Haste Makes Waste: Evaluating Planning Abilities of LLMs for Efficient and Feasible Multitasking with Time Constraints

Anonymous ACL submission

Abstract

While Large Language Model-based agents have demonstrated substantial progress in task completion, existing evaluation benchmarks tend to overemphasize single-task performance, with insufficient attention given to the crucial aspects of multitask planning and execution efficiency required in real-world scenarios. To bridge this gap, we present RECIPE2PLAN, a novel benchmark framework based on realworld cooking scenarios. Unlike conventional benchmarks, RECIPE2PLAN challenges agents to optimize cooking time through parallel task execution while respecting temporal dependencies between steps. The benchmark emphasizes the delicate balance between maximizing concurrent operations and adhering to critical timing constraints, where overly aggressive local parallelization may disrupt subsequent timesensitive steps, potentially compromising the entire cooking process. Extensive experiments with state-of-the-art models reveal challenges in maintaining this balance between efficiency and feasibility. The results highlight the need for improved temporal awareness and global multitasking capabilities in large language models. We will open-source our benchmark and code to the community.

1 Introduction

002

017

020

022

024

035

040

042

043

Large Language Models (LLMs) (OpenAI, 2023; Team and Google, 2023; Touvron et al., 2023; Qwen Team, 2024) have demonstrated the ability to plan and reason step by step (Wei et al., 2022). Leveraging this ability, LLM-based agents can automate complex real-world tasks (Yao et al., 2022b; Shinn et al., 2024; Sun et al., 2024).

The effectiveness of LLM-based agents is primarily evaluated based on the *feasibility* of their plans in the scenarios of web browsing (Yao et al., 2022a; Zhou et al., 2023), tool usage (Qin et al., 2023; Li et al., 2023), computer manipulation (Xie et al., 2024b; Gou et al., 2024) and agent navigation (Shridhar et al., 2020; Wang et al., 2022).

However, the ability to manage concurrent objectives remains an often overlooked yet crucial requirement in real-world applications, as exemplified by everyday scenarios where humans prepare multiple dishes simultaneously for a meal or conduct parallel laboratory experiments (Russell and Norvig, 2010; Zhang et al., 2024; Wu et al., 2024). Current planning benchmarks assume that models execute tasks by decomposing the overall goal into steps and achieving these subgoals sequentially, one at a time (Liu et al., 2023; Ma et al., 2024). Consequently, these datasets fail to account for the duration of an action and the potential for multitasking. The multitasking scenario proposes a different objective in addition to *feasibility*. It challenges the model to optimize the efficiency to reach multiple goals simultaneously.

044

045

046

047

051

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

081

084

Time constraints are often imposed between certain steps in the recipe of dishes or experiments, indicating specific actions must be performed within a particular time interval after the preceding step is completed. For instance, the pouring and dripping actions for pour-over coffee must be carried out in sequence without any delay, as recommended by professionals (Hoffmann, 2018). This property introduces a unique challenge for multitask planning apart from conventional benchmarks. The first plan in Figure 1 illustrates that if the agent prioritizes maximizing efficiency by rushing to multitask whenever it is idle, it may inadvertently violate future time constraints. Consequently, the agent must balance the need for efficiency with adherence to time constraints to achieve feasible multitask planning as shown in the second plan in Figure 1.

We propose a new benchmark RECIPE2PLAN based on real-world recipes and constraints to evaluate the multitasking abilities of agents. We highlight three main challenges as: (1) **Commonsense Reasoning**: The agent must identify idle periods from the recipe as opportunities for multitasking while recognizing action dependencies



Figure 1: A simplified demonstration of our benchmark RECIPE2PLAN. Actions will either occupy the agent or leave it idle. The four steps of brewing must be executed sequentially as time constraints. The goal for the agent is to plan multitasking to complete the recipes in the shortest time possible without violating any constraints. The first plan illustrates a scenario where the agent attempts always to keep the agent occupied for higher efficiency, resulting in violations of time constraints. The second plan maintains the balance between the efficiency and feasibility of the plan by leaving the agent idle on purpose to maintain the time constraints for all actions.

and physical constraints to construct feasible action sequences. (2) Dynamic Local Planning: As recipe states evolve based on executed actions, the agent must continuously determine executable actions at each timestep. Additionally, the agent must dynamically adapt its beliefs and revise the plan accordingly if its initial assumptions about properties or constraints do not align with real-world conditions. (3) Strategic Global Planning: The agent is required to allocate the use of physical objects and schedule actions on a timeline to enable efficient multitasking. It is crucial to avoid planning multitasking in a purely local and greedy manner, as this could lead to violations of time restrictions. It challenges the agent to maximize efficiency while maintaining feasibility from a global perspective.

Our benchmark provides a testbed for the efficiency of LLM-based agents, as they are approaching the upper limits of feasibility in current text-based agent benchmarks (Sun et al., 2024) and multitasking scenarios without time constraints (Table 3). By introducing time restrictions, our benchmark evaluates the planning abilities of agents to maintain a delicate balance between efficiency and feasibility, rather than simply maximizing efficiency in a greedy manner. RECIPE2PLAN aims to push the boundaries of current agent planning capabilities, making them more adept at handling complex real-world tasks. 113 While our benchmark is constructed with cooking 114

scenarios, its scope could notably extend to the design of embodied agents in real-world multitasking, such as automatic scientific discovery. The high-throughout biomedical experiments (Yang et al., 2021) share the same principles of our benchmark, requiring the agents to deliver a feasible plan that finishes all the goals efficiently. 115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

In this study, we experiment with various sizes of open-source models, such as Qwen2.5 (Qwen Team, 2024) and Llama3.1 (Dubey et al., 2024), as well as closed-source models, including Gemini-1.5-Pro (Team et al., 2024) and GPT-40 (OpenAI, 2023). Our experiments reveal that GPT-40 achieves the highest success rate of only 21.5% and the main failure source is the violation of time constraints. It suggests that current LLMs fail to deliver feasible plans while attempting efficient multitasking. We show that LLMs can deliver feasible plans if time constraints are absent. However, their efficiency still lags significantly behind a simple heuristic method (§4). We also indicate that GPT-40 can trade efficiency for success rate if focusing solely on feasibility (§5.1). Overall, we demonstrate that current LLMs struggle to balance efficiency and feasibility when multitasking with time constraints. We further analyze the commonsense reasoning, local planning, and global planning capabilities of LLMs. By isolating each ability, we identify global planning as the primary source of task failure and inefficient multitasking (§5.3).

Our contributions are as follows:

145

146

147

148

149

150

151

152

153

154

155

156

158

159

160

161

162

163

164

165

167

168

169

170

171

172

174

175

176

177

178

179

We highlight the importance of the ability of agents to plan multitasking as balancing efficiency and feasibility. Compared with existing works focusing on feasibility, our work offers a new perspective to evaluate the planning abilities of LLMs.
We construct a benchmark RECIPE2PLAN based on real-world recipes for multitask planning. It challenges the model to allocate the usage of physical objects and schedule the actions on the timeline to complete the recipes in the shortest time possible without violating time constraints.

• We evaluate open-source and closed-source models on our benchmarks. Our results show that LLMs struggle with planning multitasking under time constraints, resulting in a low success rate for the task. This highlights the need for further development in enhancing the temporal reasoning and global planning capabilities of LLM agents.

2 Related Work

Planning Benchmarks. To evaluate the planning abilities of LLM-based agents, researchers have proposed benchmarks across various domains such as web browsing (Yao et al., 2022a; Zhou et al., 2023; Deng et al., 2024), tool usage (Qin et al., 2023; Li et al., 2023), and computer manipulation (Xie et al., 2024b; Gou et al., 2024). These benchmarks assess an agent's ability to execute a sequence of actions to achieve a general goal in a partially observable environment (Liu et al., 2023; Ma et al., 2024). However, these environments do not account for the duration of each action. Additionally, they evaluate planning abilities based solely on feasibility, without comparing the efficiency of task completion between different agents.

Scheduling Benchmarks. Apart from the typi-180 cal planning task in which the agent interacts with 181 a partially observable environment without prior 182 knowledge of how to achieve the goal, the schedul-183 ing task provides the agent with a complete description of the task. The objective is to deliver 185 an action sequence from a small set of fixed actions to meet the given objectives (Pinedo and Hadavi, 1992; Smith et al., 2000; Valmeekam et al., 189 2024). Graph coloring (Stechly et al., 2024) investigates whether LLMs can self-critique their answers 190 for violations of scheduling constraints. NATU-191 RALPLAN assesses scheduling abilities in contexts such as trip planning, meeting planning, and calen-193

dar scheduling. TravelPlan (Xie et al., 2024a) deals with more complex commonsense constraints and strict restrictions. TIMEARENA primarily evaluates the multitasking capabilities of LLMs without time constraints. Our benchmark focuses on assessing the ability to balance efficiency and feasibility during multitasking with time constraints. Please refer to Table 1 for a detailed comparison. 194

195

196

197

198

199

200

202

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

229

230

231

232

233

234

235

236

238

Planning Methods. Different methods use feedback and instructions in various ways. Open-loop methods such as Chain-of-Thought (Wei et al., 2022), least-to-most (Zhou et al., 2022) and planand-solve (Wang et al., 2023) plan the action sequence without any feedback from the environment. This type of method is vulnerable to the hallucination of execution constraints and environment dynamics. Closed-loop methods such as Re-Act (Yao et al., 2022b) and Reflexion (Shinn et al., 2024) only refine local actions, which might result in global failure due to time constraints. Ada-Planner (Sun et al., 2024) refines the entire plan based on environmental feedback and past failures. LLM-Modulo framework (Kambhampati, 2024) surpasses existing baselines for complex scheduling task (Gundawar et al., 2024b) by regenerating the entire plan with detailed critiques.

3 RECIPE2PLAN

RECIPE2PLAN evaluates the planning ability of LLMs for efficient and feasible multitasking under constraints. Specifically, we provide the model with multiple goals that can be achieved by following recipes. Each recipe A is represented as a linear sequence of actions $A = (a_0, a_1, ..., a_n)$, with each action assigned a specific execution time t_n . The task is to plan the action sequence to complete all goals in the shortest time possible, adhering to the properties and constraints detailed in §3.1 and §3.2. The statistics of our benchmark are presented in Table 2. RECIPE2PLAN challenges the model to apply commonsense reasoning to infer any unwritten constraints from the recipe, including action concurrency, action dependencies, and resource limitations while planning the action sequence to minimize overall execution time.

3.1 Properties of Actions

Action Duration.This refers to the time required239for an agent to complete a specific action. For the
coffee recipes illustrated in Figure 1, the duration240

Benchmark	Commonsense Reasoning	Temporal Planning	Multitask Planning	Time Constraints	Harmonized Planning
Graph Coloring (Stechly et al., 2024)	×	×	×	×	×
NATURAL PLAN (Zheng et al., 2024)	×	1	×	×	×
TravelPlanner (Xie et al., 2024a)	1	1	×	×	×
TIMEARENA (Zhang et al., 2024)	✓	✓	✓	×	×
RECIPE2PLAN	✓	✓	✓	✓	1

Table 1: Comparison with existing *scheduling* benchmarks. Two unique properties distinguish our benchmark: (1) *Complex Planning*: Temporal multitasking recipes simultaneously require the blend of different abilities: commonsense reasoning, dynamic local planning and strategic global planning; (2) *Harmonized Planning*: The agent must balance efficiency and feasibility. While local-optimal planning for maximum efficiency might be desirable, it must not compromise time constraints, as specific actions need to be performed within particular time intervals following the preceding steps.

of actions such as pouring and dripping is fixed, and any deviation from these durations can result in spoiled flavor. Following this principle, each action in our benchmark is annotated with a specific duration. The recipe in the dataset explicitly states this duration, allowing the agent to accurately schedule the timeline.

242

243

244

246

247

251

254

258

260

261

263

264

265

269

Action Concurrency. Continuous actions, such as *pour water*, require the active involvement of the agent while the action is in progress. In contrast, autonomous actions, like *boil water*, do not require the agent's continuous attention, allowing the agent to remain idle and free to perform other tasks concurrently. Identifying autonomous actions and executing them simultaneously with other actions is the key to efficient multitasking.

Execution Interruptibility. We introduce this property by allowing actions, such as *cut onions*, to be completed by dividing the execution into two or more separate time intervals. Previous benchmarks generally do not consider duration (Ma et al., 2024) or treat the execution of actions as a certain time frame (Zhang et al., 2024). Execution interruptibility enables the model to generate a more fine-grained plan by allocating the execution for one action across different intervals, thus enhancing the efficiency and flexibility of multitasking.

3.2 Multitasking Constraints

270Action Dependencies.The dependent relation-271ships between actions are generally not explicitly272stated in the recipe. Although the actions in a recipe273are often presented in a linear sequence, the action274dependencies might form a graph structure. For275example, as illustrated in Figure 1, step 3 *pour wa-276ter depends on step 1 <i>boil water* and step 2 grind

coffee, but steps 1 and 2 can be performed independently of each other. If the agent discovers that the dependency does not align with its initial understanding during execution, it must dynamically adapt the plan based on the current status.

277

278

279

281

283

284

285

287

288

289

290

291

292

293

295

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

Resource Limitations. During planning, the agent must recognize whether an object is occupied at the current time and when it will be available again. Different recipes may require different physical objects and conditions. As shown in Figure 1, the pouring process for coffee requires water at a specific temperature, so the agent can boil water for both recipes simultaneously to speed up the process. However, if different recipes require water or an oven at different temperatures, the agent must sequentially prepare the object for each recipe based on when it becomes available. This property necessitates that the agent plan globally, scheduling the use of different objects while considering the duration of actions and specific condition requirements.

Time Constraints. This property is crucial for feasible multitasking of professional coffee preparation (Hoffmann, 2018) and delicate biomedical experiments (Itoh et al., 2021), where specific actions must be executed within a precise time interval following a preceding action. Failure to adhere to these time constraints may alter the flavor of a dish, or even cause the entire recipe or experiment to fail. This realistic property imposes a significant challenge on multitask planning. As depicted in Figure 1, the agent can not simply follow a greedy manner that prioritizes immediate actions without considering the broader temporal constraints. The incorporation of time constraints in the planning process ensures that the agent must strategically balance multitasking efficiency with feasibility.

313 3.3 Dataset Construction

Recipe Annotation. We collect and clean recipes 314 with annotated dependent relations from a website 315 for cooking¹. We ask three annotators to label the 317 properties and constraints following the pipeline in Appendix A.1. The average kappa scores among annotators are 0.78 for action concurrency, 0.50 for action interruptibility, 0.66 for time constraints, and 0.86 for resource limitations. Based on these 321 results, we explicitly list action interruptibility and time constraints in the recipes, while keeping 323 action concurrency and resource limitations as 324 implicit properties that agents need to identify through commonsense reasoning. Finally, we prompt GPT-40 (Hurst et al., 2024) to annotate action durations and time intervals for restrictions, 328 and all annotators have reviewed these annotations to ensure their reasonableness.

Combine Recipes for Multitasking. 331 We carefully select recipe combinations to evaluate planning abilities for efficient multitasking. To keep the 333 action space and context length manageable, we only combine two recipes at a time. We then adapt 335 a heuristic algorithm from Zhang et al. (2024) to plan action sequences for multitasking. Multitasking efficiency for each sequence is computed according to Equation 1. Instances are chosen for the 339 benchmark based on the following criteria: (1) Opportunities for multitasking: We include instances 341 with multitasking efficiency higher than 80% when planning without time constraints, indicating significant potential opportunities for multitasking. (2) Balance of efficiency and feasibility: We select in-345 stances in which the multitasking efficiency drops 346 347 when time constraints are considered, suggesting that an efficient greedy planning strategy would 348 likely violate these time constraints.

3.4 Environment

351

354

355

We implement an environment to provide feedback to the agent. The agent can choose to perform one action for a specified duration at a given time. If the agent determines that no action can be performed at the moment, it can choose the time for its next planned action. The environment then receives the action and checks for any constraint violations. If a constraint is violated, the environment will specify the type of violation. If the action is permissible, feedback from the environment includes the status

Recipe Statistics					
# Recipes	29				
Avg. Actions per Recipe	13.1				
Avg. Autonomous Actions per Recipe	3.4				
Avg. Interruptible Actions per Recipe	3.9				
Avg. time constraint per Recipe	3.1				
Avg. Duration per Action (min)	5.7				
Avg. Restriction Interval (min)	2.7				
Multitasking Statistics					
# Instances	65				
Avg. Executable Action per Step	3.1				
Avg. Efficiency w/o time constraint (%)	92.3				
Avg. Efficiency w/ time constraint (%)	80.1				

Table 2: Statistics of recipes and multitasking instances in our RECIPE2PLAN benchmark. The agent can choose any timestamp for the next action, expanding the search space beyond the number of executable actions solely.

of physical objects, completed actions, and ongoing autonomous actions. We present examples of observations and feedback in Appendix C.2. The agent can use this feedback to revise its global plan and decide on the next action. 361

362

363

364

366

367

368

369

370

371

372

373

374

375

376

377

378

379

381

382

383

384

387

389

390

4 **Experiments**

4.1 Baselines

Models. We evaluate several models, including the open-source Llama-3.1 with parameter sizes of 8B and 70B (Dubey et al., 2024), and Qwen2.5 with parameter sizes of 7B, 32B and 72B (Qwen Team, 2024). Additionally, we assess the closed-source models, GPT-4o-mini, GPT-4o (Hurst et al., 2024) and Gemini-1.5-Pro (Team et al., 2024). The versions of the models are detailed in Appendix C.3.

Methods. We begin by prompting the model to identify any unwritten properties and constraints from each recipe. These identified elements are then concatenated with the original description. Next, we employ a *ReAct*-style prompting method (Yao et al., 2022b) on the models to plan the action sequence. To evaluate the planning abilities and mitigate the cascading errors from commonsense reasoning, we also experiment with an oracle setting *ReAct* + *Oracle* that replaces the identified constraints with the gold annotations.

Constraint Setting. We evaluate the agent under *without time constraints* and *with time constraints* settings to study the impact of time constraints on the feasibility and efficiency of multitasking agents.

¹www.instructables.com

Model	w/o Time Constraints				w/ Time Constraints			
Mouch	Success	Progress	R-Efficiency	S×E	Success	Progress	R-Efficiency	S×E
				Re	Act			
Open-Source Models								
Qwen2.5-7B	1.5	26.5	73.6	1.3	0.0	22.4	78.0	0.0
Llama-3.1-8B	0.0	9.7	47.8	0.0	0.0	10.3	44.9	0.0
Qwen2.5-32B	80.0	96.7	60.5	51.0	15.4	57.4	86.4	9.9
Llama-3.1-70B	72.3	88.8	65.9	50.5	13.8	55.5	79.3	9.8
Qwen2.5-72B	72.3	91.2	72.3	52.3	7.7	54.9	91.6	5.8
Closed-Source Models								
GPT-4o-mini	3.1	51.8	53.6	2.0	1.5	36.0	63.4	0.4
Gemini-1.5-Pro	20.0	66.4	71.7	14.3	3.1	47.5	74.4	1.8
GPT-40	90.8	99.1	81.0	75.3	21.5	64.0	102.6	19.5
	ReAct + Oracle							
Open-Source Models								
Qwen2.5-7B	0.0	28.1	68.9	0.0	0.0	24.1	61.8	0.0
Llama-3.1-8B	0.0	10.8	60.6	0.0	0.0	10.0	49.1	0.0
Qwen2.5-32B	80.0	96.7	55.9	49.8	10.8	57.0	85.7	8.6
Llama-3.1-70B	73.8	89.0	61.6	49.9	6.2	52.9	78.9	2.5
Qwen2.5-72B	72.3	90.1	71.2	53.2	7.7	52.6	98.3	6.0
Closed-Source Models								
GPT-4o-mini	10.8	55.8	54.0	3.9	1.5	35.0	64.8	1.5
Gemini-1.5-Pro	16.9	63.6	68.4	13.7	7.7	49.1	76.8	3.8
GPT-40	95.4	99.4	78.2	76.3	27.7	60.6	104.0	20.3
Heuristics	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0

Table 3: Results of *ReAct* and *ReAct* + *Oracle* experiments on RECIPE2PLAN. We report average percentage of success rate (**Success**), progress rate (**Progress**), relative multitask efficiency (**R-Efficiency**) and multitasking score ($S \times E$). **Bold** denotes the best performance and <u>underline</u> denotes the second-best performance.

4.2 Evaluation Metrics

391

394

396

400

401

402

403

404

405

406

407

408

409

410

411

412

413

Success Rate. It measures the *feasibility* of the plans exclusively by evaluating whether the agent can deliver a plan that successfully completes all recipes. The agent might fail due to a violation of time constraints, reaching a maximum of execution errors, or being stuck in an endless dead loop.

Progress Rate. This metric measures the proportion of successfully executed actions in the recipes.It evaluates the *feasibility* of the planning process with a more fine-grained perspective.

Multitasking Efficiency. Opportunities to minimize execution time arise from autonomous actions within the recipes. We first calculate the time saved T_{save} through multitasking relative to the total duration of all completed actions. We measure the multitasking *efficiency* for successfully executed actions as the proportion of saved time and the cumulative duration of the executed autonomous actions T_{auto} similar to Zhang et al. (2024).

$$\text{Efficiency}_{agent} = \frac{T_{save}}{T_{auto}} \tag{1}$$

We observe that the progress rate potentially influences this metric. To address this, we perform a calibration based on the efficiency of plans derived from the heuristic baseline (Efficiency_{ref}), which allows us to compute the *relative multitasking efficiency*, with the rationale detailed in Appendix B.

$$R-Efficiency = \frac{Efficiency_{agent}}{Efficiency_{ref}}$$
(2)

Multitasking Score. We use this metric to present the overall *efficiency* and *feasibility* of the plans simultaneously. The score is computed as:

$$M. Score = \begin{cases} R-Efficiency & Success = 1\\ 0 & otherwise \end{cases}$$
(2)

$$(\mathbf{J})$$

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

If the plan is successfully executed (Success = 1), the multitasking score equals the R-Efficiency of the plan; otherwise, the score is 0, indicating failure. The rationale is that the agent should prioritize ensuring the successful completion of the recipes before aiming to achieve higher efficiency in task execution. The overall score is computed as the average of multitasking scores for each instance.

4.3 Main Results

LLMs can plan feasible multitasking in the absence of time constraints, but efficiency needs



Figure 2: Results of GPT-40 planning with different priority. *Balanced Priority*: Blend feasibility and efficiency as in §4. *Feasibility Priority*: Only focus on feasibility without considering efficiency.

improvement. In the *ReAct* setting, GPT-40 delivers 90.8% feasible plans for multitasking, surpassing other tested models by a large margin. This demonstrates its ability to revise beliefs about unwritten properties and constraints and to correct its actions to complete tasks. GPT-40 achieves a multitasking efficiency of 78.2%, indicating that there is still room for improvement in multitasking efficiency with LLMs.

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

LLMs face challenges in balancing efficiency and feasibility for multitasking under time constraints. The success rates and completion ratios of all models decrease significantly when multitasking with time constraints. GPT-40 only achieves a success rate of 21.5%. Interestingly, GPT-40 achieves a relative efficiency of 102.6%, which is higher than the heuristic baseline. This indicates a tendency to prioritize high efficiency during local planning. But the agent fails to maintain feasibility for time constraints from a global perspective while managing multitasking efficiency.

Open-source models exhibit diverse prefer-455 ences for feasibility and efficiency. Surpris-456 ingly, Qwen2.5-32B outperforms Qwen2.5-72B 457 and Llama-3.1-70B in success rate and progress 458 rate, whereas the latter two exhibit higher efficiency. 459 This highlights the varying preference towards fea-460 sibility and efficiency across different open-source 461 models, even without time constraints. A similar 462 trend is observed when multitasking under time 463 constraints: Qwen2.5-32B achieves the second-464 highest success rate among all tested models, while 465 Qwen2.5-72B attains the second-highest relative 466 multitasking efficiency of 91.6%. 467

468 Commonsense reasoning is not the bottleneck
469 for feasible and efficient multitasking. The mod470 els generally achieve F1 scores higher than 70% for



Figure 3: Results of prompting GPT-40 under *Oracle* setting: gold constraints, and *Oracle* + *Hint* setting: gold constraints and executable actions at each step.

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

506

commonsense reasoning as detailed in Appendix D. We investigate the impact of unidentified properties and constraints in *ReAct* + *Oracle* setting. While the success rate improves by 4.8% for GPT-40 in the w/o time constraints setting, the relative efficiency for GPT-40 decreases by 2.8%, indicating that the model struggles to formulate an efficient multitasking plan even with oracle constraints. In the w/ time constraints setting, we observe the success rate increases slightly for GPT-40 from 21.5% to 27.7%. This suggests misidentified constraints are not the primary errors leading to time constraint violations. Overall, providing oracle constraints does not result in a significant performance increase. Therefore, the main bottleneck for multitasking is the planning abilities of the agents, as we elaborate in the next section.

5 Analysis

5.1 Multitasking with Different Priority

Our experiments (§4) indicate that LLMs struggle to complete the recipes with balanced priority of feasibility and efficiency. Therefore, we evaluate if LLMs can focus solely on the *feasibility priority*, ensuring that recipes are completed without violating any constraints as detailed in Appendix C.4. LLMs can ensure more task completion by trading efficiency for feasibility. The results in Figure 2 show that under the feasibility priority setting, the success rate significantly increases from 27.7% to 49.2%, and the progress rate increases from 60.6% to 85.7%. This indicates that focusing on feasibility allows more recipes to be completed and more steps to be executed within the given time constraints. It further underscores the importance of enhancing the planning abilities of LLM agents to balance feasibility and efficiency.



Figure 4: Analysis of the distribution of invalid actions and failure source of ReAct + Oracle agents planning with time constraints. GPT-40 + Hint: We add all the executable actions in the prompt to help the agent choose the next action during dynamic local planning.

5.2 Error Analysis

507

510

511

512

514

515

516

517

518

519

520

523

524

530

531

533

535

536

In this section, we take a closer look at the dynamic local planning abilities of the agents by examining the distribution of valid and invalid actions. Invalid actions are categorized into: *action mismatch* (executing non-existing actions and repeating finished actions), and violations of other properties and constraints. We distinguish *time constraint* as a source of failure separate from invalid actions, along with other types of failure in Figure 4.

Action dependencies are the primary source of invalid actions. As illustrated in Figure 4, models with a high success rate under the *w/o time constraints* setting consistently achieve a valid action ratio above 80% under time constraints. Despite all constraints being explicitly presented in the prompt during planning, the agent frequently violates these constraints, particularly those related to action dependencies. Upon examining the reasoning traces, we observe that LLMs often breach action dependencies constraints while attempting to optimize multitasking, consequently neglecting feasibility.

Time constraints are the main sources of task failure. For the failure source of planning with time constraints, open-source models Llama-3.1 and Qwen2.5 may still get stuck in loop or exceed maximum revisions for about 10% of the instances. But the main source for the failure of planning is due to time constraints, even GPT-40 fails to maintain time constraints in 70% of the cases.

5.3 Planning with Hints of Executable Actions

As LLMs can not handle action dependencies well
while planning for efficient multitasking, we further add the executable actions for each step in the
prompt. This allows us to evaluate global planning
abilities directly.

LLMs lack global planning ability for efficient planning and maintaining time constraints. The success rate and multitask score show minimal improvement in both settings, as illustrated in Figure 3. This indicates that while agents can select a valid action for local planning (Figure 4), they fail to consider the impact of their actions on the overall feasibility and efficiency from a global perspective. Table 4 in the Appendix demonstrates a case where GPT-40 fails to estimate the priority of autonomous actions and leaves the agent idle during the execution of the last two actions. Table 5 in the Appendix provides an example where GPT-40 rushes to *heat up oil* at the beginning of the plan and executes this action concurrently with others to maximize efficiency. This plan overlooks ingredient preparation and results in the oil heated for an extended period. It does not only violate the time constraint but also risks catching fire.

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

6 Conclusions

Our paper introduces the RECIPE2PLAN benchmark, which evaluates the feasible and efficient multitasking abilities of existing LLMs. This benchmark pushes the limits of current agent planning capabilities beyond mere task completion to include the optimization of time and resource management. Our experiments reveal that while strong models like GPT-40 can generate feasible plans without time constraints, their performance decreases sharply when time constraints are imposed. This highlights a significant gap between current capabilities and the requirements for feasible and efficient multitasking. Our analysis identifies global planning as the primary area needing improvement, paving the way for future work to focus on enhancing temporal reasoning and strategic planning.

8

Limitations

579

581

582

583

584

585

586

587

588

590

594

598

604

606

610

611

612

613

614

615

616

617

620

622

624

625

629

While multitasking is a practical application for LLM agents, our text-based environment does not fully capture the complexities of real-world cooking and experimentation. Our agent does not engage in physical exploration or interact with objects in the real world, focusing solely on the temporal planning aspects of multitasking. In our setting, the agent is assumed to perform every action without delay or failure. Introducing scenarios where the agent must search for ingredients in a kitchen or lab similar to Shridhar et al. (2020) and Wang et al. (2022) would present a more realistic and challenging environment. We plan to implement such a realistic environment in future work.

The metric we use to evaluate efficiency by computing the speed of completion may be biased by the progress rate. To address this, we introduce a relative multitasking efficiency metric to calibrate our evaluation. However, the solution provided by our heuristic baseline does not guarantee the optimal plan for the task. The search space is complex because the model can choose to execute actions at arbitrary time stamps and split actions into arbitrary time intervals, making it beyond the scope of classical scheduling algorithms with time constraints (Itoh et al., 2021). While existing scheduling algorithms may take a long time to execute, our heuristic algorithm quickly identifies a feasible and efficient plan, though it may be suboptimal. We believe this heuristic can still serve as a valuable baseline for evaluating the multitasking abilities of agents. For future work, we plan to explore scheduling algorithms that can better handle the complexities of multitasking with time constraints.

References

- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. 2024.
 Mind2web: Towards a generalist agent for the web. Advances in Neural Information Processing Systems, 36.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Boyu Gou, Ruohan Wang, Boyuan Zheng, Yanan Xie, Cheng Chang, Yiheng Shu, Huan Sun, and Yu Su. 2024. Navigating the digital world as humans do: Universal visual grounding for gui agents. arXiv preprint arXiv:2410.05243.

Atharva Gundawar, Karthik Valmeekam, Mudit Verma, and Subbarao Kambhampati. 2024a. Robust planning with compound llm architectures: An llmmodulo approach. *Preprint*, arXiv:2411.14484. 630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

- Atharva Gundawar, Mudit Verma, Lin Guan, Karthik Valmeekam, Siddhant Bhambri, and Subbarao Kambhampati. 2024b. Robust planning with llm-modulo framework: Case study in travel planning. *arXiv preprint arXiv:2405.20625*.
- James Hoffmann. 2018. *The World Atlas of Coffee: From beans to brewing-coffees explored, explained and enjoyed.* Hachette UK.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-40 system card. *arXiv preprint arXiv:2410.21276*.
- Takeshi D Itoh, Takaaki Horinouchi, Hiroki Uchida, Koichi Takahashi, and Haruka Ozaki. 2021. Optimal scheduling for laboratory automation of life science experiments with time constraints. *SLAS TECHNOLOGY: Translating Life Sciences Innovation*, 26(6):650–659.
- Subbarao Kambhampati. 2024. Can large language models reason and plan? *Annals of the New York Academy of Sciences*, 1534(1):15–18.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pages 611–626.
- Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. 2023. Api-bank: A comprehensive benchmark for tool-augmented llms. In *Proceedings* of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 3102–3116.
- 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. *Preprint*, arXiv:2308.03688.
- Chang Ma, Junlei Zhang, Zhihao Zhu, Cheng Yang, Yujiu Yang, Yaohui Jin, Zhenzhong Lan, Lingpeng Kong, and Junxian He. 2024. Agentboard: An analytical evaluation board of multi-turn llm agents. *arXiv preprint arXiv:2401.13178*.
- OpenAI. 2023. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

794

795

690

691

- 733 734 736
- 737 739

- Liang-Ming Pan, Jingjing Chen, Jianlong Wu, Shaoteng Liu, Chong-Wah Ngo, Min-Yen Kan, Yugang Jiang, and Tat-Seng Chua. 2020. Multi-modal cooking workflow construction for food recipes. In Proceedings of the 28th ACM International Conference on Multimedia, pages 1132–1141.
- Michael Pinedo and Khosrow Hadavi. 1992. Scheduling: theory, algorithms and systems development. In **Operations Research Proceedings 1991: Papers of** the 20th Annual Meeting/Vorträge der 20. Jahrestagung, pages 35-42. Springer.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. arXiv preprint arXiv:2307.16789.
- Qwen Team. 2024. Qwen2.5: A party of foundation models.
- Stuart J Russell and Peter Norvig. 2010. Artificial intel*ligence a modern approach*. London.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2024. Reflexion: Language agents with verbal reinforcement learning. Advances in Neural Information Processing Systems, 36.
- 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.
- David E Smith, Jeremy Frank, and Ari K Jónsson. 2000. Bridging the gap between planning and scheduling. The Knowledge Engineering Review, 15(1):47–83.
- Kaya Stechly, Karthik Valmeekam, and Subbarao Kambhampati. 2024. On the self-verification limitations of large language models on reasoning and planning tasks. arXiv preprint arXiv:2402.08115.
- Haotian Sun, Yuchen Zhuang, Lingkai Kong, Bo Dai, and Chao Zhang. 2024. Adaplanner: Adaptive planning from feedback with language models. Advances in Neural Information Processing Systems, 36.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint arXiv:2403.05530.
- Gemini Team and Google. 2023. Gemini: A family of highly capable multimodal models. Preprint, arXiv:2312.11805.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288.

- Karthik Valmeekam, Kaya Stechly, Atharva Gundawar, and Subbarao Kambhampati. 2024. Planning in strawberry fields: Evaluating and improving the planning and scheduling capabilities of lrm o1. arXiv *preprint arXiv:2410.02162.*
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023. Planand-solve prompting: Improving zero-shot chain-ofthought reasoning by large language models. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2609–2634, Toronto, Canada. 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.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837.
- Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao Yu, and Lingpeng Kong. 2024. Os-copilot: Towards generalist computer agents with self-improvement. arXiv preprint arXiv:2402.07456.
- Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. 2024a. Travelplanner: A benchmark for real-world planning with language agents. In Forty-first International Conference on Machine Learning.
- Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. 2024b. Osworld: Benchmarking multimodal agents for openended tasks in real computer environments. arXiv preprint arXiv:2404.07972.
- Liangliang Yang, Sara Pijuan-Galito, Hoon Suk Rho, Aliaksei S Vasilevich, Aysegul Dede Eren, Lu Ge, Pamela Habibovic, Morgan R Alexander, Jan de Boer, Aurelie Carlier, et al. 2021. High-throughput methods in the discovery and study of biomaterials and materiobiology. Chemical reviews, 121(8):4561-4677.
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022a. Webshop: Towards scalable real-world web interaction with grounded language agents. Advances in Neural Information Processing Systems, 35:20744-20757.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022b. React: Synergizing reasoning and acting in language models. arXiv preprint arXiv:2210.03629.

- Yikai Zhang, Siyu Yuan, Caiyu Hu, Kyle Richardson, Yanghua Xiao, and Jiangjie Chen. 2024. Timearena: Shaping efficient multitasking language agents in a time-aware simulation. *arXiv preprint arXiv:2402.05733*.
- Huaixiu Steven Zheng, Swaroop Mishra, Hugh Zhang, Xinyun Chen, Minmin Chen, Azade Nova, Le Hou, Heng-Tze Cheng, Quoc V Le, Ed H Chi, et al. 2024.
 Natural plan: Benchmarking llms on natural language planning. *arXiv preprint arXiv:2406.04520*.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. 2022. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*.
- Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. 2023. Webarena: A realistic web environment for building autonomous agents. *arXiv e-prints*, pages arXiv–2307.

A Details for Dataset Construction

A.1 Recipe Collection

796

797

800

805

810

811

812

813

814 815

816

817

819

822

824

826

827

834

838

844

We collect cooking recipes from existing benchmark MM-Res (Pan et al., 2020). MM-Res contains 9,850 recipes from cooking websites and has annotated the dependent relationship between actions in the recipe. To curate cases from the MM-Res dataset for the purpose of our benchmark. We sample recipes that involve using a microwave, an oven or a stove and disregard those with more than 30 actions. We remove actions that is a non-cooking steps, such as introductory phrases like today we want orka. Next, we ensure there are no temporal inconsistencies between steps. Optional statements are either removed or converted into mandatory steps. For example, You can use a spoon to get all the contents if needed is revised to exclude if needed. Actions are split for clarity if needed. for instance, boil water and pour water into a cup becomes boil water and pour water into a cup, to separate the autonomous and continuous actions. Conversely, steps that describe sequential actions in separate sentences are merged; for example, use water to strain and strain until the juices are gone are combined into a single step. The dependent relationships of the revised actions are adjusted accordingly.

A.2 Details for Recipe Annotation

We recruit three graduate students with cooking experience to annotate the action properties and constraints in the recipes following the guidelines 847 in Table 9. Each student identifies whether actions 848 were autonomous or continuous, marks actions as 849 interruptible or non-interruptible, and specifies any 850 physical or time constraints associated with each 851 action. To ensure consistency and accuracy, anno-852 tations were cross-verified among the annotators, 853 with discrepancies resolved through discussion. 854

855

856

857

858

859

860

861

862

863

864

865

866

868

869

870

871

872

873

874

B Relative Multitask Efficiency



Figure 5: Demonstration of relative efficiency. The efficiency is affected by the progress rate and we use relative efficiency (relative efficiency) to calibrate the metric.

We find that the efficiency metric is influenced by the progress rate. As shown in Figure 5, if the agent aborts the interaction midway, it will achieve an efficiency of 100%. Conversely, an agent that completes all tasks in the optimal manner only obtains an efficiency of 87.5%. To address this discrepancy, we use the plan from the heuristic baseline as a reference to calibrate the efficiency. We compute the efficiency for the part of the plan that achieves the same progress rate as the agent. The relative efficiency is computed as

$$R-Efficiency = \frac{Efficiency_{agent}}{Efficiency_{ref}}$$
867

C Implementation Details

C.1 Heuristic Baseline Algorithm

We adapt the heuristic method from Zhang et al. (2024) to search for an efficient plan that is feasible. The details of the algorithm is presented in Algorithm 1.

C.2 Environment

We implement an environment to provide feedback875to the agent. The examples of feedback are listed876in Table 14. If the action fits all the constraints, the877agent receives a message of the successful execu-878tion. And an observation of the action concurrency879of the executed action, current timestamp, status880

of physical objects and the executing autonomous actions. If the action can not be executed, the environment will return an error message and detailed feedback about the violated constraint. We also provide a hint about executable actions to evaluate the global planning abilities of the agent solely in §5.3. During the interaction with the agent, the maximum number of revisions is 10. Exceeding this number will be considered as task failure. And we abort the multitasking process if the agent attempts to execute the same action three times or violates any time constraints.

C.3 Model Details

881

882

883

887

894

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

We use the Instruct version for all sizes of Qwen2.5² and Llama-3.1³ models in our study. We use vllm (Kwon et al., 2023) to deploy Qwen2.5-7B, Llama-3.1-8B and Qwen-2.5-32B on a single A800 GPU, and Llama-3.1-70B and Qwen2.5-72B on four A800 GPUs. We use the gpt-4o-2024-08-06 for GPT-40⁴, gpt-40-mini-2024-07-18 for GPT-40-mini⁵, Gemini-1.5-Pro-002 (2024-09-24) for Gemini-1.5-Pro⁶.

C.4 Planning Methods

Commonsense Reasoning We prompt the model with the same guidelines in §A.1 and one example to generate the beliefs of action concurrency, action dependency and resource limitations. The temperature are set as 0 for all models. The max tokens for generation are set as 128.

Open-Loop Planning We evaluate the open-loop planning methods to determine if current LLMs can plan action sequences without interacting with the environment. Given the complexity of our task, we implement a *Plan-and-Solve* baseline. In this approach, the model generates beliefs about unwritten properties and constraints through commonsense reasoning and creates a coarse-grained plan to perform actions simultaneously, aiming to reduce total execution time. Finally, the agent writes a fine-grained action sequence following one example for execution as shown in Table 11. The temperature is also set as 0 and the maximum generation tokens is 2048.

²https://huggingface.co/Qwen

⁴https://platform.openai.com/docs/models#
gpt-4o

Closed-Loop Planning with ReAct In this ap-924 proach, we add the beliefs of unwritten constraints 925 from commonsense reasoning to the recipe descrip-926 tion. Then the agent performs one action at a time 927 and predicts the next action based on interaction 928 with the environment. The agent receives feedback 929 after each interaction. If an action fails, detailed 930 feedback is provided, prompting the model to re-931 flect on its beliefs about unwritten properties and 932 constraints and adapt its multitasking plan dynam-933 ically. The interaction continues until the agent 934 believes all recipes are completed or the interaction 935 is aborted by the environment. We set the temper-936 ature as 0 and the maximum generation tokens as 937 512. We parse the response to get the first action 938 to avoid action trying to execute multiple actions 939 during one interaction. The prompt for the react 940 setting is detailed in Table 12. The prompt for Re-941 Act with *feasibility priority* setting is detailed in 942 Table 13. In this setting, we prompt the model to 943 finish recipes one by one to avoid violations of time 944 constraints due to multitasking. 945

D Commonsense Reasoning Evaluation

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

We present the results of our evaluation for identifying unwritten properties and constraints in Table 8. Most of the tested models display an F1 score above 80% for identifying action dependency and object occupancy, with GPT-40 demonstrating robust performance by achieving F1 scores of 90.34% and 91.12%, respectively. Qwen2.5-32B and Qwen2.5-72B also exhibit strong commonsense capabilities in action dependency and resource limitations. However, the task of identifying autonomous actions poses a greater challenge. GPT-4o-mini achieves the highest recall at 82.18%, while Gemini-1.5-Pro exhibits the highest precision at 88.33%. While the models perform commendably in identifying action dependencies and object occupancy, there is a clear need for improvement in identifying autonomous actions, which present significant opportunities for multitasking.

E Open-Loop Planning Results

LLMs cannot plan feasible multitasking without environmental feedback. In the Plan-and-Solve setting, the model is prompted to plan an action schedule without feedback from the environment. The results in Table 7 show that even GPT-40 achieves only a success rate of 3.1% and a complete ratio of 21.7% in the *w/o time constraints*

³https://huggingface.co/meta-llama

⁵https://platform.openai.com/docs/models#
gpt-4o-mini

⁶https://ai.google.dev/gemini-api/docs/models/ gemini#gemini-1.5-pro

scenario. When time constraints are added, the
success rate and complete ratio drop further. Many
generated plans attempt to execute continuous actions simultaneously, leading to plan failure. This
suggests that current LLMs lack the planning ability to schedule multitasking without environmental
feedback.

F Planning Multitasking with Iterations

982

983

984

987

990

991

994

995

998

1000

1001

1002

1003

1004

The LLM-Modulo framework (Gundawar et al., 2024a; Kambhampati, 2024) has demonstrated that large language models (LLMs) can effectively plan complex schedules with the help of critics, as seen in benchmarks like TravelPlanner (Xie et al., 2024a), which includes multiple soft and hard constraints. In our experiments described in §4, the agent cannot recover from time constraint violations. This section evaluates whether the LLM-Modulo framework can improve model performance by providing detailed critiques of the entire plan, thus enabling more robust global planning, correcting time constraint violations, and achieving efficient and feasible multitasking.

Rather than using a step-by-step evaluation as in §4, we design different critics to assess each property and constraint outlined in §3.1 and §3.2 for detailed feedback on the whole plan including the actions after the step that the plan fails. Following the approach in Gundawar et al. (2024b), we use a reformatter to control the output format and concatenate the feedback from these critics and incorporate it into the initial prompt, performing 10 iterations. We evaluate GPT-40 and the latest o1-mini under this framework.

Iterations of critics can help LLMs plan feasi-1006 ble multitasking. We sample 20 cases from our 1007 benchmark which GPT-40 fails to solve in ReAct 1008 setting. To evaluate the global planning abilities 1009 with iterations fairly, we prompt the models with 1010 the same plan and let the model revise the plan. We 1011 find that GPT-40 and 01-mini can deliver 3 feasible 1012 plans after 10 iterations. Out of the feasible plans, o1-mini achieves a relative efficiency of 98.89% 1014 which is very close to a heuristic baseline. 1015

1016LLMs can not adjust a feasible plan for higher1017multitasking efficiency without breaking feasi-1018bility We prompt the model with 18 feasible plans1019from GPT-40 under ReAct setting, we then prompt1020the model to determine whether the plan can be fur-1021ther optimized for higher multitasking efficiency.1022GPT-40 and 01-mini can only maintain the feasi-

bility of half of the plans and the relative efficiency1023of these also decreases compared with the initial1024plans.1025

Algorithm 1: Heuristic Algorithm for Multitasking with time constraints **Input:** Set of actions \mathcal{A} , Durations \mathcal{T} , Dependencies $p(\mathcal{A})$. **Output:** Heuristic minimal time \mathcal{T}_{min} . 1 Define autonomous actions \mathcal{A}^* and continuous actions \mathcal{A}' from \mathcal{A} . ² Sort \mathcal{A}^* by \mathcal{T} in descending order. $\mathcal{A} \leftarrow \text{concatenate}(\mathcal{A}^*, \mathcal{A}').$ 4 Initialize Action_list as an empty list. **5** foreach $a_i \in \mathcal{A}$ do $P \leftarrow \text{BFS}(a_i, p(a_i))$ to collect 6 prerequisites. foreach $p_i \in P$ do 7 if $p_i \in \mathcal{A}$ then 8 Action_list.append (p_i) . 9 Remove p_i from \mathcal{A} . 10 end 11 end 12 Action_list.append(a_i). 13 14 end 15 Define function $DFS(\mathcal{A}^*, \mathcal{A}', \mathcal{T}_{\min})$: if Action_list is empty then return \mathcal{T}_{min} . 16 17 end 18 foreach $a_i \in Action_list$ do if $check_constraint(a_i)$ then 19 if $a_i \in \mathcal{A}^*$ then 20 $\mathcal{A}^* \leftarrow \mathcal{A}^* \setminus \{a_i\}.$ 21 $\mathcal{T}_{min} \leftarrow \mathcal{T}_{min}.$ 22 end 23 else 24 $\mathcal{A}' \leftarrow \mathcal{A}' \setminus \{a_i\}.$ 25 $\mathcal{T}_{\min} \leftarrow \mathcal{T}_{\min} + \mathcal{T}(a_i).$ 26 end 27 result $\leftarrow DFS(\mathcal{A}^*, \mathcal{A}', \mathcal{T}_{\min})$ 28 if result is not failure then 29 return result. 30 end 31 end 32 33 end 34 return failure. 35 $\mathcal{T}_{\min} \leftarrow 0$. 36 result $\leftarrow DFS(\mathcal{A}^*, \mathcal{A}', \mathcal{T}_{\min})$ 37 if result is failure then 38 return "No feasible schedule found." 39 end 40 else return result. 41 42 end

Recipe 1:Tacos

Step 0 (3 min): Place the fish in a cooking pot with enough water to cover them.

Step 1 (20 min): Let it boil for around 20 minutes.

Step 2 (3 min): Drain them and put them in a bowl.

Step 3 (3 min): Use a fork to smash them.

Step 4 (5 min): Chop the onion and tomato in little squares.

Step 5 (2 min): Put the onion in a pan with a little oil.

Step 6 (5 min): Let it cook a little for 5 minutes.

Step 7 (2 min): Place the tomato and stir.

Step 8 (5 min): Let it cook for 5 minutes stirring constantly. Step 9 (5 min): When it's ready put the fish in the pan and mix well.

Step 10 (1 min): Add salt and pepper to taste.

Step 11 (15 min): Let it cool down.

Step 12 (1 min): Warm the tortillas for 1 minute in the microwave.

Step 13 (2 min): Put the fish we made in the middle of the tortilla using a spoon and make sure it doesn't reach the edge. Step 14 (2 min): Fold the tortilla in half and put one toothpick on each side to hold it closed.

Step 15 (10 min): Fry them in a pan with oil

Step 16 (2 min): Take the toothpicks off and serve when they are still warm.

Interrutable steps: 3, 4.

Autonomous actions: step 1, 6, 11, 12, 15.

Action Dependency: 0->1, 1->2, 2->3, 4->5, 5->6, 6->7, 7->8, 3->9, 8->9, 9->10, 10->11, 12->13, 11->13, 13->14, 14->15, 15->16.

Steps 1, 6, 7, 8, 9, 10, 15 require stove, Steps 12 requires microwave.

Recipe 2:Smore-Bars

Step 0 (10 min): Preheat your oven to 350 degrees fahrenheit. Step 1 (3 min): Grease a 9x13 inch pan.

Step 2 (1 min): Melt your 1 cup of butter in the microwave until it is completely melted.

Step 3 (5 min): crush 2 cups (approximately 2 sleeves) of graham crackers.

Step 4 (3 min): Mix the melted butter and crushed graham crackers together.

Step 5 (5 min): Take about 3/4 (doesn't need to be exact) of your butter/graham cracker mixture and press into the bottom of your greased pan.

Step 6 ($2 \min$): Unwrap your candy bars and arrange them.

Step 7 (3 min): Evenly spread out your bag of mini marshmallows across entire pan.

Step 8 (2 min): Sprinkle your remaining butter/graham cracker mixture across pan.

Step 9 (15 min): Place pan in oven for 15 minutes.

Step 10 (2 min): Cut and Enjoy!

Interrutable steps: 1, 3, 5, 6, 7, 8, 10.

Autonomous actions: step 0, 2, 9.

Action Dependency: 2->4, 3->4, 1->5, 4->5, 5->7, 6->7, 7->8, 0->9, 8->9, 9->10.

Steps 0, 9 require oven, Steps 2 requires microwave.



Table 4: Case study of GPT-40 planning without time constraints. GPT-40 result in a lower efficiency compared with the heuristic baseline. The primary difference is that GPT-40 prioritizes different autonomous actions and leaves the agent idle during the execution of the last two actions

Recipe 2:Daikon-Radish Step 0 (5 min): Peel the skin from the radishes and rinse them. ## Recipe 1: Vada Step 1 (5 min): Slice them into thin circular slices. Step 0 (5 min): First, start with blending the chilies, ginger Step 2 (5 min): Chop napa cabbage into thin slices. and coriander along with some cumin seeds. Step 3 (5 min): Peel the skin off the onions and slice into Step 1 (3 min): Partially blend the chana dal with the mixture cubes. from the previous step. Step 4 (3 min): Thinly slice the green onions. Step 2 (5 min): Slice the onions. Step 5 (8 min): In a large non-stick pan, sauté the onion until Step 3 (3 min): Once the mixture is ready, you need to mix slightly golden brown. Step 6 (2 min): Once the onion is cooked, place the sliced the mixture with cut onion. Step 4 (3 min): Later add some coriander and curry leaves napa cabbage onto the pan. and continue mixing. Step 7 (2 min): Stir fry for about 2 minutes. Step 5 (5 min): Heat up some oil in a pan. Step 8 (8 min): Add the radish into the pan and stir fry until Step 6 (5 min): Shape the paste into circular disk-shaped soft. Step 9 (2 min): To season and garnish, add a couple of Step 7 (5 min): Deep fry the shaped chunks in the hot oil. teaspoons of soy sauce, half a teaspoon of sesame oil, and a pinch of salt and pepper. Step 8 (10 min): Fry the vada in oil until it turns golden Step 10 (3 min): Place the bacon slices onto a pan. Step 9 (5 min): Serve the dish hot and along with some Step 11 (10 min): Cook until crispy. ketchup and some mint chutney. Step 12 (5 min): Slice into smaller pieces to make bacon bits. Step 13 (2 min): Top the radish dish with some bacon bits. Interruptible steps: 0, 1, 2, 3, 4, 6. Autonomous actions: step 5, 7, 8. Interruptible steps: 0, 1, 2, 3, 4, 12, 13. Autonomous actions: step 5, 11. Action Dependency: 0->1, 1->3, 2->3, 3->4, 4->6, 5->7, 6->7, 7->8, 8->9. Action Dependency: 0->1, 1->2, 2->3, 3->4, 4->5, 5->6, 6->7, Time Constraints: 5-7 (5 min), 7->8 (5 min), 8->9 (5 min). 7->8, 8->9, 10->11, 11->12, 12->13, 9->13. Steps 5, 7, 8 require stove. Time Constraints: 5->6 (2 min), 6->7 (1 min), 7->8 (3 min), 8->9 (2 min). Steps 5, 6, 7, 8, 9, 11 require stove.

chunks.

brown.



Table 5: Case study of planning with time constraints. We only show part of the plan to focus on the source of task failure. GPT-40 plans to execute step 5 of Vada first to maximize local efficiency which leads to violation of time constraints between step 5 and step 7 of Vada. The heuristic plan strategically starts step 5 of Vada until other prerequisite actions of step 7 are nearly finished.

Model	Pass Cases	R-Efficency *
Fix Failed	d Plans	
GPT-40 o1-mini	3 / 20 3 / 20	84.50 98.89
Optimize	Feasible Plans	
GPT-40 o1-mini	10 / 18 8 / 18	93.53 (-2.62) 89.41 (-3.87)

Table 6: Results of LLM-Modulo framework on generating feasible or more efficient plans. We report a number of feasible cases and the relative efficiency. *: Average R-Efficiency for feasible plans.

.

Model	w/o time constraint				w/ time constraint			
	Success	Progress	R-Efficiency	S×E	Success	Progress	R-Efficiency	S×E
				Plan-ar	nd-Solve			
Llama-3.1-8B	0.0	6.2	30.3	1.9	0.0	6.6	35.7	2.3
Qwen2.5-7B	0.0	7.5	46.1	3.5	0.0	6.3	36.2	2.3
Qwen2.5-32B	0.0	8.6	54.5	4.7	0.0	8.8	54.7	4.8
Llama-3.1-70B	0.0	11.3	55.0	6.2	0.0	10.1	55.3	5.6
Qwen2.5-72B	1.5	17.2	65.6	11.3	0.0	13.3	63.6	8.4
GPT-4o-mini	1.5	13.4	50.0	6.7	0.0	10.3	56.7	5.8
GPT-40	3.1	21.7	70.2	15.2	1.5	17.4	69.7	12.2
Heuristics	100	100	100	100	100	100	100	100

Table 7: Results of Plan-and-Solve setting. We report Success Rate (**Success**), Average Progress Rate (**Progress**), Relative Multitask Efficiency (**R-Efficiency**) and Multitasking Ability ($S \times E$).

Constraint	Model	Recall	Precision	F1
	Qwen2.5-7B	57.68	62.98	60.21
	Llama-3.1-8B	58.26	50.87	54.31
	Qwen2.5-32B	67.92	79.96	73.45
Astion Consumances	Llama-3.1-70B	81.09	75.83	78.37
Action Concurrency	Qwen2.5-72B	74.04	70.92	72.44
	GPT-4o-mini	82.18	68.27	74.58
	Gemini-1.5-Pro	40.48	88.33	55.51
	GPT-40	73.10	79.39	76.12
Action Dependency	Qwen2.5-7B	69.67	78.15	73.67
	Llama-3.1-8B	79.68	82.61	81.12
	Qwen2.5-32B	87.41	92.18	89.73
	Llama-3.1-70B	91.57	92.58	92.07
	Qwen2.5-72B	92.18	92.82	92.50
	GPT-4o-mini	78.38	82.82	80.54
	Gemini-1.5-Pro	88.36	90.09	89.22
	GPT-40	89.94	92.33	91.12
Resource Limitations	Qwen2.5-7B	69.82	91.26	79.12
	Llama-3.1-8B	85.31	88.96	87.10
	Qwen2.5-32B	91.28	96.15	93.65
	Llama-3.1-70B	88.30	97.21	92.54
	Qwen2.5-72B	92.22	95.04	93.61
	GPT-4o-mini	92.10	89.56	90.81
	Gemini-1.5-Pro	79.95	95.45	87.02
	GPT-40	85.30	96.02	90.34

Table 8: Results of commonsense reasoning for unwritten properties and constraints.

Action Concurrency

Please identify if the action is autonomous or continuous.

- Autonomous Action: The action can be performed alongside other actions, allowing the agent to perform multiple tasks simultaneously. (e.g. preheat oven).

- Continuous Action: The step requires active involvement of the agent to complete and must be executed independently without overlapping with other tasks (e.g., 'Crack 3 eggs into a bowl').

Execution Interruptibility

A step classified as non-interruptible means that it cannot be split into two separate periods, and no other actions can be started during the execution of this action. Identify whether an action in a process can be interrupted or not

- If the action is logically interruptible (e.g., 'Dice the onions'), classify it as interruptible.

- If the action requires the agent to finish in one go(e.g., 'Keep stirring...'), classify it as non-interruptible.

- If the action involves heating (e.g., 'Melt the chocolate over low heat'), classify it as non-interruptible to ensure that the heating time is not extended.

- If the action can be executed in a short time (e.g., 'Pour water into a cup' or 'Add something into something'), classify it as non-interruptible.

Resource Limitations

Annotate the steps that use one of the following physical objects.

- Oven: You should always preheat the oven to a specific temperature before using it. If the oven is already preheated by a previous step, you can skip the preheat action.

- Microwave: Use this tool to heat something quickly. You can only microwave for one recipe at the same time.

- Stove: Use the heater to warm your pan or pot for cooking.

time constraints

Identify pairs of actions if there is a time constraint between them.

- If the object of action has been heated, the time interval between steps should be some value to avoid extending the heating time (e.g., 'Fry the okra' -> 'Mix the onion with okra'). Steps involving cooling allow for more flexible time intervals.

- If the state of an object will change over time (e.g., 'Melt butter' -> 'Mix with something'), the next step should occur within a specific time frame to ensure the desired outcome

- Please only consider the actions with the direct dependent relationship. And you do not need to specify the time interval.

Table 9: Guidelines for recipe annotation.

Recipe 1:Baked-Potato

Step 0 (10 min): Preheat the oven to 425 degrees.

Step 1 (2 min): Pierce the potato several times with a fork.

Step 2 (5 min): Bake the potato in the preheated oven.

Step 3 (1 min): Melt butter in the microwave.

Step 4 (10 min): Remove potato from the oven and use a sharp knife to make decorative cuts on the top of the potato.

Step 5 (1 min): Pour melted butter over the potato and serve.

- You can minimize the execution time based on the following properties:

You can execute only part of the action duration to pause steps 1, 4 for more efficient multitasking. But other actions must be finished without interruption.

- Do not violate any following constraints when executing this recipe: Step 5 must be performed within 2 min after Step 3 is finished.

Thoughts on the recipe:

The agent can perform autonomous actions step 0, 2, 3 in parallel with other actions to speed up the process.

The action before the arrow must be completed before the action after the arrow can be started: 0-2, 1-2, 2-24, 3-55, 4-55. The following actions would occupy the corresponding physical objects. The agent can not perform the action if the object is occupied. The properties such as volume and temperature of the object should also match the requirement of the recipe: Steps 0, 2 require oven, Steps 3 requires microwave

Table 10: Demonstration of recipe description. For the oracle setting, we replace the thoughts on the unwritten thoughts with the actual properties and constraints from the recipe

You are a multitask planner. You will plan an action sequence to finish some recipes as quickly as possible without violating any constraints.

Recipes

Each recipe is a sequence of actions designed to achieve a specific goal. Each action is a textual description companied with the duration to finish the action. Each recipe has autonomous actions such as boiling water that let the agent be idle during execution. They can be executed in parallel with other actions to speed up the process Continuous actions such as pouring water occupy the agent and only one continuous action can be executed at the same time across all recipes.

Task Description

Your task is to complete all of the recipes as quick as possible while following the recipe. The key to success is to follow the recipe and constraints, then complete the steps in the correct order while minimizing the execution time by executing the autonomous actions concurrently. First, let's analyze the recipe and create a concise plan on how to perform actions simultaneously to reduce the total execution time. Then write your action sequence following the plan. Your action should be a list of 'Step(step_num, recipe_name, time, timestamp)' which indicates the performing the given step for the given time at the timestamp. Your time and timestamp should be written as HH:MM:SS.

/*

Recipe 1:Baked-Potato Description for baked potatoes

Recipe 2:Cheese-Sandwich Description for cheese sandwiches

Plan

Start by preheating the oven for both recipes simultaneously, since they require the same temperature and preheating is an autonomous action. This allows the agent to perform other actions while the oven is preheating.

While the oven is preheating, start the preparation steps that do not need the oven. I can pierce the potato (step 1 for Baked-Potato) and prepare the Cheese-Sandwich by buttering the bread (Step 1) and then placing the cheese between the slices (Step 2). These steps only take 5 minutes, then wait for the oven to preheat.

At 00:10:00, the oven should be preheated. Bake the potato (Step 2 for Baked-Potato) and heat soup (Step 4 for Cheese-Sandwich) simultaneously. And wait for the autonomous actions to finish at 00:15:00 since the oven and microwave are occupied and no other actions can be performed.

At 00:15:00, start baking the sandwich (Step 3 for Cheese-Sandwich) and making decorative cuts on the potato (Step 4 for Baked-Potato) simultaneously. Step 4 of Baked-Potato is interruptible, so the agent can pause this action and perform step 5 of Cheese-Sandwich at 00:20:00 to serve the sandwich with the soup without violating the constraints. And the agent can continue to finish the potato with executing autonomous action step 3 and continuous action step 4 simultaneously and serve the potato at 00:26:00.

Action Sequence

Step(0, Baked-Potato, 10 min, 00:00:00), Step(0, Cheese-Sandwich, 10 min, 00:00:00), Step(1, Baked-Potato, 2 min, 00:00:00), Step(1, Cheese-Sandwich, 2 min, 00:02:00), Step(2, Cheese-Sandwich, 1 min, 00:04:00), Step(2, Baked-Potato, 5 min, 00:10:00), Step(4, Cheese-Sandwich, 5 min, 00:10:00), Step(3, Cheese-Sandwich, 5 min, 00:15:00), Step(4, Baked-Potato, 5 min, 00:15:00), Step(5, Cheese-Sandwich, 1 min, 00:20:00), Step(4, Baked-Potato, 4 min, 00:21:00), Step(3, Baked-Potato, 1 min, 00:25:00), Step(4, Baked-Potato, 1 min, 00:25:00), Step(4, Baked-Potato, 1 min, 00:26:00), Step(5, Cheese-Sandwich, 1 min, 00:25:00), Step(5, Baked-Potato, 1 min, 00:26:00), Step(5, Step(5, Step(4, Step(4, Step(5, Step(5,

Please follow the example to generate the action sequence for the following recipes.

/*

Recipe 1: Description

Recipe 2: Description

Plan

Table 11: Prompt for Plan-and-Solve setting.

You are a multitask planner. You will plan an action sequence to finish some recipes as quickly as possible without violating any constraints.

Recipes

Each recipe is a sequence of actions designed to achieve a specific goal. Each action is a textual description companied with the duration to finish the action. Each recipe has autonomous actions such as boiling water that let the agent be idle during execution. They can be executed in parallel with other actions to speed up the process Continuous actions such as pouring water occupy the agent and only one continuous action can be executed at the same time across all recipes.

Task Description

You are required to analyze the current status of the environment and decide the next action to take so that you can finish the recipes in the shortest time without violating constraints. Give your thoughts on the given status, action history and observation. Your analysis should be within 100 words starting with 'Thought:'. Then you should choose your next action. If you think you have already completed all the recipes, please output 'Action: Finish'. Write your action as 'Action: Step(step_num, recipe_name, time, timestamp)' - Perform the given step for the given time at the timestamp. Your time and timestamp should be written as HH:MM:SS. You can only perform one action each time. Do not repeat actions that are already in progress. If you choose to wait for current actions to finish, please state the time you will wait for. Then write your next action as 'I will wait and perform the next action at HH:MM:SS. Action: Step(step_num, recipe_name, time, timestamp).

/*

Recipe 1:Baked-Potato Description for baked potatos

Recipe 2:Cheese-Sandwich Description for cheese sandwichs

Action Sequence

Observation: Continuous action executed successfully. The current timestamp is 0:5:0. Status of physical objects: oven is occupied, temperature is 425.0; microwave is not occupied; stove is not occupied; You are currently executing the following autonomous actions: step 0 of Baked-Potato (will finish at 0:10:0), step 0 of Cheese-Sandwich (will finish at 0:10:0) Thought: I have completed the preparatory steps for the Cheese-Sandwich. Since both recipes are still preheating in the oven, I should wait for the preheating to complete. After preheating, both recipes are ready for the next baking step. Step 4 of the Baked-Potato is a long continuous step that can be executed in parallel with steps 3 and 4 of the Cheese-Sandwich. So I should bake the Baked-Potato first.

I will wait and perform the next action at 00:10:00. Action: Step(2, Baked-Potato, 00:05:00, 00:10:00)

•••

Observation: Autonomous action started successfully. The current timestamp is 0:15:0. Status of physical objects: oven is occupied, temperature is 425.0; microwave is not occupied; stove is not occupied; You are currently executing the following autonomous actions: step 4 of Cheese-Sandwich (will finish at 0:20:0)

Thought: With the Cheese-Sandwich currently baking in the oven, I can now focus on the the Baked-Potato. I will remove the Baked-Potato from the oven and prepare it for serving. I should not start melting the butter for the Baked-Potato before nearly finishing the Step 4 of Baked-Potato to avoid the butter cooling down. So I will start making decorative cuts on the top of the potato first. And I will stop at 0:20:00 to serve the Cheese-Sandwich to avoid violating the time constraint. Action: Step(4, Baked-Potato, 00:05:00, 00:15:0)

•••

Observation: Continuous action executed successfully. The current timestamp is 0:21:0. Status of physical objects: oven is occupied, temperature is 425.0; microwave is not occupied; stove is not occupied

Thought: I have served the Cheese-Sandwich and can now continue with the Baked-Potato. I will prepare the melted butter and pour it over the potato to complete the recipe. Since the step 5 must be performed within 2 min after Step 3, I will start melting the butter when step 4 is nearly finished in 1 minute. Now I would continue with the decorative cuts on the top of the potato.

Action: Step(4, Baked-Potato, 00:04:00, 00:21:0)

Observation: Continuous action executed successfully. The current timestamp is 0:27:0. Status of physical objects: oven is occupied, temperature is 425.0; microwave is occupied; stove is not occupied Action: Finish

*/

Please follow the example to generate the action sequence for the following recipes.

/* ## Recipe 1:

Description

Recipe 2: Description

Action Sequence

You are a multitask planner. You will plan an action sequence to finish some recipes without violating any constraints.

Recipes

Each recipe is a sequence of actions designed to achieve a specific goal. Each action is a textual description companied with the duration to fininsh the action. Each recipe has autonomous actions such as boiling water that lets the agent be idle during execution. They can be executed in parallel with other actions to speed up the process Continuous actions such as pouring water occupy the agent and only one continuous action can be executed at the same time across all recipes.

Task Description

You are required to analyze the current status of the environment and decide the next action to take so that you can finish the recipes without violating constraints. Do not rush to complete the recipes. Take your time to execute the actions without violating constraints. especially the time constraints. You should prioritize on not violating the time constraints and ensure the recipe is completed successfully.

Give your thoughts on the given status, action history and observation. If you find your initial thoughts of the recipes does not align with the current status, you can revise your initial thoughts. Your analysis should starts with 'Thought:'. Then you should choose your next action. If you think you have already completed all the recipes, please output 'Action: Finish'. Write your action as 'Action: Step(step_num, recipe_name, time, timestamp)' - Perform the given step for the given time at timestamp. Your time and timestamp should be written as HH:MM:SS. You can only perform one action each time. Do not repeat actions that are already in progress. If you choose to wait for current actions to finish, please state the time you will wait for. Then write your next action as 'I will wait and perform the next action at HH:MM:SS. Action: Step(step_num, recipe_name, time, timestamp).

/*

Recipe 1:Baked-Potato

Description for baked potatos

Observation: The current timestamp is 00:00:00. Status of physical objects: oven is not occupied, is not preheated; microwave is not occupied; stove is not occupied

Thought: I will preheat the oven. This will allow me to perform other actions while the oven is preheating.

Action: Step(0, Baked-Potato, 00:10:00, 00:00:00)

Observation: The current timestamp is 00:00:00. Status of physical objects: oven is occupied, temperature is 425.0; microwave is not occupied; stove is not occupied; You are currently executing the following autonomous actions: step 0 of Baked-Potato (will finish at 00:10:00)

Thought: Since the oven is preheating, I can start piercing the potato for the Baked-Potato recipe. And this does not risk violating any time constraints.

Action: Step(1, Baked-Potato, 00:02:00, 00:02:00)

Observation: The current timestamp is 00:04:00. Status of physical objects: oven is occupied, temperature is 425.0; microwave is not occupied; stove is not occupied; You are currently executing the following autonomous actions: step 0 of Baked-Potato (will finish at 00:10:00), step 0 of Cheese-Sandwich (will finish at 00:10:00)

Thought: Step 1 of the Baked-Potato is complete. I will now move on to step 2 of the Baked-Potato. This does not risk violating any time constraints.

Action: Step(2, Baked-Potato, 00:05:00, 00:10:00)

Observation: The current timestamp is 00:10:00. Status of physical objects: oven is occupied, temperature is 425.0; microwave is not occupied; stove is not occupied; You are currently executing the following autonomous actions: step 2 of Baked-Potato (will finish at 00:15:00)

Thought: Step 2 of the Baked-Potato is complete. Since Step 5 depends on both Step 3 and Step 4 and Step 5 must be performed within 2 minutes after Step 3 is complete, and. Executing Step 3 first will result in violation of time constraints since Step 4 is a long continuous step. I will now first finish Step 4 of the Baked-Potato to ensure the time constraints are followed.

Action: Step(4, Baked-Potato, 00:10:00, 00:15:00)

•••

Observation: The current timestamp is 00:26:00. Status of physical objects: oven is not occupied, temperature is 425.0; microwave is not occupied; stove is not occupied Thought: The recipe is completed.

Action: Finish

/ Please follow the example to generate the action sequence for the following recipes. / ## Recipe 1: Description

Action Sequence

Table 13: Prompt for ReAct-style action sequence generation with *feasibility priority*.

	Feedback				
Wrong Recipe	Recipe C is not one of our goals. Please select actions from our recipes A and B				
Wrong Action	There is no step X in recipe A.				
Time Error	The current time is HH:MM:SS. You can not perform any actions before the current times-				
	tamp.				
Infeasible Multitask	You can not perform step X of Recipe A and step Y of Recipe B simultaneously since they are all continuous actions.				
Object Occupancy	Object M is currently occupied.				
Object Temperature	The Property of the Object is α , but step X of recipe A needs β .				
Object Volume	The object only has a volume of water. It is not enough for this action.				
Dependency	Step Y of recipe A can not be performed because prerequisite step X is not completed. (The				
	expected finish time is HH:MM:SS)				
Repeated Action	Prerequisite step X is already used for the next action step Y in recipe A. You should not				
	execute the same step twice. If you insist, please complete all the previous steps first.				
time constraint	The time interval between Step X and Step Y in Recipe A exceeds the allowed time limit t				
Action Duration	Your plan execution time t min exceeds the time needed to perform the action.				
Execution Interruptibility	Step X of Recipe A is not interruptable. You should finish the action in one go.				
	Observation				
Success Execution	Autonomous / Continuous action executed successfully. Stove is not occupied; Oven is not				
	occupied, temperature is t . You are executing step X of recipe A.				
Failed Execution	Step X of Action A can not be executed.				
Hint					
Executable Actions	The following actions are ready to be executed after HH:MM:SS, Step X of Recipe A, Step Y of Recipe B.				

Table 14: Examples of observation, feedback and hints from the environment.

Critic Example

Critic for the plan:

Plan Completeness: The following actions are missing in your plan: step 2 of Cobbler; step 3 of Cobbler; step 4 of Cobbler; step 5 of Cobbler; step 6 of Pancakes; step 7 of Pancakes. Include them in the plan to complete the recipe.

Action Duration: The duration of the following actions do not align with the action duration: Pancakes step 9; Pancakes step 10; Pancakes step 12; Make sure the duration of the actions are correct.

Action Concurrency: You can not start another action while executing a continuous action. In your plan, the following actions can not be performed simultaneously with each other: step 0 of Pancakes and step 1 of Pancakes; step 1 of Pancakes and step 5 of Pancakes; step 6 of Cobbler and step 3 of Cobbler . Please adjust the timeline to avoid the conflict.

Action Interruption: The following actions should not be interrupted in your plan: step 14 of Pancakes. Make sure they are finished in one go.

Action Dependency: step 1 of Pancakes can be performed only after prerequisite action step 0 of Pancakes is finished. You should complete the prerequisites before performing the next action.

Time constraint: The following action pairs violate the time constraint: step 5 of Cobbler should start within 2 min after step 3 of Cobbler is finished; step 14 of Pancakes should start within 2 min after step 12 of Pancakes is finished. Reschedule the actions to meet the time constraint.

Physical Object: Step 14 of Pancakes can not be performed at time 00:02:00 due to Object stove is occupied. Adjust the use of the physical objects to meet the requirements.

Multitasking Efficiency The plan is feasible. The agent is idle during the following timestamps: HH:MM:SS and HH:MM:SS. You can assign continuous actions to the agent to optimize the plan for a shorter execution time. If you think the plan is optimal, you can answer Action: Done to finish the task.

Table 15: Critic example for LLM-Modulo framework.