2023). Distinguished from other datasets and benchmarks (Chu et al., 2023), our TIMEARENA creates a dynamic and interactive simulated environment. Agents need to interact with TIMEARENA over time and decide the relationship of actions in the evolving environment.

3 TIMEARENA

We create TIMEARENA, a textual simulated environment to evaluate the efficient multitasking capabilities of language agents. To help illustrate, we

first show an overview and an example run of how an agent interacts with the TIMEARENA environment (§ 3.1), and then describe the design of the simulation environment in more detail (§ 3.2-3.3).

3.1 Overview of TIMEARENA

TIMEARENA challenges agents to complete multiple tasks strategically in the shortest possible time. This simulation emphasizes the importance of understanding, performing, and optimizing actions within a constrained timeframe, mirroring practical scenarios involving time management.

Central to TIMEARENA are **Tasks**, **Objects**, and **Actions**. **Tasks** define the objectives for the agents, **Objects** represent elements in the environment that agents will encounter and interact with, and **Actions** are the means to accomplish these tasks. Real-time feedback and scoring mechanisms are integral to the environment, assessing agent performance and adding to the simulation’s complexity and realism. Unique features like the duration and occupancy of actions and strategic resource utilization distinguish TIMEARENA from other environments.

An Example Run As in Figure 2, consider an agent tasked with “make tea” (Task 1), “wash clothes” (Task 2) and “wash bed sheet” (Task 3). The agent starts by decomposing the task into actions like **boil water**. In TIMEARENA, all actions have a duration (e.g., **Boil water** needs 8 minutes.) and dependencies (e.g., At T=4min, **make tea** violates the dependency because **wash**

teapot and **boil water** are not completed yet.). The agent then interacts with objects (e.g., **wash clothes in washing machine**), which become **occupied** during the process. The agent can engage in non-occupied actions simultaneously (e.g., **wash teapot**) while others (e.g., **boil water**) are in progress. Environmental feedback guides the agent, indicating the legitimacy of actions and the completion of tasks. For example, if the **washing machine** is occupied, the agent adjusts its strategy. The agent’s goal is to complete all tasks efficiently, with performance evaluated based on progress and completion time.

This dynamic interaction in TIMEARENA fosters an environment where strategic planning, resource management, and adaptability are key to an agent’s success.

3.2 Components of TIMEARENA

Tasks In TIMEARENA, We design tasks within three distinct *scenarios* or simulated settings, namely, *household activities*, *cooking*, and *laboratory work*. Each scenario represents a specific context or environment where multitasking is an integral part of the activities involved.² For example, one can do **sweep floor** while doing **boil water**, and do **wash dishes** while doing **cook soup**. Each scenario contains 10 tasks, and some actions and objects are shared across multiple tasks of a scenario. Each task requires multiple actions to be executed, which manipulates the objects in the environment for task completion. In the beginning, TIMEARENA gives a list of tasks to the agent, with a comprehensive task instruction consisting of a *task description*, an *action space*, and an *object set*:

- **Task Description:** Introduces task objectives (e.g., “Task 1: Make a dish of beef fried rice, which consists of cooked rice and fried beef.”);
- **Action Space:** Lists the valid actions for the tasks (e.g., **chop**, **wash**);
- **Object Set:** Lists the available objects in the environment for the tasks (e.g., **pot**, **beaker**).

At every timestep t , the agent needs to generate valid actions on the objects and receive feedback from the environment.

Objects Objects are integral to completing tasks and situating within the environment. In TIMEARENA, there are 71 different objects for

²Details of tasks are in Appendix A.1.

all the tasks. Every task involves a list of objects, which might overlap with other tasks of the same scenario. To mimic the resource limitation in real-world parallel processing, we introduce:

- **Object Occupancy:** the state of the object involved in an action is set to be *occupied*, e.g., **wash cup** will cause the object **cup** to be occupied. This object cannot be processed until the involved action is completed (after some time). Then, this object is reset as non-occupied and waits for another action.

Actions We design a total of 45 actions for all 30 tasks. Each action consists of a detailed description (e.g., **chop OBJ**, **chop the whole item into sliced pieces.**), showing a change of states the action will cause to an object.³ Different from existing text-based simulations (Wang et al., 2022; Gong et al., 2023; Shridhar et al., 2020), in our case, an action has a duration of time and may occupy the agent from performing other actions, to the passage of time. In detail:

- **Action Dependency:** An action within the same task might depend on completing other actions within the same task. As depicted in Figure 2, **make tea** is dependent on **wash teapot**.
- **Duration of Time:** Each action holds a time-frame in the timeline, ranging from 1 to 10 minutes. In practice, agents only have an educated guess of the time duration of each action until actually interacting with TIMEARENA.
- **Agent Occupancy:** One key to parallel processing is agent occupancy, which prevents agents from performing other tasks. Therefore, we consider two types of actions based on agent occupancy: Type 1 action occupies the agent til completion (e.g., **wash teapot**); and Type 2 action lets agents be idle, allowing to perform other actions (e.g., **boil water**).

3.3 The Interaction between Agent and Environment

Environmental Feedback The feedback from a textual environment is important to simulate and implement the constraints in TIMEARENA using only textual messages. We define feedback as the response from the environment following an action by an agent. A feedback message could be of multiple types, including:

³All the actions are listed in Appendix A.2.

- **Invalid Action:** An action attempt that does not match the required format, *e.g.*, “**clean teapot** is invalid”.
- **Action on Non-existing Object:** An action attempt that visits objects that are not in the object set, *e.g.*, “**pan** is non-existent”.
- **Wrong Action Input:** An action attempt that the prerequisite action has not completed (*e.g.*, “Cannot perform action **add to** on object **shrimp**. Because **shrimp** is raw.”) or has been completed (*e.g.*, “**wash beaker** has been completed”).
- **Action on Mismatched Object:** An action attempt that does not match the object, *e.g.*, “You cannot perform **read** on **potato**.”
- **Action on Occupied Object:** An action attempt on occupied objects, *e.g.*, “Object **pot** is being occupied by another action”.

Correspondingly, valid actions will trigger environmental feedback of the following types:

- **Action Start:** Avoiding previous errors, valid actions will receive a feedback message containing the specific performing time, marking the start of the action, *e.g.*, *You are doing wash cup, it will take 9 minutes.*
- **Action Completion:** When an action is completed, the environment will send a message, *e.g.*, *cup is clean*, and reset the occupying state of the object (**cup**).

Progress Score The progress score, denoted as a percentage, reflects the agent’s completion rate of required actions within the environment, where the total duration for all actions is considered as 100%. Each action’s contribution to the progress score is proportionate to its duration. Specifically, if an action’s duration is t_i minutes, its contribution to the progress score is calculated as $s_i = \left(\frac{t_i}{\sum_{j=1}^n t_j} \right) \times 100\%$, with n representing the total number of actions. For instance, an action lasting 5 minutes in a total action duration of 20 minutes contributes 25% to the progress score.

4 Experiments

4.1 Experiment Settings

Task Set Construction In our experiments, we design three categories of task combinations based on the number of tasks: # **Task=1**, # **Task=2** and # **Task=3** scenarios. In # **Task=1** scenario, agents focus on completing one task (*e.g.*, **make tea**). For

Scenario	# Actions	# Objects	Time (min)
Cooking	5.6	5.5	18.9
Household Activity	4.1	3.5	12.8
Laboratory Work	5.3	2.7	16.1

Table 1: Average number of actions and objects per task in each scenario, and the average shortest completion time for these tasks.

the other two scenarios, we combine either two or three tasks from 10 single tasks (*e.g.*, **make tea** and **wash clothes**). Then, we randomly select 10 combined tasks for each scenario.⁴

Interaction Initially, the environment provides a comprehensive task instruction that details the task, action space, and object set. Subsequently, the agent produces an action based on this instruction, adhering to a prescribed format specified in the action space; any deviation is considered invalid. To facilitate action recognition by the environment, regular expressions are employed to parse actions from responses (*e.g.*, extracting **wash clothes** from “*I will wash clothes*”). For each action execution, the agent must incorporate task instructions, previous actions, and feedback from the environment into LLMs as context.⁵

Maximum Time Each combined task is allocated a maximum completion time. We set the time limit for completing a single task at 40 minutes, which exceeds the total time required for all actions in any given task. For tasks that are combined, the time limit is proportionally increased by the number of tasks involved.

Oracle Performance As shown in Tabel 2, *Oracle* represents the optimal performance, including the shortest completion time and the fastest completion rate, which are manually calculated. Specifically, we calculate oracle performance using a greedy strategy: always start the longest non-occupied actions as early as possible and avoid idleness when there are actions to perform.⁶

Finishing The interaction finishes under any of the following conditions: 1) Agents have completed all the actions that solve the tasks (*i.e.*, the progress score reaches 100%) 2) Time has run out;

⁴Appendix B.1 shows examples of single and combined tasks.

⁵Appendix B.2 gives an example of interaction between the agent and the environment.

⁶Appendix A.4 shows our algorithm for calculating the oracle performance.

	Model	# Task=1				# Task=2				# Task=3			
		AS ↑	CS ↑	CR ↑	CT ↓	AS ↑	CS ↑	CR ↑	CT ↓	AS ↑	CS ↑	CR ↑	CT ↓
Cooking	Mistral-7B	63.70	3.59	30.00	25.67	42.20	1.49	0	-	39.40	1.06	0	-
	OpenChat-3.5	76.30	<u>3.89</u>	30.00	<u>20.33</u>	37.10	1.80	0	-	<u>41.00</u>	1.17	0	-
	Vicuna-13B	84.60	4.10	60.00	21.83	48.80	1.76	0	-	26.00	1.03	0	-
	Mixtral-8x7B	50.80	3.81	<u>10.00</u>	19.00	40.10	1.99	0	-	27.60	1.17	0	-
	Gemini Pro	78.30	3.57	50.00	24.60	31.00	1.75	0	-	18.50	1.26	0	-
	GPT-3.5	77.70	3.61	30.00	24.33	52.30	1.87	0	-	33.10	<u>1.23</u>	0	-
	GPT-4	98.70	3.48	90.00	28.22	93.50	1.83	70.00	<u>52.57</u>	82.50	<u>1.21</u>	40.00	76.25
	+ Self-plan	<u>89.00</u>	3.83	<u>60.00</u>	26.50	<u>64.90</u>	2.05	<u>10.00</u>	37.00	26.20	1.15	0	-
Oracle	100	5.31	100	18.90	100	2.85	100	35.00	100	1.94	100	52.50	
Household Activity	Mistral-7B	64.80	6.00	20.00	15.50	45.30	2.46	0	-	49.90	1.78	0	-
	OpenChat-3.5	70.50	5.34	30.00	15.67	68.20	2.73	0	-	44.30	<u>1.83</u>	0	-
	Vicuna-13B	69.50	5.94	40.00	14.25	45.90	2.34	0	-	24.90	1.69	0	-
	Mixtral-8x7B	68.80	6.08	40.00	<u>15.00</u>	51.60	2.85	10.0	<u>31.00</u>	60.20	<u>1.83</u>	10.00	58.00
	Gemini Pro	68.10	5.92	40.00	16.50	60.50	3.02	10.00	25.00	40.30	1.93	0	-
	GPT-3.5	<u>87.40</u>	5.98	70.00	16.71	63.80	2.57	10.00	36.00	45.30	1.82	0	-
	GPT-4	100	5.81	100	17.20	100	<u>2.89</u>	100	34.50	98.40	1.82	90.00	<u>54.78</u>
	+ Self-plan	87.20	6.01	<u>80.00</u>	16.37	84.50	2.80	<u>50.00</u>	35.20	<u>95.30</u>	1.93	<u>60.00</u>	50.16
Oracle	100	7.81	100	12.80	100	4.23	100	23.60	100	2.82	100	35.40	
Laboratory Work	Mistral-7B	70.80	4.39	30.00	21.67	47.10	2.27	0	-	38.40	1.37	0	-
	OpenChat-3.5	65.50	5.07	30.00	13.33	45.80	2.10	0	-	27.50	1.30	0	-
	Vicuna-13B	59.60	3.94	20.00	26.00	20.80	1.87	0	-	22.90	1.40	0	-
	Mixtral-8x7B	64.10	4.57	40.00	24.25	41.80	2.43	0	-	32.40	1.58	0	-
	Gemini Pro	88.00	<u>5.17</u>	70.00	19.57	57.50	2.64	<u>20.00</u>	35.50	25.70	1.61	0	-
	GPT-3.5	71.50	4.52	30.00	22.00	47.60	2.17	0	-	37.90	1.52	0	-
	GPT-4	97.50	5.32	90.00	18.67	85.30	2.61	50.00	39.20	83.10	<u>1.71</u>	60.00	<u>60.33</u>
	+ Self-plan	<u>95.30</u>	5.09	<u>80.00</u>	20.12	<u>83.00</u>	2.79	50.00	<u>36.40</u>	<u>70.00</u>	1.87	60.00	54.66
Oracle	100	6.21	100	16.1	100	4.14	100	24.60	100	2.84	100	35.50	

Table 2: Model performance under different task combination settings in TIMEARENA. We report Average Progress Score (AS), Completion Speed (CS), Task Completion Rate (CR), and Average Completion Time (CT). #Task= n represents that agents are required to do n tasks altogether. We also list the Oracle result for comparison. The best results are **bolded**, and the second best ones are underlined.

3) Agents who have performed incorrect actions 5 times in a row are considered to fail the task.

Model Choice We select a variety of language models to drive the agent: Mistral-7B (Jiang et al., 2023) by MistralAI, OpenChat-3.5 (Wang et al., 2023a) which is fine-tuned from a 7B Mistral model fine-tuned with C-RLFT (openchat_3.5), Vicuna-13B (Chiang et al., 2023) which is fine-tuned from a 13B LLaMA model (Touvron et al., 2023c) with instructions and user-shared conversations, Mixtral-8x7B (Mistral AI team, 2023) by MistralAI, which is a Mixture-of-Expert version of Mistral, Gemini Pro by Google (Team et al., 2023), GPT-3.5 (OpenAI, 2022) by OpenAI (gpt-3.5-turbo-1106) and GPT-4 (OpenAI, 2023) by OpenAI (gpt-4-1106-preview).

4.2 Evaluation Metrics

To comprehensively evaluate the ability of agents to multitask, we consider both time and score and

design the following four metrics:

- **Average Progress Score (score, AS):** The average highest progress score achievable by an agent:

$$AS = \frac{\sum_{i \in N} P_i}{N}$$

where P_i denotes the maximum progress score of i -th task that agents can reach, and N denotes the number of all tasks.

- **Completion Speed (score per minute, CS):** The average of the highest score divided by the time taken to achieve it:

$$CS = \frac{\sum_{i \in N} P_i}{\sum_{i \in N} T_i}$$

where T_i denotes the time required to reach P_i of i -th task.

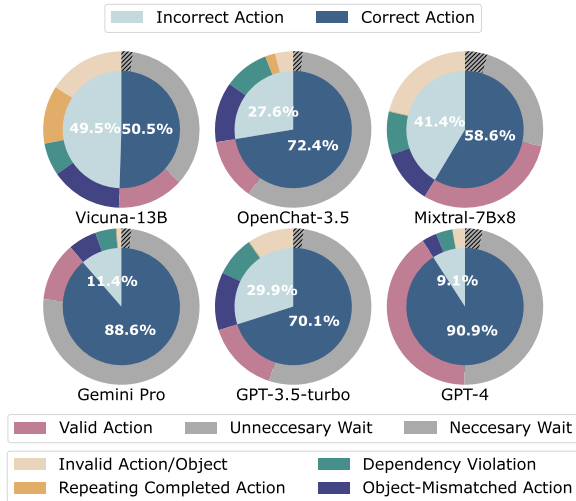


Figure 3: The proportions of correct and incorrect actions of each language agent.

- **Task Completion Rate (% , CR):** The rate of successfully completed tasks:

$$CR = \frac{S}{N}$$

where S denotes the tasks completed successfully. Notably, when combining tasks, a combined task counts as one task.

- **Average Completion Time (minutes, CT):** The average time taken for completing tasks successfully:

$$CT = \frac{\sum_{i \in S} T_i}{S}$$

4.3 Main Results

Based on the performance of language agents in Table 2, GPT-4 achieves the best performance across different task combinations. Moreover, the combined tasks are more challenging than single tasks despite the longer time given. Apart from GPT-4, most models fail to complete 2 or 3 tasks, showing their limited multitasking abilities and the challenging nature of our environment.

For open-source models, OpenChat-3.5 and Vicuna-13B are even better than GPT-3.5, demonstrating the potential of open-sourced models to develop multitasking capabilities. However, although OpenChat-3.5 and Vicuna-13B exhibit faster completion speed and shorter average completion time than GPT-4, a lower task completion rate indicates that these models quickly complete simple actions initially but then encounter difficulties; they either

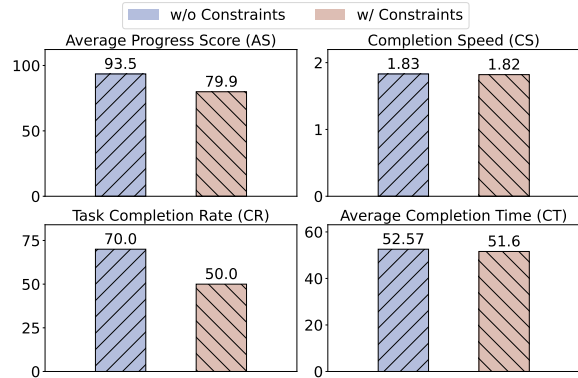


Figure 4: Comparison of the performance of GPT-4 with and without resource constraints. We imposed constraints by limiting to a single instance each of **pot**, **fryer**, and **oven**.

get caught in repetitive actions or fail to properly segment subsequent tasks, which significantly impacts task performance. For example, initially, **potato** is unpicked, so the agent first performs **pick potato**. Subsequently, the agent mistakenly opts for **cook potato in pot** rather than the correct **chop potato**, because it incorrectly decomposes the task.

To explore the potential of heuristic algorithms in improving model performance, we introduce *self-plan prompting* to GPT-4, as illustrated in Appendix B.3. Under this method, the model initially discovers the dependencies among actions, task descriptions, and objects and estimates the duration of each action. It then adopts a greedy strategy similar to **Oracle Performance**, favoring selecting the longest-duration actions that do not require continuous engagement from the agent in the task model to formulate a plan. Then, the agent executes this plan through interactions with the environment. However, the results indicate that self-plan prompting is outperformed by vanilla GPT-4. There are three possible reasons for such performance: 1) The difficulty in accurately parsing actions and identifying their dependencies; 2) The reliance on estimating action durations might introduce cascading errors, leading to inaccurate results of the greedy strategy; 3) The rigid adherence to flawed plans, without adapting to the dynamic nature of interactions with the environment, leads to its failure.

4.4 Analysis

Can language agents master multitasking? To further explore whether language agents can master multitasking, we conduct detailed analyses to inves-

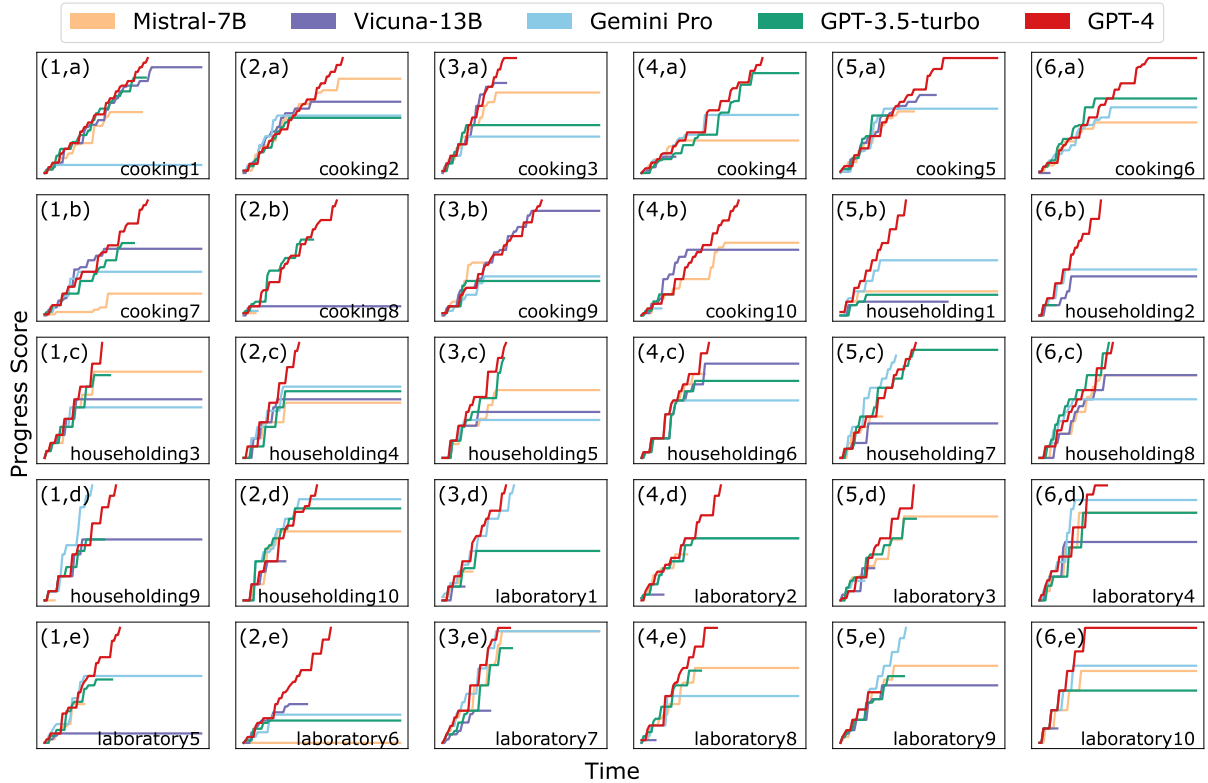


Figure 5: Task progress score curves of language agents on two task combinations in TIMEARENA. The names at the bottom-right indicate the scenario and task number. For example, cooking1 represents the first combination of tasks in the cooking scenario.

tigate the types of actions these models can perform. Based on environmental feedback, the actions are classified into two categories: correct ones and incorrect ones. We further define six fine-grained types of actions: 1) **Correct Actions**: Valid Action, Wait 2) **Incorrect Actions**: Invalid Action/Object, Dependency Violation, Repeating Completed Action and Object-Mismatched Action.⁷

We calculate the frequency of these actions of each agent throughout their interactions. The results in Figure 3 show that a significant proportion of invalid actions are due to dependency violations and mismatches with objects. Multitasking involves performing several tasks simultaneously. As the number of tasks increases, the complexity of objects and actions escalates, leading to intricate dependencies between actions. Thus, the high proportion of actions that violate dependencies and mismatch objects suggests that language agents face challenges in managing complex action interdependencies during multitasking, indicating a limitation in their multitasking capabilities.

⁷Detailed description of different types of actions can be found in Appendix A.3.

Are language agents aware of parallel processing? Parallel processing can significantly reduce the time required for efficient multitasking. As stated in § 3.2, we consider the duration of time for each action in TIMEARENA to explore the parallel processing capabilities of language agents. If an agent is capable of parallel processing, it can engage in additional actions instead of unnecessary waiting for the current action to complete. To answer this question, we decompose **wait** action into two types: **necessary wait** and **unnecessary wait**. The former represents that no actions can currently be performed, requiring waiting for other actions to complete. In particular, we report the maximum number of **necessary wait**. **Unnecessary wait** indicates that there are other action options available. Figure 3 shows that **wait** actions constitute over half of the valid actions performed by different LLMs, and **necessary wait** only accounts for a small part of it. This indicates a tendency for agents to engage in unnecessary waiting, showing their ignorance of parallel processing and inability to complete tasks in minimal time (Table 2).

Do resource constraints affect the multitasking of language agents?

Resource constraints refer to limitations in the availability of resources (*e.g.*, the number of objects) necessary for task completion, which is rather common in real life. To design resource constraints, we first select three objects: **pot**, **fryer** and **oven** in the cooking scenario, and choose # Task=2 setting in Table 2. Then, we set that there is only one instance of each of the three objects, simulating the limitation of resources in the environment. Figure 4 compares GPT-4’s performance before and after applying these constraints. We find that the constraints do not affect the task completion time or speed, revealing that GPT-4 rarely attempts to process tasks in parallel. However, a noticeable decline in both completion rate and progress score indicates that the constraints prevent the models from better comprehending and decomposing multiple tasks.

Language agents trapped in an infinite loop.

To delve into why language agents struggle with multiple tasks, we analyze the progress score changes over time for various models. As illustrated in Figure 5, Vicuna, Mistral, Gemini and GPT-3.5 often cease scoring without completing all tasks, maintaining low scores until time runs out (*e.g.*, (5,b), (2,c) and (6,d)). We further examine their actions during these periods and find that they always perform **incorrect actions and waiting** alternately. Since **wait** is a valid action, repeatedly alternating between waiting and incorrect actions does not lead to task failure, but neither does it contribute to an increase in scores. To find out whether agents wait for good reasons, we ask them to explain each action via the chain-of-thought prompting strategy, and they often believe **wait** can pause incorrect actions. However, they find it hard to adjust their incorrect actions based on feedback after waiting, resulting in them being trapped in infinite loops.

5 Conclusion

In this paper, we introduce TIMEARENA, a text-based simulated environment designed to incorporate the notion of *time*. TIMEARENA extends beyond simply acknowledging the dependency of actions by also considering their duration, an essential factor in time modeling. Using TIMEARENA, we evaluate the multitasking and parallel processing capability of language agents. Our findings indicate that as tasks become more complex, the

models struggle to complete them and often fail to recognize opportunities for parallel processing. This reveals that language agents still have significant room for improvement when completing multiple tasks in dynamic environments, highlighting an area for future research.

Limitations

In TIMEARENA, we implement detailed descriptions of tasks and environments, along with fine-grained textual feedback to simulate interactions. However, TIMEARENA is still designed as a textual simulation for LLMs, lacking visual information that might be necessary for agents to succeed in real-world tasks. For example, in the laboratory work scenario, it is challenging to completely represent chemical reactions through text due to their complexity. The number of tasks and scenarios is limited, while the number of multitasking scenarios that allow parallel processing is large in real life. Moreover, in TIMEARENA, agents interact with the environment only through actions that are explicitly presented in action prompts, rather than exploring freely. Also, whether an action occupies an agent sometimes depends on specific conditions. For instance, the action **cook beef** is classified as non-occupying in TIMEARENA, implying that it does not engage agents continuously. Yet, in reality, this action requires attention, such as turning the beef to prevent burning, a detail TIMEARENA overlooks, potentially reducing the realism of our simulation.

Ethical Statement

We hereby acknowledge that all authors of this work are aware of the provided ACL Code of Ethics and honor the code of conduct.

Use of Human Annotations Our institution recruited three annotators to implement the task creation for three scenarios. We ensure the privacy rights of the annotators are respected during the annotation process. The annotators receive compensation exceeding the local minimum wage and have consented to tasks generated for TIMEARENA for research purposes.

Risks The TIMEARENA in our experiment is created by human annotators, and we further examine them to guarantee that they are devoid of socially harmful or toxic language. However, evaluating the data quality of tasks is based on common sense,

which can vary among individuals from diverse backgrounds.

Acknowledgement

We would like to thank Xintao Wang, Ruihan Yang, Tinghui Zhu from Fudan University for their valuable comments and suggestions for the manuscript. We would also like to thank Peter Jansen from University of Arizona and Bodhisattwa Prasad Majumder from Allen Institute for AI for fruitful discussions that helped shape this project at an early stage.

References

- Saaket Agashe, Yue Fan, and Xin Eric Wang. 2023. [Evaluating multi-agent coordination abilities in large language models](#).
- Elif Akata, Lion Schulz, Julian Coda-Forno, Seong Joon Oh, Matthias Bethge, and Eric Schulz. 2023. [Playing repeated games with large language models](#).
- James F Allen. 1991. Planning as temporal reasoning. *KR*, 91:3–14.
- Jiangjie Chen, Siyu Yuan, Rong Ye, Bodhisattwa Prasad Majumder, and Kyle Richardson. 2023a. Put your money where your mouth is: Evaluating strategic planning and execution of llm agents in an auction arena. *arXiv preprint arXiv:2310.05746*.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. 2023b. [Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors](#).
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023).
- Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Haotian Wang, Ming Liu, and Bing Qin. 2023. [Timebench: A comprehensive evaluation of temporal reasoning abilities in large language models](#).
- Bhuvan Dhingra, Jeremy R. Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W. Cohen. 2022. [Time-Aware Language Models as Temporal Knowledge Bases](#). *Transactions of the Association for Computational Linguistics*, 10:257–273.
- Ran Gong, Qiuyuan Huang, Xiaojian Ma, Hoi Vo, Zane Durante, Yusuke Noda, Zilong Zheng, Song-Chun Zhu, Demetri Terzopoulos, Li Fei-Fei, and Jianfeng Gao. 2023. [Mindagent: Emergent gaming interaction](#).
- Xinyu Hua, Zhe Hu, and Lu Wang. 2019. [Argument generation with retrieval, planning, and realization](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2661–2672, Florence, Italy. Association for Computational Linguistics.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Xiangzhe Kong, Jialiang Huang, Ziquan Tung, Jian Guan, and Minlie Huang. 2021. [Stylized story generation with style-guided planning](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2430–2436, Online. Association for Computational Linguistics.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. 2023. [Agentbench: Evaluating llms as agents](#).
- Yujie Lu, Weixi Feng, Wanrong Zhu, Wenda Xu, Xin Eric Wang, Miguel Eckstein, and William Yang Wang. 2022. Neuro-symbolic causal language planning with commonsense prompting. *arXiv preprint arXiv:2206.02928*.
- Puneet Mathur, Rajiv Jain, Franck Dernoncourt, Vlad Morariu, Quan Hung Tran, and Dinesh Manocha. 2021. [TIMERS: Document-level temporal relation extraction](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 524–533, Online. Association for Computational Linguistics.
- Puneet Mathur, Vlad Morariu, Verena Kaynig-Fittkau, Jiuxiang Gu, Franck Dernoncourt, Quan Tran, Ani Nenkova, Dinesh Manocha, and Rajiv Jain. 2022. [DocTime: A document-level temporal dependency graph parser](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 993–1009, Seattle, United States. Association for Computational Linguistics.
- Timothy Miller, Steven Bethard, Dmitriy Dligach, Chen Lin, and Guergana Savova. 2015. [Extracting time expressions from clinical text](#). In *Proceedings of BioNLP 15*, pages 81–91, Beijing, China. Association for Computational Linguistics.
- Mistral AI team. 2023. [Mixtral of experts](#). Accessed: 2023-12-15.
- Gabriel Mukobi, Hannah Erlebach, Niklas Lauffer, Lewis Hammond, Alan Chan, and Jesse Clifton. 2023. Welfare diplomacy: Benchmarking language model cooperation. *arXiv preprint arXiv:2310.08901*.

- Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. 2021. Gpt3-to-plan: Extracting plans from text using gpt-3. *FinPlan 2021*, page 24.
- OpenAI. 2022. [Chatgpt](#).
- OpenAI. 2023. [Gpt-4 technical report](#).
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). In *Advances in Neural Information Processing Systems*.
- Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. [Generative agents: Interactive simulacra of human behavior](#).
- Baolin Peng, Xiujun Li, Jianfeng Gao, Jingjing Liu, and Kam-Fai Wong. 2018. [Deep Dyna-Q: Integrating planning for task-completion dialogue policy learning](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2182–2192, Melbourne, Australia. Association for Computational Linguistics.
- Lianhui Qin, Aditya Gupta, Shyam Upadhyay, Luheng He, Yejin Choi, and Manaal Faruqui. 2021. [TIME-DIAL: Temporal commonsense reasoning in dialog](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7066–7076, Online. Association for Computational Linguistics.
- Jingqing Ruan, Yihong Chen, Bin Zhang, Zhiwei Xu, Tianpeng Bao, Guoqing Du, Shiwei Shi, Hangyu Mao, Ziyue Li, Xingyu Zeng, and Rui Zhao. 2023. [Tptu: Large language model-based ai agents for task planning and tool usage](#).
- Stuart J Russell and Peter Norvig. 2010. *Artificial intelligence a modern approach*. London.
- Roger C Schank and Robert P Abelson. 1975. Scripts, plans, and knowledge. In *IJCAI*, volume 75, pages 151–157.
- Roger C Schank and Robert P Abelson. 2013. *Scripts, plans, goals, and understanding: An inquiry into human knowledge structures*. Psychology Press.
- Chao Shang, Guangtao Wang, Peng Qi, and Jing Huang. 2022. [Improving time sensitivity for question answering over temporal knowledge graphs](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8017–8026, Dublin, Ireland. Association for Computational Linguistics.
- Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang, Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng Li, and Yueting Zhuang. 2023. [Taskbench: Benchmarking large language models for task automation](#).
- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. 2020. AlfworlD: Aligning text and embodied environments for interactive learning. *arXiv preprint arXiv:2010.03768*.
- Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M. Sadler, Wei-Lun Chao, and Yu Su. 2023. [Llm-planner: Few-shot grounded planning for embodied agents with large language models](#). In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Oliviero Stock. 1998. *Spatial and temporal reasoning*. Springer Science & Business Media.
- Qingyu Tan, Hwee Tou Ng, and Lidong Bing. 2023. [Towards benchmarking and improving the temporal reasoning capability of large language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14820–14835, Toronto, Canada. Association for Computational Linguistics.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Gemini Team and Google. 2023. [Gemini: A family of highly capable multimodal models](#).
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu,

- Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023c. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Siddharth Vashishtha, Benjamin Van Durme, and Aaron Steven White. 2019. [Fine-grained temporal relation extraction](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2906–2919, Florence, Italy. Association for Computational Linguistics.
- Lluís Vila. 1994. A survey on temporal reasoning in artificial intelligence. *Ai Communications*, 7(1):4–28.
- Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. 2023a. Openchat: Advancing open-source language models with mixed-quality data. *arXiv preprint arXiv:2309.11235*.
- Haoyu Wang, Hongming Zhang, Yuqian Deng, Jacob Gardner, Dan Roth, and Muhao Chen. 2023b. [Extracting or guessing? improving faithfulness of event temporal relation extraction](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 541–553, Dubrovnik, Croatia. Association for Computational Linguistics.
- Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. 2022. [ScienceWorld: Is your agent smarter than a 5th grader?](#) In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11279–11298, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Shenzhi Wang, Chang Liu, Zilong Zheng, Siyuan Qi, Shuo Chen, Qisen Yang, Andrew Zhao, Chaofei Wang, Shiji Song, and Gao Huang. 2023c. [Avalon’s game of thoughts: Battle against deception through recursive contemplation](#).
- Yuqing Wang and Yun Zhao. 2023. Tram: Benchmarking temporal reasoning for large language models. *arXiv preprint arXiv:2310.00835*.
- Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. 2023d. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. *arXiv preprint arXiv:2302.01560*.
- Yue Wu, Xuan Tang, Tom M. Mitchell, and Yuanzhi Li. 2023a. [Smartplay: A benchmark for llms as intelligent agents](#).
- Zhenyu Wu, Ziwei Wang, Xiuwei Xu, Jiwen Lu, and Haibin Yan. 2023b. Embodied task planning with large language models. *arXiv preprint arXiv:2305.03716*.
- Zelai Xu, Chao Yu, Fei Fang, Yu Wang, and Yi Wu. 2023. [Language agents with reinforcement learning for strategic play in the werewolf game](#).
- Sen Yang, Xin Li, Lidong Bing, and Wai Lam. 2023. [Once upon a time in graph: Relative-time pretraining for complex temporal reasoning](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11879–11895, Singapore. Association for Computational Linguistics.
- Siyu Yuan, Jiangjie Chen, Ziquan Fu, Xuyang Ge, Soham Shah, Charles Jankowski, Yanghua Xiao, and Deqing Yang. 2023. [Distilling script knowledge from large language models for constrained language planning](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4303–4325, Toronto, Canada. Association for Computational Linguistics.
- Hongxin Zhang, Weihua Du, Jiaming Shan, Qinrong Zhou, Yilun Du, Joshua B. Tenenbaum, Tianmin Shu, and Chuang Gan. 2023. [Building cooperative embodied agents modularly with large language models](#).
- Michael Zhang and Eunsol Choi. 2021. [SituatQA: Incorporating extra-linguistic contexts into QA](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7371–7387, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ben Zhou, Kyle Richardson, Qiang Ning, Tushar Khot, Ashish Sabharwal, and Dan Roth. 2021. Temporal reasoning on implicit events from distant supervision. *Proceedings of NAACL*.
- Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haofei Yu, Zhengyang Qi, Louis-Philippe Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, and Maarten Sap. 2023. [Sotopia: Interactive evaluation for social intelligence in language agents](#).

A TIMEARENA Details

A.1 Tasks

TIMEARENA contains 30 tasks in cooking, household activity, and laboratory work scenarios. To illustrate how to complete a task, we show the flow chart for each task in Figure 6, Figure 7 and Figure 8.

A.2 Actions

The environment implements 45 actions, and each action has a description. We show the details of these actions in Table 3.

A.3 Action Types

As shown in Table 4, we define 4 incorrect action types and 2 correct action types for analyzing why agents fail in multitasking.

A.4 Greedy Algorithm

We show the greedy algorithm in Algorithm 1.

B Examples of TIMEARENA

B.1 Tasks

Table 5, 6 and 7 present some examples of task combinations in TIMEARENA for a better understanding.

B.2 Interaction

Table 5 shows an example of interaction between an agent and environment in the cooking scenario.

B.3 Self-plan

Table 9 shows the prompt of the self-plan method.

Algorithm 1: Greedy Algorithm for Minimal Time Calculation

Input: Set of actions \mathcal{A} , Durations \mathcal{T} , Dependencies $p(\mathcal{A})$.

Output: Minimal time \mathcal{T}_{\min} .

```
1 Define non-occupied actions  $\mathcal{A}^*$  and
  occupied actions  $\mathcal{A}'$  from  $\mathcal{A}$ .
2 Sort  $\mathcal{A}^*$  by  $\mathcal{T}$  in descending order.
3  $\mathcal{A} \leftarrow \text{concatenate}(\mathcal{A}^*, \mathcal{A}')$ .
4 Initialize Action_list as an empty list.
5 foreach  $a_i \in \mathcal{A}$  do
6    $P \leftarrow \text{BFS}(a_i, p(a_i))$  to collect
   prerequisites.
7   foreach  $p_i \in P$  do
8     if  $p_i \in \mathcal{A}$  then
9       Action_list.append( $p_i$ ).
10      Remove  $p_i$  from  $\mathcal{A}$ .
11    end
12  end
13  Action_list.append( $a_i$ ).
14 end
15  $\mathcal{T}_{\min} \leftarrow 0$ .
16 while not empty  $\mathcal{A}^*$  or  $\mathcal{A}'$  do
17   foreach  $a_i \in \text{Action\_list}$  do
18     if check_dependency( $a_i$ ) then
19       if  $a_i \in \mathcal{A}^*$  then
20          $\mathcal{T}_{\min} \leftarrow \mathcal{T}_{\min} + 1$ .
21         Remove  $a_i$  from  $\mathcal{A}^*$ .
22       else
23          $\mathcal{T}_{\min} \leftarrow \mathcal{T}_{\min} + \mathcal{T}(a_i)$ .
24         Remove  $a_i$  from  $\mathcal{A}'$ .
25       end
26     break.
27   end
28 end
29 Increment  $\mathcal{T}_{\min}$  by 1 if no action is
  performed.
30 end
```

Action	Description
<code>pick OBJ</code>	Pick the unpicked item
<code>cook OBJ1 in OBJ2</code>	Cook the raw item until it's cooked through
<code>chop OBJ</code>	Chop the whole item into sliced pieces
<code>fry OBJ1 in OBJ2</code>	Fry the raw item until it is fried to perfection
<code>wash OBJ</code>	Wash the dirty item to make clean
<code>bake OBJ1 in OBJ2</code>	Bake the raw item in the oven until it's roasted
<code>activate OBJ</code>	Activate the inactive device to turn it active
<code>pour OBJ1 into OBJ2</code>	Pour the liquid in item into the empty container until it is full
<code>brew OBJ1 with OBJ2</code>	Brew the dry item leaves with the container until they're steeped
<code>gather OBJ</code>	Gather the scattered items until it is collected
<code>scrape OBJ1 into OBJ2</code>	Scrape the contents from the full item into th empty item
<code>place OBJ1 into OBJ2</code>	Place the unplaced item into the right place
<code>fill OBJ1 with OBJ2</code>	Fill the container with something
<code>hoe OBJ</code>	Hoe the uncultivated item until it is cultivated and ready for planting
<code>weed_with OBJ</code>	Weed with the item
<code>set_up OBJ</code>	Set up the item that is not set yet until it is already set
<code>iron OBJ</code>	Iron the wrinkled item until they are smooth
<code>put OBJ1 on OBJ2</code>	Put the item on the right place
<code>add OBJ1 to OBJ2</code>	Add one item to the container
<code>rinse OBJ</code>	Rinse the dry item
<code>find OBJ</code>	Find the missed item so that it is found and can be used
<code>heat OBJ</code>	Heat the cool item until it is hot
<code>dilute OBJ</code>	Dilute the concentrated item until it is diluted
<code>cut OBJ</code>	Cut the whole item into divided pieces
<code>dissolve OBJ1 in OBJ2</code>	Dissolve the solid item in the liquid until it is dissolved
<code>polish OBJ</code>	Polish the rusty item until it is polished
<code>empty OBJ</code>	Empty the full item until it is empty
<code>hanging OBJ</code>	Hang the item
<code>water OBJ1 by OBJ2</code>	Water the item by something
<code>trim OBJ</code>	Trim the overgrown item
<code>plant OBJ</code>	Plant the uncultivated item until it is planted
<code>store OBJ</code>	Store the unstored item
<code>stir OBJ1 with OBJ2</code>	Stir the separate liquid in item with something until it is homogeneous
<code>soak OBJ1 in OBJ2</code>	Soak the dry item in something until it is wet
<code>mop OBJ</code>	Mop the dirty item until it is clean
<code>read OBJ</code>	Read the unknown item
<code>fold OBJ</code>	Fold the spread item until it is tidy
<code>crush OBJ</code>	Crush the intact item until it is crushed
<code>cool OBJ</code>	Cool the hot item until it is cool
<code>dry OBJ</code>	Dry the item until it is dry
<code>wipe OBJ</code>	Wipe the dirty item until it is clean
<code>put OBJ1 in OBJ2</code>	Put the item in something
<code>label OBJ</code>	Give the ambiguous item a label
<code>crystallize OBJ</code>	Crystallize the fluid item until it is crystallized
<code>filter OBJ</code>	Filter the mixed item until it is refined

Table 3: Details of actions with descriptions.

Type	Subtype	Explanation	Example: Make tea
Incorrect Actions	Invalid Action/Object	An action does not in the action space or non-existent objects are visited.	<Valid Actions> activate; wash; brew with; pour into
	Repeating Completed Action	An action is in the action space and matches the objects, but it has already been completed.	<Objects> tea(dry); kettle(inactive); teapot(dirty); cup(dirty)
	Dependency Violation	An action is in the action space and matches the objects, but the necessary prerequisite actions have not been completed.	<Trajectory> T=1: clean teapot T=2: brew tea with teapot T=3: wash teapot T=4: wash kettle T=5: wash teapot T=6: activate kettle T=7: wait
	Object-Mismatched Action	An action is in the action space and the object is available, but they do not match.	...
Correct Actions	Valid Action	An action is in the action space and matches the objects.	
	Wait	An action is used to pass the current time.	

Table 4: Action types and their explanations with an example.

As an AI agent, your objective is to efficiently complete a series of tasks as described. You must adhere to the specific requirements and constraints of each task, including dependencies and timing. Efficiency is key; complete all tasks in the shortest possible time. I will provide instructions regarding actions and objects.

****Action Protocol**:**

- You can perform only one action at a time.
- After each observation from the environment, output an action based on that observation and the instructions.
- Actions fall into two categories:
 - Continuous Actions: Perform these actions until completion (e.g., "wash OBJ").
 - Autonomous Actions: These progress over time, allowing simultaneous tasks (e.g., "heat OBJ").
- Follow the "Valid Actions" format for your output (e.g., "wash cup").
- If no action is required, use "wait" to skip the current time.
- Output the action explicitly (e.g., "wash cup").
- Select object names (OBJ) from the list of Available Objects (e.g., use "rice" instead of "cooked rice").

****Task 1****

<Description>

- Prepare a noodle dish, which consists of cooked noodle, fried mushrooms and shrimp.

<Valid Actions and Usages>

- pick OBJ: Pick the unpicked item.
- cook OBJ1 in OBJ2: Cook the raw item until it's cooked through.
- chop OBJ: Chop the whole item into sliced pieces.
- fry OBJ1 in OBJ2: Fry the raw item until it is fried to perfection.
- add OBJ1 to OBJ2: Add one item to the container.
- wash OBJ: Wash the dirty item to make clean.
- wait: pass the current time without doing anything.

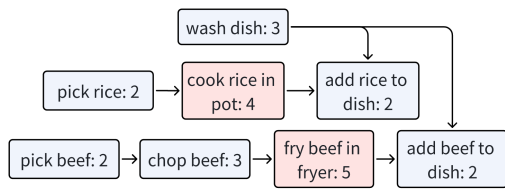
****All Available Objects (OBJ)****

noodle; mushroom; shrimp; fryer; pot; dish

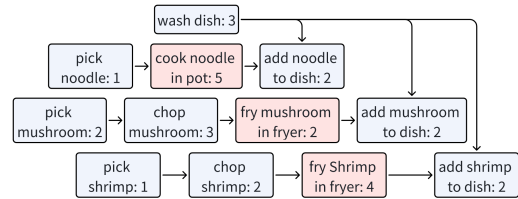
****The Initial States of Objects****

noodle: unpicked; mushroom: unpicked; shrimp: unpicked; fryer: empty; pot: empty; dish: dirty

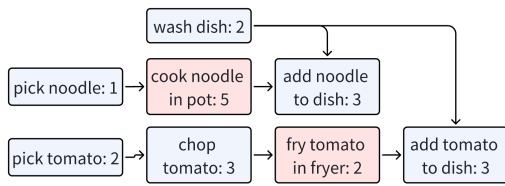
Table 5: An example of # Task=1 scenario.



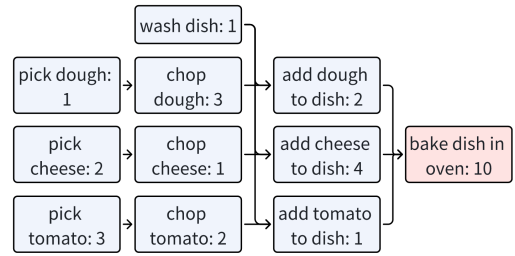
(a) The first task in the cooking scenario.



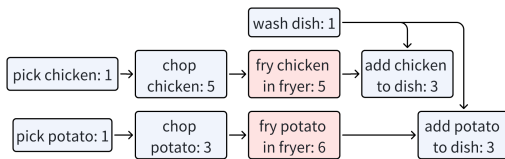
(b) The second task in the cooking scenario.



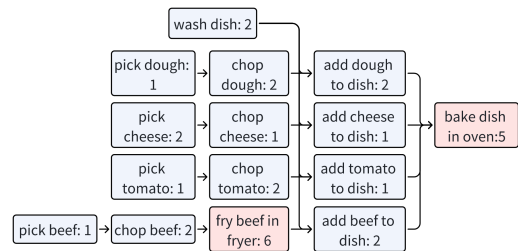
(c) The third task in the cooking scenario.



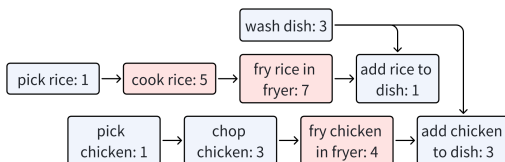
(d) The fourth task in the cooking scenario.



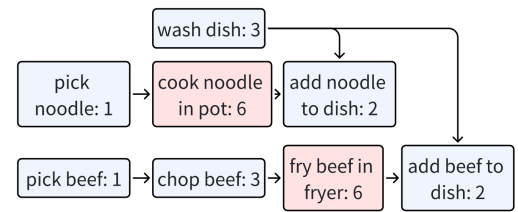
(e) The fifth task in the cooking scenario.



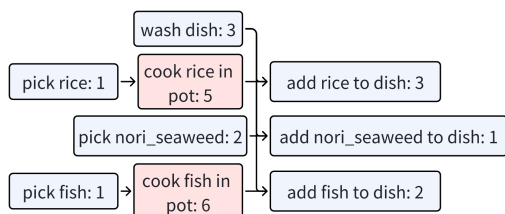
(f) The sixth task in the cooking scenario.



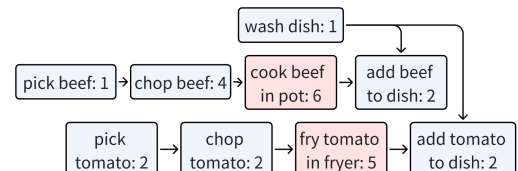
(g) The seventh task in the cooking scenario.



(h) The eighth task in the cooking scenario.

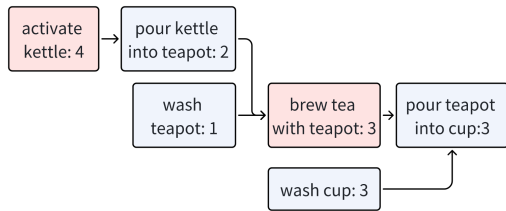


(i) The ninth task in the cooking scenario.

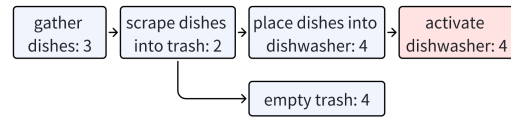


(j) The tenth task in the cooking scenario.

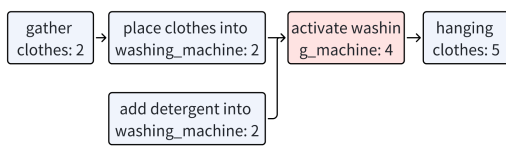
Figure 6: The action dependencies and durations for the ten tasks in the cooking scenario. Actions that occupy the agent, preventing them from doing anything else, are indicated with a blue background. In contrast, actions not occupying the agent, allowing for parallel tasks, are marked with a red background.



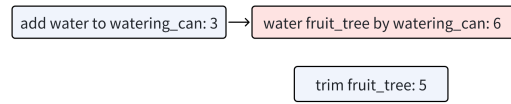
(a) The first task in the household activity scenario.



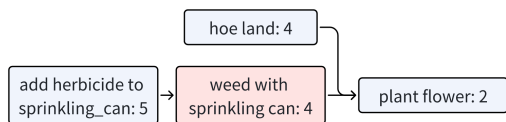
(b) The second task in the household activity scenario.



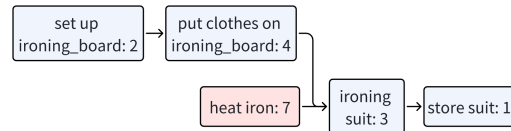
(c) The third task in the household activity scenario.



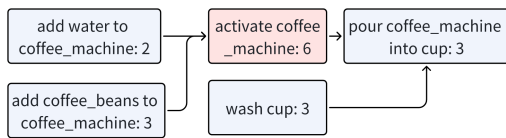
(d) The fourth task in the household activity scenario.



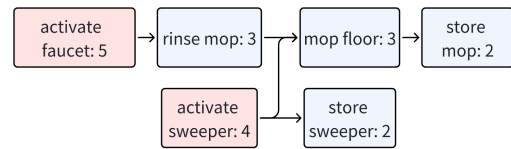
(e) The fifth task in the household activity scenario.



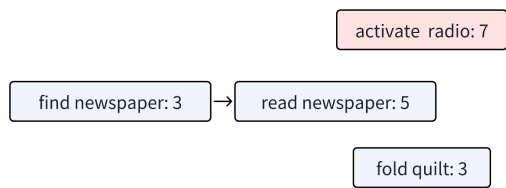
(f) The sixth task in the household activity scenario.



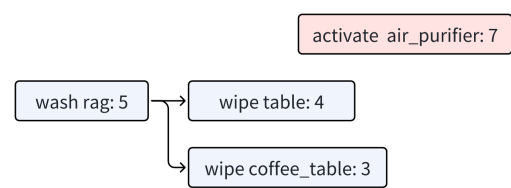
(g) The seventh task in the household activity scenario.



(h) The eighth task in the household activity scenario.



(i) The ninth task in the household activity scenario.



(j) The tenth task in the household activity scenario.

Figure 7: The action dependencies and durations for the ten tasks in the household activity scenario. Actions that occupy the agent, preventing them from doing anything else, are indicated with a blue background. In contrast, actions that do not occupy the agent, allowing for parallel tasks, are marked with a red background.

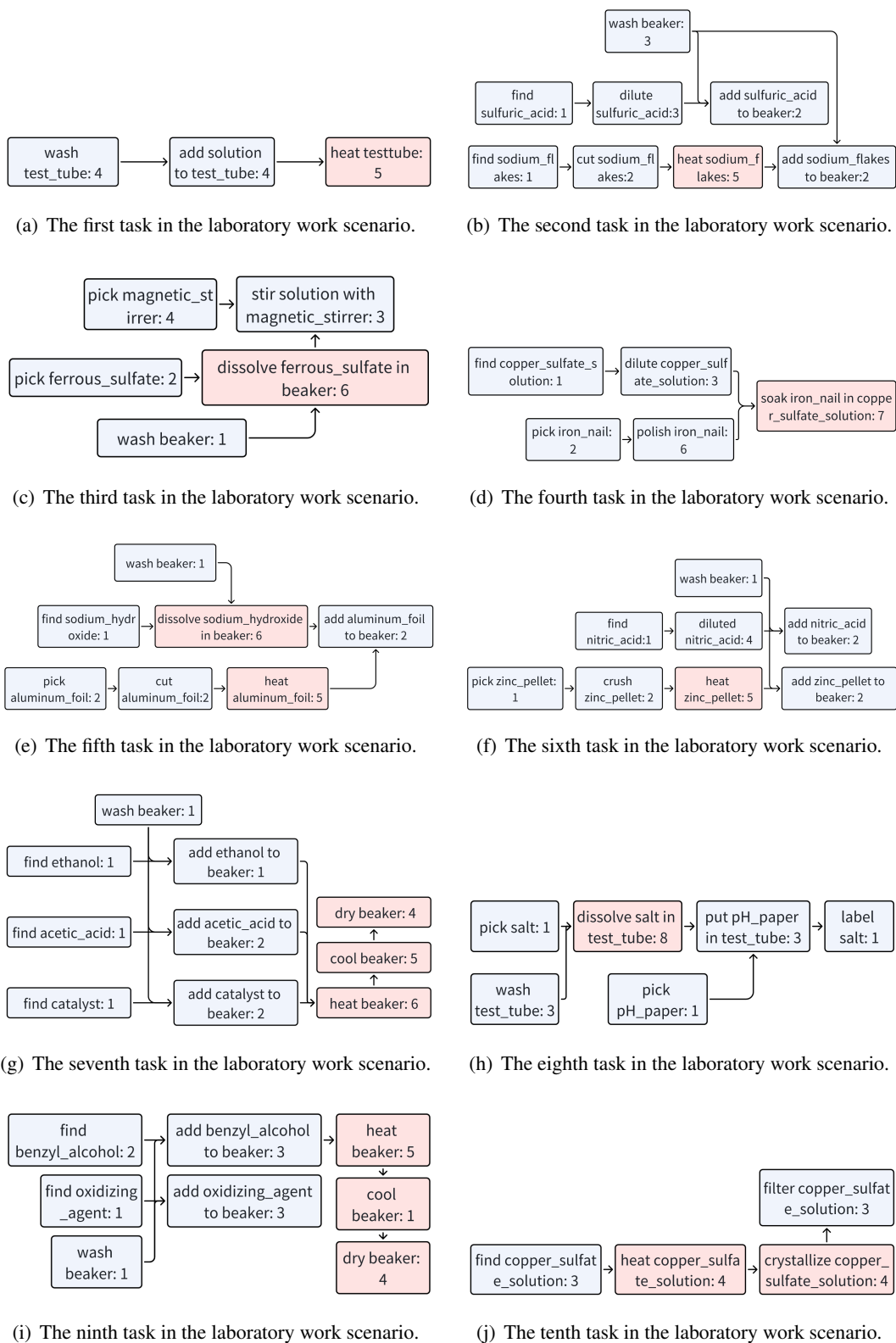


Figure 8: The action dependencies and durations for the ten tasks in the laboratory work scenario. Actions that occupy the agent, preventing them from doing anything else, are indicated with a blue background. In contrast, actions that do not occupy the agent, allowing for parallel tasks, are marked with a red background.

As an AI agent, your objective is to efficiently complete a series of tasks as described. You must adhere to the specific requirements and constraints of each task, including dependencies and timing. Efficiency is key; complete all tasks in the shortest possible time. I will provide instructions regarding actions and objects.

****Action Protocol**:**

- You can perform only one action at a time.
- After each observation from the environment, output an action based on that observation and the instructions.
- Actions fall into two categories:
 - Continuous Actions: Perform these actions until completion (e.g., "wash OBJ").
 - Autonomous Actions: These progress over time, allowing simultaneous tasks (e.g., "heat OBJ").
- Follow the "Valid Actions" format for your output (e.g., "wash cup").
- If no action is required, use "wait" to skip the current time.
- Output the action explicitly (e.g., "wash cup").
- Select object names (OBJ) from the list of Available Objects (e.g., use "rice" instead of "cooked rice").

****Task 1****

<Description>

- Prepare and bake a cheese and tomato pizza

<Valid Actions and Usages>

- pick OBJ: Pick the unpicked item.
- chop OBJ: Chop the whole item into sliced pieces.
- wash OBJ: Wash the dirty item to make clean.
- add OBJ1 to OBJ2: Add one item to the container.
- bake OBJ1 in OBJ2: Bake the raw item in the oven until it's roasted.
- wait: pass the current time without doing anything.

****Task 2****

<Description>

- Prepare chicken and potato stir-fry, which consists of fried chicken and fried potato.

<Valid Actions and Usages>

- pick OBJ: Pick the unpicked item.
- chop OBJ: Chop the whole item into sliced pieces.
- fry OBJ1 in OBJ2: Fry the raw item until it is fried to perfection.
- add OBJ1 to OBJ2: Add one item to the container.
- wash OBJ: Wash the dirty item to make clean.
- wait: pass the current time without doing anything.

****All Available Objects(OBJ)****

dish_1; dish_2; dough; cheese; tomato; oven; chicken; potato; fryer

****The Initial States of Objects****

dish_1: dirty; dish_2: dirty; dough: unpicked; cheese: unpicked; tomato: unpicked; oven: empty; chicken: unpicked; potato: unpicked; fryer: empty

Table 6: An example of # Task=2 scenario.

As an AI agent, your objective is to efficiently complete a series of tasks as described. You must adhere to the specific requirements and constraints of each task, including dependencies and timing. Efficiency is key; complete all tasks in the shortest possible time. I will provide instructions regarding actions and objects.

****Action Protocol**:**

- You can perform only one action at a time.
- After each observation from the environment, output an action based on that observation and the instructions.
- Actions fall into two categories:
 - Continuous Actions: Perform these actions until completion (e.g., "wash OBJ").
 - Autonomous Actions: These progress over time, allowing simultaneous tasks (e.g., "heat OBJ").
- Follow the "Valid Actions" format for your output (e.g., "wash cup").
- If no action is required, use "wait" to skip the current time.
- Output the action explicitly (e.g., "wash cup").
- Select object names (OBJ) from the list of Available Objects (e.g., use "rice" instead of "cooked rice").

****Task 1****

<Description>

- Prepare a garden bed for planting flowers by using sprinkling can filled with herbicide, hoeing, and weeding

<Valid Actions and Usages>

- add OBJ1 to OBJ2: Add one item to the container.
- weed_with OBJ: Weed with the item.
- hoe OBJ: Hoe the uncultivated item until it is cultivated and ready for planting.
- plant OBJ: Plant the uncultivated item until it is planted
- wait: pass the current time without doing anything.

****Task 2****

<Description>

- Iron a suit and store it properly

<Valid Actions and Usages>

- set_up OBJ: Set up the item that is not set yet until it is already set.
- put OBJ1 on OBJ2: Put the item on the right place.
- heat OBJ: Heat the cool item until it is hot.
- iron OBJ: Iron the wrinkled item until they are smooth.
- store OBJ: Store the unstored item\nwait: pass the current time without doing anything.

****Task 3****

<Description>

- Make a cup of coffee

<Valid Actions and Usages>

- add OBJ1 to OBJ2: Add one item to the container.
- activate OBJ: Activate the inactive device to turn it active.
- wash OBJ: Wash the dirty item to make clean.
- pour OBJ1 into OBJ2: Pour the liquid in item into the empty container until it is full.
- wait: pass the current time without doing anything.

****All Available Objects(OBJ)****

sprinkling_can; herbicide; land; flower; ironing_board; suit; iron; coffee_beans; coffee_machine; water; cup

****The Initial States of Objects****

sprinkling_can: empty; herbicide: not added; land: uncultivated; flower: uncultivated; ironing_board: not set yet; suit: not put on right place; iron: cool; coffee_beans: not added; coffee_machine: empty; water: not added; cup: dirty

Table 7: An example of # Task=3 scenario.

```

<|Environment|>:
As an AI agent, your objective is to efficiently complete a series of tasks as
described. You must adhere to the specific requirements and constraints of each task
, including dependencies and timing. Efficiency is key; complete all tasks in the
shortest possible time. I will provide instructions regarding actions and objects.

**Action Protocol**
- You can perform only one action at a time.
- After each observation from the environment, output an action based on that
observation and the instructions.
- Actions fall into two categories:
- Continuous Actions: Perform these actions until completion (e.g., "wash OBJ").
- Autonomous Actions: These progress over time, allowing simultaneous tasks (e.g., "
heat OBJ").
- Follow the "Valid Actions" format for your output (e.g., "wash cup").
- If no action is required, use "wait" to skip the current time.
- Output the action explicitly (e.g., "wash cup").
- Select object names (OBJ) from the list of Available Objects (e.g., use "rice"
instead of "cooked rice").

**Task 1**
<Description>
- Make a dish of beef fried rice, which consists of cooked rice and fried beef.
<Valid Actions and Usages>
- wash OBJ: Wash the dirty item to make clean.
- pick OBJ: Pick the unpicked item.
- cook OBJ1 in OBJ2: Cook the raw item until it's cooked through.
- chop OBJ: Chop the whole item into sliced pieces.
- fry OBJ1 in OBJ2: Fry the raw item until it is fried to perfection.
- add OBJ1 to OBJ2: Add one item to the container.
- wait: pass the current time without doing anything.

**All Available Objects(OBJ)**
rice; beef; pot; fryer; dish

**The Initial States of Objects**
rice: unpicked; beef: unpicked; pot: empty; fryer: empty; dish: dirty

<|Environment|>:
In t=1, your action is:

<|Agent|>:
pick dish

<|Environment|>:
You are doing ``wash dish``, it will take 2 minutes.
rice is picked.

In t=3, your action is:

<|Agent|>:
cook rice in pot

<|Environment|>:
You are doing ``cook rice in pot``, it will take 4 minutes.

In t=4, your action is:

<|Agent|>:
wait

<|Environment|>:
You wait for one minute.

...

```

Table 8: An example of interaction in cooking scenario.

```

<|Environment|>:
As an AI agent, your objective is to efficiently complete a series of tasks as
described. You must adhere to the specific requirements and constraints of each task
, including dependencies and timing. Efficiency is key; complete all tasks in the
shortest possible time. I will provide instructions regarding actions and objects.

**Action Protocol**
- You can perform only one action at a time.
- After each observation from the environment, output an action based on that
observation and the instructions.
- Actions fall into two categories:
- Continuous Actions: Perform these actions until completion (e.g., "wash OBJ").
- Autonomous Actions: These progress over time, allowing simultaneous tasks (e.g., "
heat OBJ").
- Follow the "Valid Actions" format for your output (e.g., "wash cup").
- If no action is required, use "wait" to skip the current time.
- Output the action explicitly (e.g., "wash cup").
- Select object names (OBJ) from the list of Available Objects (e.g., use "rice"
instead of "cooked rice").

**Task 1**
<Description>
- Make a dish of beef fried rice, which consists of cooked rice and fried beef.
<Valid Actions and Usages>
- wash OBJ: Wash the dirty item to make clean.
- pick OBJ: Pick the unpicked item.
- cook OBJ1 in OBJ2: Cook the raw item until it's cooked through.
- chop OBJ: Chop the whole item into sliced pieces.
- fry OBJ1 in OBJ2: Fry the raw item until it is fried to perfection.
- add OBJ1 to OBJ2: Add one item to the container.
- wait: pass the current time without doing anything.

**All Available Objects(OBJ)**
rice; beef; pot; fryer; dish

**The Initial States of Objects**
rice: unpicked; beef: unpicked; pot: empty; fryer: empty; dish: dirty

Given the list of valid actions, available objects, and the task descriptions (goal
), please perform the following steps:
- Identify and list all of the necessary actions required to accomplish the task's
goal.
- For each action, determine and note the specific objects that are required.
- Assess and map out any dependencies between actions, indicating which actions must
precede others.
- Arrange the actions in a logical sequence that respects the dependencies and leads
efficiently towards completing the task.
- If any action has multiple dependencies, list them in order of priority based on
the task's constraints and goal.
- Present the final action sequence in a clear and ordered list, ensuring that the
progression of steps will achieve the task's objective.

The key to efficiency:
- When completing tasks, some actions are non-occupied actions(Type 2), meaning you
can perform other actions simultaneously.
- To maximize efficiency, adhere to the following principle: always start the non-
occupied action you anticipate will be the most time-consuming as early as possible.
- You should perform actions during idle times as much as possible to minimize the
time spent doing nothing.

...

```

Table 9: Prompt of self-plan method in cooking scenario.