

An LLM Compiler for Parallel Function Calling

Sehoon Kim^{*1} Suhong Moon^{*1} Ryan Tabrizi¹ Nicholas Lee¹
Michael W. Mahoney^{1,2,3} Kurt Keutzer¹ Amir Gholami^{1,2}

¹ UC Berkeley ² ICSI ³ LBNL

{sehoonkim, suhong.moon, rtabrizi, nicholas_lee, mahoneymw, keutzer, amirgh}@berkeley.edu

Abstract

Recent language models have shown remarkable results on various complex reasoning benchmarks. The reasoning capabilities of LLMs enable them to execute external function calls to overcome their inherent limitations, such as knowledge cutoffs, poor arithmetic skills, or lack of access to private data. This development has allowed LLMs to select and coordinate multiple functions based on the context to tackle more complex problems. However, current methods for multiple function calling often require sequential reasoning and acting for each function which can result in high latency, cost, and sometimes inaccurate behavior. To address this, we introduce `LLMCompiler`, which executes functions in parallel to efficiently orchestrate multiple function calling. Drawing from the principles of classical compilers, `LLMCompiler` streamlines parallel function calling with three components: (i) an LLM Planner, formulating execution plans; (ii) a Task Fetching Unit, dispatching function calling tasks; and (iii) an Executor, executing these tasks in parallel. `LLMCompiler` automatically generates an optimized orchestration for the function calls and can be used with both open-source and closed-source models. We have benchmarked `LLMCompiler` on a range of tasks with different patterns of function calling. We observe consistent latency speedup of up to 3.7 \times , cost savings of up to 6.7 \times , and accuracy improvement of up to \sim 9% compared to ReAct. Our code is available at <https://github.com/SqueezeAILab/LLMCompiler>.

1 Introduction

Recent advances in the reasoning capability of Large Language Models (LLMs) have expanded the applicability of LLMs beyond content generation to solving complex problems [51, 23, 57, 5, 50, 63, 8, 55, 12]; and recent works have shown how this reasoning capability could be helpful in improving accuracy for solving complex and logical tasks. The reasoning capability has also allowed function (i.e., tool) calling capability where LLMs can invoke provided functions and use the function outputs to help complete their tasks. These functions range from a simple calculator that can invoke arithmetic operations to more complex LLM-based functions.

The ability of LLMs to integrate various tools and function calls could enable a fundamental shift in how we develop LLM-based software. However, this brings up an important challenge: *what is the most effective approach to incorporate multiple function calls?* A notable approach has been introduced in ReAct [58], where the LLM calls a function, analyzes the outcomes, and then reasons about the next action, which involves a subsequent function call. For a simple example illustrated in Figure 1 (Left), where the LLM is asked if Scott Derrickson and Ed Wood have the same nationality, ReAct initially analyzes the query and decides to use a search tool to search for Scott Derrickson. The result of this search (i.e., observation) is then concatenated back to the original prompt for the LLM to reason about the next action, which invokes another search tool to gather information about Ed Wood.

ReAct has been a pioneering work in enabling function calling, and it has been integrated into several frameworks [26, 31]. However, scaling this approach for more complex applications requires considerable optimizations. This is due to the sequential nature of ReAct, where it executes function calls and reasons about their observations one after the other. This approach, along with the agent systems that extend ReAct [20, 57, 41, 42, 47], may lead to inefficiencies in latency and cost due to the sequential function calling and repetitive LLM invocations for each reasoning and action step. Furthermore, while dynamic reasoning about the observations has benefits in certain cases,

^{*}Equal contribution

HotpotQA Question: Were Scott Derrickson and Ed Wood of the same nationality?

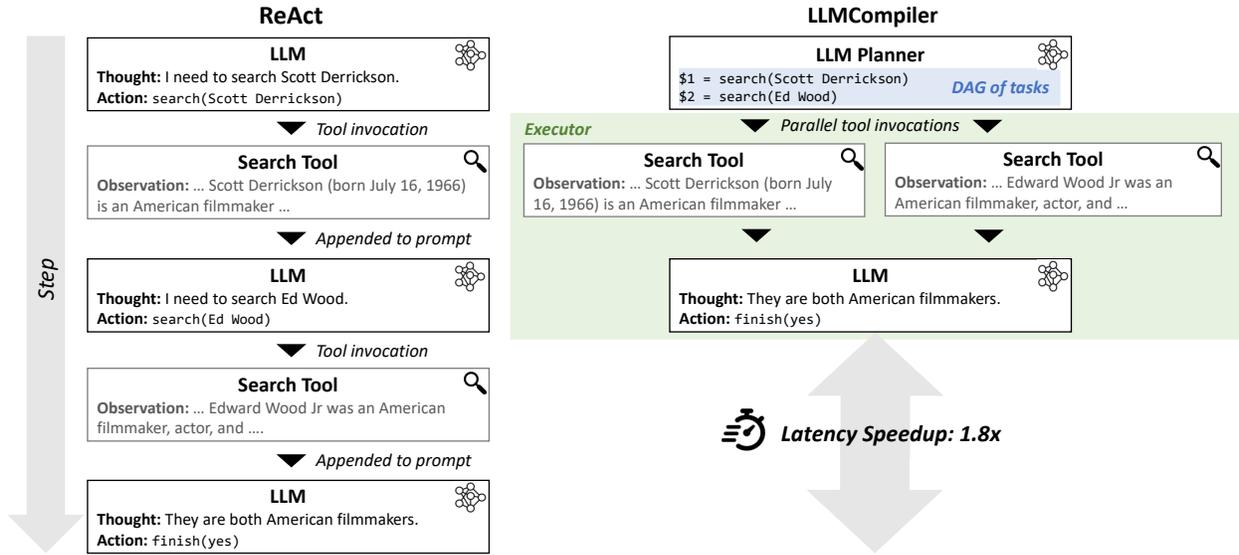


Figure 1: An illustration of the runtime dynamics of LLMCompiler, in comparison with ReAct [58], given a sample question from the HotpotQA benchmark [54]. In LLMCompiler (Right), the Planner first decomposes the query into several tasks with inter-dependencies. The Executor then executes multiple tasks in parallel, respecting their dependencies. Finally, LLMCompiler joins all observations from the tool executions to produce the final response. In contrast, sequential tool execution of the existing frameworks like ReAct (Left) leads to longer execution latency. In this example, LLMCompiler attains a latency speedup of $1.8\times$ on the HotpotQA benchmark. While a 2-way parallelizable question from HotpotQA is presented here for the sake of simple visual illustration, LLMCompiler is capable of managing tasks with more complex dependency patterns (Fig. 2 and Sec. 4).

concatenating the outcomes of intermediate function calls could disrupt the LLM’s execution flow, potentially reducing accuracy [53]. Common failure cases include repetitive invocation of the same function, which is also highlighted in the original paper [58], and early stopping based on the partial intermediate results, as will be further discussed in Sec. 4.1 and Appendix. A.1.

To address this challenge, we draw inspiration from classical compilers, where optimizing instruction executions in traditional programming languages has been extensively explored. A key optimization technique in compilers involves identifying instructions that can be executed in parallel and effectively managing their dependencies. Similarly, one can envision a compiler tailored for LLM function calling, which can efficiently orchestrate various function calls and their dependencies. This shares a similar philosophy with the recent studies that align LLMs with computer systems [19, 37]. To this end, we introduce LLMCompiler, a novel framework that enables parallel multi-tool execution of LLMs across different models and workloads. To the best of our knowledge, LLMCompiler is the first framework to optimize the orchestration of LLM function calling that can not only improve latency and cost, but also accuracy by minimizing interference from the outputs of intermediate function calls. In more detail, we make the following contributions:

- We introduce LLMCompiler, an LLM compiler that optimizes the parallel function calling performance of LLMs. At a high level, this is achieved by introducing three key components: (i) an LLM Planner (Sec. 3.1) that identifies an execution flow; (ii) a Task Fetching Unit (Sec. 3.2) that dispatches the function calls in parallel; (iii) an Executor (Sec. 3.3) that executes the dispatched tasks using the associated functions.
- We evaluate LLMCompiler on *embarrassingly parallel* patterns using HotpotQA [54] and Movie Recommendation [4], where we observe $1.80\times/3.74\times$ speedup and $3.37\times/6.73\times$ cost reduction compared to ReAct (Sec. 4.1).
- To test the performance on more complex patterns, we introduce a new benchmark called ParallelQA which includes various non-trivial function calling patterns. We show up to $2.27\times$ speedup, $4.65\times$ cost reduction, and 9% improved accuracy compared to ReAct (Sec. 4.2).

- We evaluate `LLMCompiler`'s capability in dynamic replanning, which is achieved through a feedback loop from the Executor back to our LLM Planner. For the Game of 24 [57], which requires repeated replanning based on the intermediate results, `LLMCompiler` demonstrates a $2\times$ speedup compared to Tree-of-Thoughts (Sec. 4.3).
- We showcase that `LLMCompiler` can explore the interactive decision-making environment effectively and efficiently. On WebShop, `LLMCompiler` achieves up to $101.7\times$ speedup and 25.7% improved success rate compared to the baselines. (Sec. 4.4)

2 Related Work

2.1 Latency Optimization in LLMs

Various studies have focused on optimizing model design [21, 11, 30, 9, 24, 10, 22, 6, 28] and systems [25, 59, 1, 2] for efficient LLM inference. Optimizations at the application level, however, are less explored. This is critical from a practical point of view for situations involving black-box LLM models and services where modifications to the models and the underlying inference pipeline are highly restricted.

Skeleton-of-Thought [34] recently proposed to reduce latency through application-level parallel decoding. This method involves a two-step process of an initial skeleton generation phase, followed by parallel execution of skeleton items. However, it is primarily designed for embarrassingly parallel workloads and does not support problems that have inherently interdependent tasks, as it assumes no dependencies between skeleton tasks. This limits its applicability in complex scenarios such as coding [7, 33, 14, 3] or math [16, 15] problems, as also stated in the paper [34]. `LLMCompiler` addresses this by translating an input query into a series of tasks with inter-dependencies, thereby expanding the spectrum of problems it can handle.

Concurrently to our work, OpenAI has recently introduced a parallel function calling feature in their 1106 release, enhancing user query processing through the simultaneous generation of multiple function calls [36]. Despite its potential for reducing LLM execution time, this feature has certain limitations as it is exclusively available for OpenAI's proprietary models. However, there is a growing demand for using open-source models driven by the increasing number of open-source LLMs as well as parameter-efficient training techniques [27, 18, 17] for finetuning and customization. `LLMCompiler` enables efficient parallel function calling for open-source models, and also, as we will show later in Sec. 4, it can potentially achieve better latency and cost.

2.2 Plan and Solve Strategy

Several studies [52, 38, 40, 63, 13] have explored prompting methods of breaking down complex queries into various levels of detail to solve them, thereby improving LLM's performance in reasoning tasks. Specifically, Decomposed Prompting [20] tackles complex tasks by decomposing them into simpler sub-tasks, each optimized through LLMs with dedicated prompts. Step-Back Prompting [61] enables LLMs to abstract high-level concepts from details to enhance reasoning abilities across various tasks. Plan-and-Solve Prompting [49] segments multi-step reasoning tasks into subtasks to minimize errors and improve task accuracy without manual prompting. However, these methods primarily focus on improving the accuracy of reasoning benchmarks. In contrast, `LLMCompiler` uses a planner to identify parallelizable patterns within queries, aiming to reduce latency while maintaining accuracy.

A notable work is ReWOO [53] which employs a planner to separate the reasoning process from the execution and observation phases to decrease token usage and cost as compared to ReAct. Our approach is different from ReWOO in multiple aspects. First, `LLMCompiler` allows parallel function calling which can reduce latency as well as cost. Second, `LLMCompiler` supports dynamic replanning which is important for problems whose execution flow cannot be determined statically in the beginning (Sec. 4.3).

2.3 Tool-Augmented LLMs

A notable work is Toolformer [43], which produces a custom LLM output to let the LLM decide what the inputs for calling the functions should be and where to insert the result. This approach has inspired various tool calling frameworks [29, 44]. ReAct [58] proposed to have LLMs interact with external environments through reasoning and action generation for improved performance. Gorilla [39] introduced a finetuned LLM designed for function calling, and ToolLLM [41] and RestGPT [46] have extended LLMs to support real-world APIs. Moreover, OpenAI [35] released their own function calling capabilities, allowing their LLMs to return formatted JSON for execution.

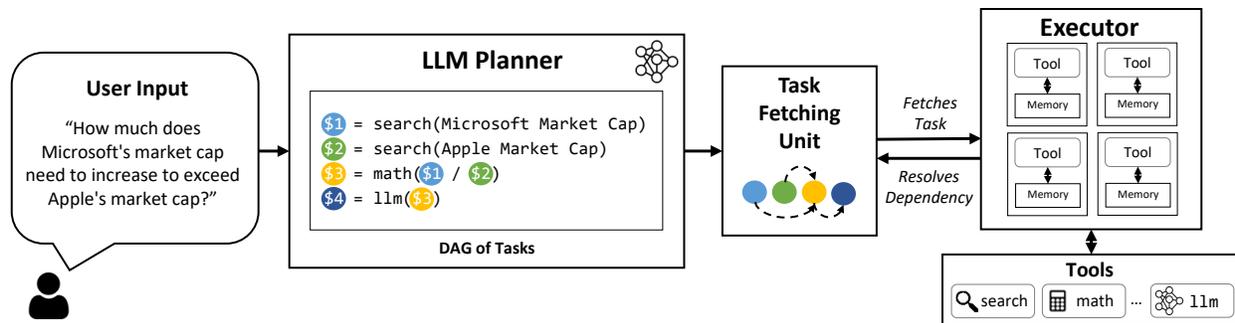


Figure 2: Overview of the `LLMCompiler` framework. The LLM Planner generates a DAG of tasks with their inter-dependencies. These tasks are then dispatched by the Task Fetching Unit to the Executor in parallel based on their dependencies. In this example, Task \$1 and \$2 are fetched together for parallel execution of two independent search tasks. After each task is performed, the results (i.e., observations) are forwarded back to the Task Fetching Unit to unblock the dependent tasks after replacing their placeholder variables (e.g., the variable \$1 and \$2 in Task \$3) with actual values. Once all tasks have been executed, the final answer is delivered to the user.

3 Methodology

To illustrate the components of `LLMCompiler`, we use a simple 2-way parallel example in Figure 2. To answer “How much does Microsoft’s market cap need to increase to exceed Apple’s market cap?”, the LLM first needs to conduct web searches for both companies’ market caps, followed by a division operation. While the existing frameworks, including ReAct, perform these tasks sequentially, it is evident that they can be executed in parallel. The key question is how to automatically determine which tasks are parallelizable and which are interdependent, so we can orchestrate the execution of the different tasks accordingly. `LLMCompiler` accomplishes this through a system that consists of the following three components: an LLM Planner (Section 3.1) that generates a sequence of tasks and their dependencies; a Task Fetching Unit (Section 3.2) that replaces arguments based on intermediate results and fetches the tasks; and an Executor (Section 3.3) that executes the tasks with associated tools. To use `LLMCompiler`, users are only required to provide tool definitions, and optional in-context examples for the Planner, as will be further discussed in Appendix A.3.2.

3.1 LLM Planner

The LLM Planner is responsible for generating a sequence of tasks to be executed along with any dependency among them. For instance, Tasks \$1 and \$2 in Figure 2 are two independent searches that can be performed in parallel. However, Task \$3 has a dependency on the outcomes of the first and second searches. Therefore, the Planner’s role is to automatically identify the necessary tasks, their input arguments, as well as their inter-dependencies using the sophisticated reasoning capability of LLMs, essentially forming a directed acyclic graph of task dependencies. If a task is dependent on a preceding task, it incorporates a placeholder variable, such as \$1 in Task 3 of Figure 2, which will later be substituted with the actual output from the preceding task (Sec. 3.2).

The Planner in `LLMCompiler` leverages LLMs’ reasoning capability to decompose tasks from natural language inputs. To achieve this, the Planner LLM incorporates a pre-defined prompt that guides it on how to create dependency graphs and to ensure correct syntax (see Appendix A.7 for details). Besides this, users also need to supply tool definitions and optional in-context examples for the Planner. These examples provide detailed demonstrations of task decomposition specific to a problem, helping the Planner to better understand the rules. Further details on user-supplied information for `LLMCompiler` are elaborated in Appendix A.3.2. In Appendix A.3.1, we introduce an additional optimization for the Planner that streams tasks as soon as they are created, instead of waiting to complete the entire planning process.

3.2 Task Fetching Unit

The Task Fetching Unit, inspired by the instruction fetching units in modern computer architectures, fetches tasks to the Executor as soon as they are ready for (parallel) execution based on a greedy policy. Another key functionality

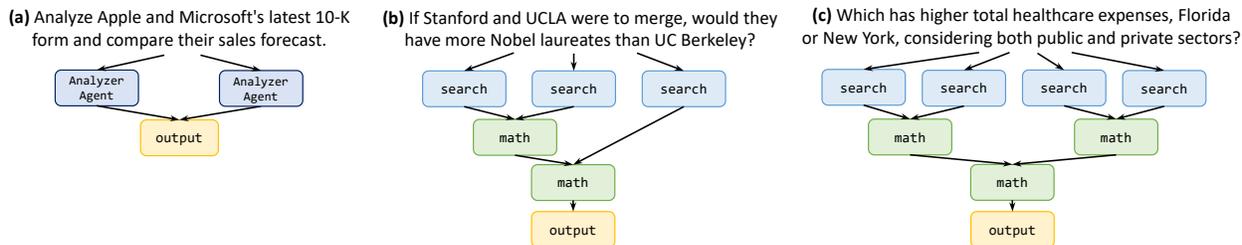


Figure 3: Examples of questions with different function calling patterns and their dependency graphs. HotpotQA and Movie Recommendation datasets exhibit pattern (a), and ParallelQA dataset exhibits patterns (b) and (c), among other patterns. In (a), we need to analyze each company’s latest 10-K. In (b), we need three searches for each school, followed by one addition and one comparison operation. In (c), we need to search for each state’s annual healthcare spending in each sector, sum each state’s spending, and then perform a comparison.

is to replace variables with the actual outputs from preceding tasks, which were initially set as placeholders by the Planner. For the example in Figure 2, the variable $\$1$ and $\$2$ in Task $\$3$ would be replaced with the actual market cap of Microsoft and Apple after the search tasks are done. This can be implemented with a simple fetching and queuing mechanism without a dedicated LLM.

3.3 Executor

The Executor asynchronously executes tasks fetched from the Task Fetching Unit. As the Task Fetching Unit guarantees all the tasks dispatched to the Executor are independent, it can simply execute them concurrently. The Executor is equipped with user-provided tools, and it delegates the task to the associated tool. These tools can be simple functions like a calculator, Wikipedia search, or API calls, or they can even be LLM agents that are tailored for a specific task. As depicted in the Executor block of Figure 2, each task has dedicated memory to store its intermediate outcomes, similar to what typical sequential frameworks do when aggregating observations as a single prompt [58]. Upon completion of the task, the final results are forwarded as input to the tasks dependent on them.

3.4 Dynamic Replanning

In various applications, the execution graph may need to adapt based on intermediate results that are a priori unknown. A similar analogy in programming is branching, where the path of execution is determined only during runtime, depending on which branch conditions are satisfied. Such dynamic execution patterns can also appear with LLM function calling. For simple branching (e.g., if-else statements) one could statically compile the execution flow and choose the right dynamically based on the intermediate results. However, for more complex branching it may be better to do a recompilation or replanning based on the intermediate results. In replanning, the Executor sends the intermediate results back to our LLM Planner. Based on that, the Planner produces a new set of tasks with their associated dependencies and dispatches them to the Task Fetching Unit and then the Executor. This process is repeated until the final result is achieved. We show an example use case of this in Sec. 4.3 for solving the Game of 24 using the Tree-of-Thoughts approach.

4 Results

In this section, we evaluate LLMCompiler using a variety of models and problem types. We use both the proprietary GPT models and the open-source LLaMA-2 model, with the latter demonstrating LLMCompiler’s capability in enabling parallel function calling in open-source models. Furthermore, there are various types of parallel function calling patterns that can be addressed with LLMs. This ranges from embarrassingly parallel patterns, where all tasks can be executed in parallel without any dependencies between them, to more complex dependency patterns, as illustrated in Figure 3. Significantly, we also assess LLMCompiler on the Game of 24 benchmark involving dynamic replanning based on intermediate results, highlighting its adaptability to dynamic dependency graphs. Finally, we apply LLMCompiler to the WebShop benchmark to showcase its potential in decision-making tasks. In this section, we start presenting results for simple execution patterns, and then we move to more complex ones.

Table 1: Accuracy and latency comparison of LLMCompiler compared to the baseline on different benchmarks, including HotpotQA, Movie Recommendation, our custom dataset named ParallelQA, and the Game of 24. For HotpotQA and Movie Recommendation, we frequently observe looping and early stopping (Sec. 4.1). We incorporated ReAct-specific prompting to minimize these behaviors as much as possible, which we denote as ReAct[†]. ReAct indicates the original results without this prompting. We do not include the latency for the original ReAct since looping and early stopping make precise latency measurement difficult.

Benchmark	Method	GPT (Closed-source)			LLaMA-2 70B (Open-source)		
		Accuracy (%)	Latency (s)	Speedup	Accuracy (%)	Latency (s)	Speedup
HotpotQA	ReAct	61.52	-	-	54.74	-	-
	ReAct [†]	62.47	7.12	1.00×	54.40	13.44	1.00×
	OAI Parallel Function	62.05	4.42	1.61×	-	-	-
	LLMCompiler	62.00	3.95	1.80×	57.83	9.58	1.40×
Movie Rec.	ReAct	68.60	-	-	70.00	-	-
	ReAct [†]	72.47	20.47	1.00×	70.60	33.37	1.00×
	OAI Parallel Function	77.00	7.42	2.76×	-	-	-
	LLMCompiler	77.13	5.47	3.74×	77.80	11.83	2.82×
ParallelQA	ReAct	89.09	35.90	1.00×	59.59	15.47	1.00×
	OAI Parallel Function	87.32	19.29	1.86×	-	-	-
	LLMCompiler	89.38	16.69	2.15×	68.14	26.20	2.27×
Game of 24	Tree-of-Thoughts	74.00	241.2	1.00×	30.00	952.06	1.00×
	LLMCompiler	75.33	83.6	2.89×	32.00	456.02	2.09×

4.1 Embarrassingly Parallel Function Calling

The simplest scenario involves an LLM using a tool repeatedly for independent tasks such as conducting parallel searches or analyses to gather information on different topics like the pattern depicted in Figure 3 (a). While these tasks are independent of each other and can be executed in parallel, ReAct, along with other LLM solutions as they stand, would need to run sequentially. This leads to increased latency and token consumption due to its frequent LLM invocations for each tool usage, as also illustrated in Figure 1. In this section, we demonstrate how LLMCompiler can identify parallelizable patterns and execute independent tasks concurrently to resolve this issue using the following two benchmarks:

- **HotpotQA:** A dataset that evaluates multi-hop reasoning [54]. We only use the comparison dev set. This contains 1.5k questions comparing two different entities, thus exhibiting a 2-way embarrassingly parallel execution pattern. An example question is shown in Figure 1.
- **Movie Recommendation:** A dataset with 500 examples that asks to identify the most similar movie out of four options to another set of four movies, exhibiting an 8-way embarrassingly parallel pattern [4].

Experimental Setups. As a baseline method, we compare LLMCompiler with ReAct. We follow the ReAct [58] setup using the same Wikipedia search tool that LLMs can use to search for information. We did not include the lookup tool since it is not relevant to our problem setting. We have optimized the prompt and in-context examples for both ReAct and LLMCompiler to the best of our abilities. For all experiments across these datasets, we use gpt-3.5-turbo (1106 release). For the experiments using GPT, we additionally report the results using OpenAI’s parallel function calling capability, which was announced concurrently with our work. We also show how LLMCompiler can be effectively combined with the open-source LLaMA-2 70B model to provide the model with parallel function calling capabilities. For all experiments, we have measured accuracy, end-to-end latency, as well as input and output token usage. See A.4 for more details on experimental setups.

Accuracy and Latency. We report the accuracy, end-to-end latency, and relative speed-up of LLMCompiler compared to ReAct in Table 1. First, we observe that ReAct consistently achieves lower accuracy compared to OpenAI parallel function calling and LLMCompiler. We identify two main failure modes in ReAct: (1) the tendency for redundant generation of prior function calls, a point also noted in the original ReAct paper [58]; and (2) premature early stopping based on the incomplete intermediate results. In Appendix. A.1, we offer a detailed analysis demonstrating how these two prevalent failure cases significantly hurt ReAct’s accuracy, and how they can be resolved with LLMCompiler, leading to an accuracy enhancement of up to 7 – 8%. Furthermore, we have conducted interventional

Table 2: Input and output token consumption as well as the estimated cost on HotpotQA, Movie Recommendation, and our custom dataset named ParallelQA. The cost is computed based on the pricing table of the GPT models used for each benchmark.

Benchmark	Method	In. Tokens	Out. Tokens	Cost (\$/1k)	Cost Red.
HotpotQA	ReAct	2900	120	5.00	1.00×
	OAI Para. Func.	2500	63	2.66	1.87×
	LLMCompiler	1300	80	1.47	3.37×
Movie Rec.	ReAct	20000	230	20.46	1.00×
	OAI Para. Func.	5800	160	6.14	3.33×
	LLMCompiler	2800	115	3.04	6.73×
ParallelQA	ReAct	46000	470	480	1.00×
	OAI Para. Func.	25000	370	260	1.81×
	LLMCompiler	9200	340	103	4.65×

experiments in which we incorporated ReAct-specific prompts to avoid repetitive function calls and early stopping. ReAct[†] in Table 1 refers to ReAct *with* this ReAct-specific prompt. The ReAct-specific prompt yields a general accuracy improvement with ReAct[†] as compared to the original ReAct. Nevertheless, LLMCompiler still demonstrates on-par and better accuracy than ReAct[†] as such prompting does not serve as a perfect solution to completely avoiding the erroneous behavior of ReAct.

Additionally, when compared to ReAct[†], LLMCompiler demonstrates a noticeable speedup of 1.80× and 1.40× on the HotpotQA benchmark with GPT and LLaMA, respectively. Similarly, LLMCompiler demonstrates 3.74× and 2.82× speedup on the Movie Recommendation benchmark with each model. Note that we benchmark the latency of LLMCompiler against that of ReAct[†] since the repeating and early stopping behavior of the original ReAct as discussed above makes its latency unpredictable and unsuitable for a fair comparison. LLMCompiler demonstrates a speedup of up to 35% compared to OpenAI parallel function calling whose latency gain over ReAct is 1.61× and 2.76× on each benchmark¹.

Costs. Another important consideration of using LLMs is cost, which depends on the input and output token usage. The costs for GPT experiments are provided in Table 2. LLMCompiler is more cost-efficient than ReAct for cost as it involves less frequent LLM invocations. Interestingly, LLMCompiler also outperforms the recent OpenAI parallel function calling in cost efficiency. This is because LLMCompiler’s planning phase is more prompt length efficient than that of OpenAI parallel function calling since our Planner’s in-context examples are rather short and only include plans, not observations (see Appendix A.7).

4.2 Parallel Function Calling with Dependencies

The cases considered above are rather simple, as only one tool is used and all tasks can be executed independently of one another. However, similar to code execution in traditional code blocks, we may encounter function calling scenarios that involve more complex dependencies. To systematically evaluate the capability to plan out function calling in scenarios that involve complex task dependencies, we have designed a custom benchmark called ParallelQA. This benchmark is designed to incorporate non-trivial function calling patterns, including three different types of patterns in Figure 3 (b) and (c). Inspired by the IfQA benchmark [60], ParallelQA contains 113 examples that involve mathematical questions on factual attributes of various entities. In particular, completing the task requires using two tools (i.e., search and math tools), with the second tool’s argument depending on the result of the first tool’s output. We have meticulously included questions that are only answerable with information from Wikipedia’s first paragraph, effectively factoring out the failure cases due to unsuccessful searches. See Appendix A.8 for more details in ParallelQA.

Experimental Setups. Similar to Sec. 4.1, we use ReAct [58] as the main baseline. Here, both LLMCompiler and ReAct are equipped with two tools: (1) the search tool, identical to the one mentioned in Sec.4.1; and (2) the math

¹ Unfortunately, we are unable to conclude why this is the case, as OpenAI has not publicly disclosed any details about their function calling mechanism. One speculation is that there might be additional overheads to validate the function and argument names and to convert them into a system prompt. Nevertheless, we have seen a consistent trend with multiple runs over several days.

tool, which solves mathematical problems. The math tool is inspired by the Langchain [26]’s `LLMMathChain`, which uses an LLM as an agent that interprets input queries and invokes the `numexpr` function with the appropriate formula. This enables the math chain to address a broad spectrum of math problems that are written both in mathematical and verbal form. See Appendix A.4 for more details on experimental setups.

Accuracy and Latency. As shown in the `ParallelQA` row of Table 1, `LLMCompiler` arrives at the final answer with an average speedup of $2.15\times$ with `gpt-4-turbo` and $2.27\times$ with `LLaMA-2 70B` by avoiding sequential execution of the dependency graphs. Beyond the latency speedup, we observe higher accuracy of `LLMCompiler` with the `LLaMA-2` model as compared to that of `ReAct`, due to the reasons discussed in Sec. 4.1. Particularly in the `LLaMA-2` experiment, where `LLMCompiler` achieves around a 9% increase in accuracy, we noted that $\sim 20\%$ of the examples experienced repetitive function calls with `ReAct`, aligning with our observations from the accuracy analysis detailed in Appendix A.1. Additionally, a comprehensive analysis of `LLMCompiler`’s failure cases is provided in Appendix A.2, where we note minimal `Planner` failures, highlighting `LLMCompiler`’s effectiveness in breaking down problems into complex multi-task dependencies.

Cost. Similar to Sec.4.1, `LLMCompiler` demonstrates substantial cost reductions of $4.65\times$ and $2.57\times$ compared to `ReAct` and `OpenAI`’s parallel function calling, respectively, as indicated in Table 2. This efficiency stems from `LLMCompiler`’s reduced frequency of LLM invocations, which is also the case with `OpenAI`’s parallel function calling, which is limited to planning out immediate parallelizable tasks, not the entire dependency graph. For example, in Figure 3 (c), `OpenAI`’s method would necessitate three distinct LLM calls for initial search tasks, following math tasks, and the final math task. In contrast, `LLMCompiler` achieves this with a single LLM call, planning all tasks concurrently.

4.3 Parallel Function Calling with Replanning

In the previous sections, we have discussed cases in which dependency graphs can be determined statically. However, there are cases where dependency graphs need to be constructed dynamically depending on intermediate observations. Here, we consider one such dynamic approach in the context of the `Game of 24` with the `Tree-of-Thoughts (ToT)` strategy proposed in [57]. The `Game of 24` is a task to generate 24 using a set of four numbers and basic arithmetic operations. For example, from the numbers 2, 4, 4, and 7, a solution could be $4 \times (7 - 4) \times 2 = 24$. ToT approaches this task through two iterative LLM processes: (i) the thought proposer generates candidate partial solutions by selecting two numbers and applying an operation (e.g. 2, 3, 7 from 2, 4, 4, 7 by calculating $7 - 4$); (ii) the state evaluator assesses the potential of each candidate. Only the promising candidates are then processed in subsequent iterations of the thought proposer and state evaluator until 24 is reached. Details about the `Game of 24` benchmark and the ToT strategy can be found in Appendix A.9.

While ToT achieves significant improvement at solving the `Game of 24`, its sequential, breadth-first search approach through the state tree can be time-consuming. `LLMCompiler` offers a faster alternative by enabling parallel execution of the thought proposer and the subsequent feasibility evaluator, akin to a parallel beam search method.

Experimental Setups. Although `LLMCompiler` offers latency advantages, solving this problem with a single static graph is not feasible, as the `Planner` cannot plan out the thought proposing stage before identifying the selected candidates from the state evaluator of the previous iteration. Consequently, the `Planner` is limited to planning only within one iteration at a time. To address this, we resort to `LLMCompiler`’s replanning capability. In particular, `LLMCompiler` is equipped with three tools: `thought_proposer` and `state_evaluator`, which are both LLMs adapted from the original ToT framework, and `top_k_select`, which chooses the top k candidates from the `thought_proposer` based on the `state_evaluator`’s assessment. After all these tools are executed, `LLMCompiler` can decide to “replan” if no proposal reaches 24, triggering the `Planner` to devise new plans using the shortlisted states from `top_k_select` of the previous iteration. In this way, `LLMCompiler` can dynamically regenerate plans of each iteration, being able to tackle highly complex tasks that require iterative replanning based on the outcomes of previous plans.

To evaluate `LLMCompiler`’s performance on the `Game of 24`, we use 100 different instances of the game. For each problem, we consider the output as successful if its operations are valid and yield 24 while also using the provided numbers exactly once each. Further details on experiment setups are outlined in Appendix A.4.

Success Rate and Latency. In the last two rows of Table 1, we explore the latency and success rate of `LLMCompiler` in comparison to the baseline described in [57] on the Game of 24 benchmark. With the `gpt-4` model, `LLMCompiler` demonstrates a $2.89\times$ enhancement in latency while slightly improving the success rate compared to the baseline. Similarly, when applied with the `LLaMA-2` model, `LLMCompiler` shows a $2.01\times$ improvement in latency, again without compromising on success rate. These results demonstrate not only a significant latency reduction without quality degradation, but also the replanning capability of `LLMCompiler` for solving complex problems.

4.4 Application: `LLMCompiler` in Interactive Decision Making Tasks

In this section, we demonstrate that `LLMCompiler` can explore language-based interactive environments effectively by benchmarking `LLMCompiler` on `WebShop` [56]. More details of the `WebShop` environment are provided in Appendix A.10. As highlighted in [45, 58, 56], `WebShop` exhibits considerable diversity, which requires extensive exploration to purchase the most appropriate item. While recent work feature advanced exploration strategies and show promising results [62, 32], their approaches are largely based on a sequential and extensive tree search that incurs significant latency penalties. Here, `LLMCompiler` showcases an exploration strategy that is both effective and efficient with the use of parallel function calling. Our method enables broader exploration of items in the environment, which improves success rate compared to `ReAct`. At the same time, this exploration can be parallelized, yielding up to $101.7\times$ speedup against the baselines that perform sequential exploration.

Experimental Setups. We evaluate `LLMCompiler` against three baselines on this benchmark, `ReAct` [58], `LATS` [62], and `LASER` [32], using 500 `WebShop` instructions. The evaluation metrics are success rate, average score, and latency (more details in Appendix A.10). For this experiment, `LLMCompiler` is equipped with two tools: `search` and `explore`. The `search` function triggers the model to generate and dispatch a query that returns a list of typically ten items from the `Webshop` environment. The `explore` function then clicks through links for each of the found items and retrieves information about options, prices, attributes, and features that are available. Finally, based on the gathered information, `LLMCompiler` decides on the item that best matches the input instruction for purchasing. Further details on experiments can be found in Appendix A.4.

Performance and Latency. Our approach significantly outperforms all baseline models as shown in Table 3. When using `gpt-3.5-turbo`, `LLMCompiler` achieves a 25.7% and 7.5% improvement in success rate against `ReAct` and `LATS`; with `gpt-4`, our method improves upon `ReAct` and `LASER` by 17.6% and 2.8%, respectively. In terms of latency, `LLMCompiler` exhibits a $101.7\times$ and $2.69\times$ speedup against `LATS` and `LASER`. While we note that `LLMCompiler` execution is slightly slower than `ReAct` on this benchmark, mainly due to the `Planner` overhead, we also highlight that the gains in success rate far outweigh the minor latency penalty.

We further delve into why `LLMCompiler` attains such an improved success rate and score compared to `ReAct`. Based on our observations, we discover that the `ReAct` agent tends to commit to a decision with imperfect information, a scenario that can arise when the agent has not gathered sufficient details about the features and options available for items. This observation was also noted in [45] – without exploring more items in the environment, the agent struggles to differentiate between seemingly similar choices, ultimately failing to make the correct decision. In contrast, `LLMCompiler` undergoes further exploration by visiting all ten items found by `search` and retrieving relevant information about each item. We find that employing an effective search strategy is critical to decision-making tasks such as the `WebShop` benchmark.

The relatively high performance of `LATS` can also be explained in terms of its exploration scheme. In this framework, the agent executes a brute-force search through the state and action space of `Webshop`, exploring as many as 30 trajectories before making the final purchase. While this approach provides richer information for decision-making, the end-to-end execution becomes prohibitively slow.

Finally, we report that the average scores of both `LATS` and `LASER` fall within the range of standard deviation of our method’s. The average score estimates with standard deviation for `LLMCompiler` are 72.1 ± 4.01 and 75.0 ± 1.43 for `gpt-3.5-turbo` and `gpt-4`, respectively. Further note that while the performance differences are marginal, our method exhibits significant execution speedup, $101.7\times$ over `LATS` and $2.69\times$ over `LASER`.

Table 3: Performance and Latency Analysis for WebShop. We evaluate `LLMCompiler` with two models: `gpt-4` and `gpt-3.5-turbo` and compare `LLMCompiler` against three baselines: `ReAct`, `LATS`, and `LASER`. We report success rate and average score in percentage. We reproduce the success rate and average score for `ReAct`, while those for `LATS` and `LASER` are from their papers. `N` denotes the number of examples used for evaluation.

Model	Method	Succ. Rate	Score	Latency (s)	N
gpt-3.5-turbo	ReAct	19.8	54.2	5.98	500
	LATS	38.0	75.9	1066	50
	<code>LLMCompiler</code>	44.0	72.8	10.72	50
	<code>LLMCompiler</code>	45.5	72.1	10.48	500
gpt-4-0613	ReAct	35.2	58.8	19.90	500
	LASER	50.0	75.6	72.16	500
	<code>LLMCompiler</code>	52.8	75.0	26.73	500

5 Conclusions

Existing methods for invoking multiple functions with LLMs resort to sequential and dynamic reasoning, and as a result, they suffer from inefficiencies in latency, cost, and accuracy. As a solution, we introduced `LLMCompiler`, a compiler-inspired framework that enables efficient parallel function calling across various LLMs, including open-source models like `LLaMA-2` and OpenAI’s `GPT` series. By decomposing user inputs into tasks with defined inter-dependencies and executing these tasks concurrently through its `Planner`, `Task Fetching Unit`, and `Executor` components, `LLMCompiler` demonstrates substantial improvements in latency (up to $3.7\times$), cost efficiency (up to $6.7\times$), and accuracy (up to $\sim 9\%$), even outperforming OpenAI’s parallel function calling feature in latency gains. Further exploration that builds upon `LLMCompiler` will enhance thereby revolutionizing the future development of LLM-based applications. We look forward to future work building upon our framework that will improve both the capabilities and efficiencies of LLMs in executing complex, large-scale tasks, thus transforming the future development of LLM-based applications.

Acknowledgements

We appreciate the valuable feedback from Minwoo Kang. We acknowledge gracious support from Furiosa team. We also appreciate the support from Microsoft through their Accelerating Foundation Model Research, including great support from Sean Kuno. Furthermore, we appreciate support from Google Cloud, the Google TRC team, and specifically Jonathan Caton, and Prof. David Patterson. Prof. Keutzer’s lab is sponsored by the Intel corporation, Intel One-API, Intel VLAB team, the Intel One-API center of excellence, as well as funding through BDD and BAIR. We also appreciate support from Samsung including Dongkyun Kim, and David Thorsley. We appreciate great feedback and support from Ellick Chan, Saurabh Tangri, Andres Rodriguez, and Kittur Ganesh. Sehoon Kim and Suhong Moon would like to acknowledge the support from the Korea Foundation for Advanced Studies (KFAS). Amir Gholami was supported through funding from Samsung SAIT. Michael W. Mahoney would also like to acknowledge a J. P. Morgan Chase Faculty Research Award as well as the DOE, NSF, and ONR. Our conclusions do not necessarily reflect the position or the policy of our sponsors, and no official endorsement should be inferred.

References

- [1] <https://github.com/nvidia/tensorrt-llm>.
- [2] <https://huggingface.co/text-generation-inference>.
- [3] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large language models, 2021.
- [4] BIG bench authors. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023.

- [5] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. Graph of thoughts: Solving elaborate problems with large language models, 2023.
- [6] Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John Jumper. Accelerating large language model decoding with speculative sampling, 2023.
- [7] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. 2021.
- [8] Wenhui Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks, 2023.
- [9] Tim Dettmers, Ruslan Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh Ashkboos, Alexander Borzunov, Torsten Hoefler, and Dan Alistarh. Spqr: A sparse-quantized representation for near-lossless llm weight compression, 2023.
- [10] Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot, 2023.
- [11] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. GPTQ: Accurate post-training compression for generative pretrained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- [12] Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. *arXiv preprint arXiv:2211.10435*, 2022.
- [13] Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. Reasoning with language model is planning with world model, 2023.
- [14] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding challenge competence with apps. *NeurIPS*, 2021.
- [15] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [16] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.
- [17] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for NLP. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 2790–2799. PMLR, 09–15 Jun 2019.
- [18] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- [19] Andrej Karpathy. Intro to large language models, 2023.

- [20] Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. In *The Eleventh International Conference on Learning Representations*, 2023.
- [21] Sehoon Kim, Coleman Hooper, Amir Gholami, Zhen Dong, Xiuyu Li, Sheng Shen, Michael W. Mahoney, and Kurt Keutzer. Squeezellm: Dense-and-sparse quantization, 2023.
- [22] Sehoon Kim, Karttikeya Mangalam, Suhong Moon, Jitendra Malik, Michael W. Mahoney, Amir Gholami, and Kurt Keutzer. Speculative decoding with big little decoder, 2023.
- [23] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners, 2023.
- [24] Woosuk Kwon, Sehoon Kim, Michael W. Mahoney, Joseph Hassoun, Kurt Keutzer, and Amir Gholami. A fast post-training pruning framework for transformers, 2022.
- [25] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- [26] Langchain. <https://github.com/langchain-ai/langchain>.
- [27] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning, 2021.
- [28] Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative decoding, 2023.
- [29] Yaobo Liang, Chenfei Wu, Ting Song, Wenshan Wu, Yan Xia, Yu Liu, Yang Ou, Shuai Lu, Lei Ji, Shaoguang Mao, Yun Wang, Linjun Shou, Ming Gong, and Nan Duan. Taskmatrix.ai: Completing tasks by connecting foundation models with millions of apis, 2023.
- [30] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for llm compression and acceleration, 2023.
- [31] Jerry Liu. LlamaIndex, 11 2022.
- [32] Kaixin Ma, Hongming Zhang, Hongwei Wang, Xiaoman Pan, and Dong Yu. Laser: Llm agent with state-space exploration for web navigation, 2023.
- [33] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhunoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback, 2023.
- [34] Xuefei Ning, Zinan Lin, Zixuan Zhou, Zifu Wang, Huazhong Yang, and Yu Wang. Skeleton-of-thought: Large language models can do parallel decoding, 2023.
- [35] OpenAI. Gpt-4 technical report, 2023.
- [36] OpenAI. New models and developer products announced at devday, 2023.
- [37] Charles Packer, Vivian Fang, Shishir G. Patil, Kevin Lin, Sarah Wooders, and Joseph E. Gonzalez. Memgpt: Towards llms as operating systems, 2023.
- [38] Pruthvi Patel, Swaroop Mishra, Mihir Parmar, and Chitta Baral. Is a question decomposition unit all we need? In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4553–4569, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- [39] Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large language model connected with massive apis, 2023.

- [40] Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. Measuring and narrowing the compositionality gap in language models, 2023.
- [41] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*, 2023.
- [42] Yangjun Ruan, Honghua Dong, Andrew Wang, Silviu Pitis, Yongchao Zhou, Jimmy Ba, Yann Dubois, Chris J. Maddison, and Tatsunori Hashimoto. Identifying the risks of lm agents with an lm-emulated sandbox. *arXiv preprint arXiv:2309.15817*, 2023.
- [43] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023.
- [44] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face, 2023.
- [45] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023.
- [46] Yifan Song, Weimin Xiong, Dawei Zhu, Wenhao Wu, Han Qian, Mingbo Song, Hailiang Huang, Cheng Li, Ke Wang, Rong Yao, Ye Tian, and Sujian Li. Restgpt: Connecting large language models with real-world restful apis, 2023.
- [47] Theodore R. Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L. Griffiths. Cognitive architectures for language agents, 2023.
- [48] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu 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. Llama 2: Open foundation and fine-tuned chat models, 2023.
- [49] Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. *arXiv preprint arXiv:2305.04091*, 2023.
- [50] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models, 2023.
- [51] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc., 2022.
- [52] Tomer Wolfson, Mor Geva, Ankit Gupta, Matt Gardner, Yoav Goldberg, Daniel Deutch, and Jonathan Berant. Break it down: A question understanding benchmark. *Transactions of the Association for Computational Linguistics*, 2020.
- [53] Binfeng Xu, Zhiyuan Peng, Bowen Lei, Subhabrata Mukherjee, Yuchen Liu, and Dongkuan Xu. Rewoo: Decoupling reasoning from observations for efficient augmented language models, 2023.

- [54] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.
- [55] Zonglin Yang, Li Dong, Xinya Du, Hao Cheng, Erik Cambria, Xiaodong Liu, Jianfeng Gao, and Furu Wei. Language models as inductive reasoners, 2022.
- [56] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents, 2023.
- [57] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of Thoughts: Deliberate problem solving with large language models, 2023.
- [58] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- [59] Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and Byung-Gon Chun. Orca: A distributed serving system for {Transformer-Based} generative models. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*, pages 521–538, 2022.
- [60] Wenhao Yu, Meng Jiang, Peter Clark, and Ashish Sabharwal. Ifqa: A dataset for open-domain question answering under counterfactual presuppositions, 2023.
- [61] Huaixiu Steven Zheng, Swaroop Mishra, Xinyun Chen, Heng-Tze Cheng, Ed H. Chi, Quoc V Le, and Denny Zhou. Take a step back: Evoking reasoning via abstraction in large language models, 2023.
- [62] Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language agent tree search unifies reasoning acting and planning in language models, 2023.
- [63] Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations*, 2023.

A Appendix

A.1 Accuracy Analysis: ReAct vs. LLMCompiler

In this section, we delve into a detailed analysis that compares the accuracy of both ReAct and LLMCompiler, highlighting two failure cases that are prevalent in ReAct: (i) premature early stopping, and (ii) repetitive function calls. Furthermore, we demonstrate that while those failure cases negatively impact the ReAct accuracy, they can be effectively addressed by LLMCompiler, thereby yielding the improved accuracy of our framework. We analyze two specific scenarios: the Movie Recommendation evaluation with GPT, where ReAct often prematurely stops, leading to significantly lower accuracy compared to LLMCompiler (68.60 vs. 77.13 in Table 1); and the HotpotQA evaluation with LLaMA-2 70B, where ReAct’s repetitive function calls result in a notable accuracy degradation compared to LLMCompiler (70.00 vs. 77.80 in Table 1).

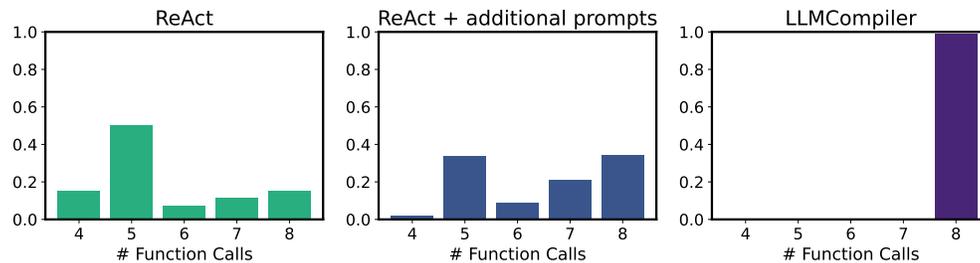


Figure A.1: Distributions of the number of function calls when running the Movie Recommendation benchmark on ReAct (Left), ReAct with specific prompts to avoid early stopping (Middle, corresponding to ReAct[†] in Table 1), and LLMCompiler (Right). LLMCompiler (Right) consistently completes the search for all 8 movies, whereas ReAct (Left) often exit early, demonstrated by about 85% of examples stopping early. Although the custom prompts shift ReAct’s histogram to higher function calls (Middle), they still fall short of ensuring comprehensive searches for all movies. gpt-3.5-turbo is used for the experiment.

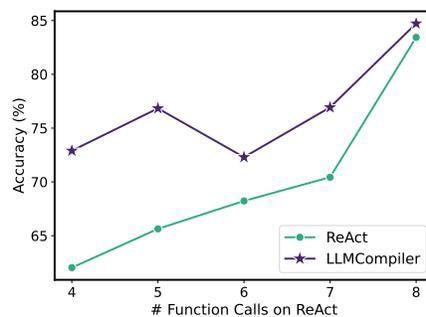


Figure A.2: The Movie Recommendation accuracy of the examples that are categorized by the number of function calls on ReAct, measured both on ReAct and LLMCompiler. The plot indicates that in ReAct, a decrease in the number of function calls correlates with lower accuracy, indicating that premature exits lead to reduced accuracy. In contrast, when the same examples are evaluated using LLMCompiler, which ensures complete searches for all eight movies before reaching a decision, they achieve consistently higher and more consistent accuracy than those processed by ReAct. gpt-3.5-turbo is used for the experiment and the results are averaged over 3 different runs.

Premature Early Stopping of ReAct. ReAct frequently suffers from premature early stopping, ceasing function calls too early, and making decisions based on incomplete information. A clear example of this is observed in the Movie Recommendation benchmark, where ReAct often searches for fewer than the required 8 movies before delivering its final answer. In Figure A.1 (Left), we illustrate the distribution of the number of function calls within ReAct (using GPT) across HotpotQA comparison benchmark. Here, we observe around 85% of the examples exhibit early stopping, making decisions without completing all 8 movie searches. This contrasts with LLMCompiler (Right), where almost all examples (99%) complete the full search of 8 movies. Although adding specific prompts to

prevent early stopping shifts the distribution towards more function calls (Figure A.1, Middle), resulting in an accuracy improvement from 68.60 to 72.47 (ReAct[†] in Table 1), it is nevertheless an imperfect solution.

To further assess how early stopping negatively impacts accuracy, we categorize Movie Recommendation benchmark examples by their number of function calls in ReAct. We then evaluated these groups using `LLMCompiler`, ensuring complete search results for all 8 movies. Table A.2 reveals that fewer function calls in ReAct correlate with lower average accuracy (green line). Conversely, if these examples were processed through `LLMCompiler`, with complete searches for all eight movies, they consistently attained higher accuracy (purple line). This not only indicates that ReAct struggles with premature exits (which is not fully addressed by prompting), but the earlier it stops, the greater the decline in accuracy, contributing to the overall accuracy drop observed in Table 1. In contrast, `LLMCompiler` effectively addresses this issue.

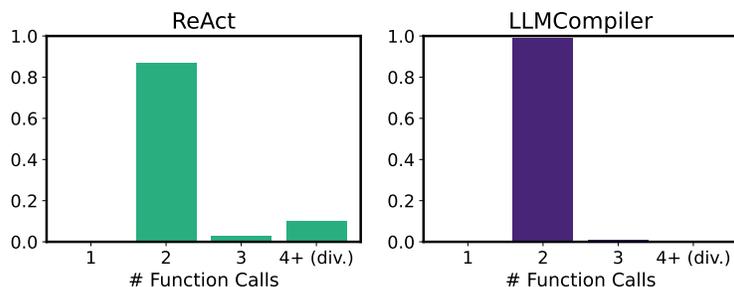


Figure A.3: Distributions of the number of function calls when running the HotpotQA benchmark on ReAct (Left) and `LLMCompiler` (Right). While `LLMCompiler` (Right) consistently completes the task within 2 function calls, which is expected as HotpotQA exhibits a 2-way parallelizable pattern, ReAct (Left) shows that around 10% of the examples undergo repetitive (>4) function calls, resulting in a diverging behavior of the framework. LLaMA-2 70B is used for the experiment.

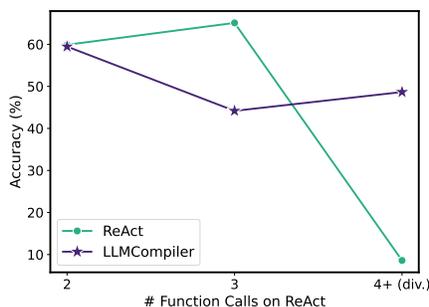


Figure A.4: The HotpotQA accuracy of the examples that are categorized by the number of function calls on ReAct, measured both on ReAct and `LLMCompiler`. The plot indicates that in ReAct, repetitive function calls of more than or equal to four times can result in a significant accuracy degradation due to its infinite looping and diverging behavior. On the other hand, when the same examples are evaluated using `LLMCompiler`, which ensures only two searches per example, they achieve a higher of around 50%. LLaMA-2 70B is used for the experiment.

Repetitive Function Calls of ReAct. Another common failure case of ReAct is its tendency for repetitive function calls, often leading to infinite loops or exceeding the context length limit. This problem is particularly noticeable in the HotpotQA benchmark where ReAct repeatedly calls the same function if the Wikipedia search returns insufficient information about the searched entity. Although HotpotQA is inherently 2-way parallelizable, as illustrated in Figure A.3, we observe that about 10% of its examples require more than four function calls in ReAct, usually resulting in an infinite loop or a divergent behavior. In contrast, `LLMCompiler` executes only two function calls for most examples.

To show how the repetitive function calls impact the overall accuracy, we conduct an accuracy analysis similar to the previous case. In Figure A.4, we categorize HotpotQA benchmark examples by the number of function calls in ReAct and then compare their accuracy on both ReAct and `LLMCompiler`. The analysis reveals that examples that launch two function calls in ReAct maintain the same accuracy in `LLMCompiler`. However, cases with more

Table A.1: A latency comparison between using and not using streaming in the Planner. Streaming yields consistent latency improvement across different benchmarks, as it enables the Task Fetching Unit to start task execution immediately as each task is produced by the Planner. The impact of streaming is especially notable in the ParallelQA benchmark, where tool execution times are long enough to effectively hide the Planner’s execution time.

Benchmark	w/o streaming (s)	w/ streaming (s)	Latency speedup
HotpotQA	4.00	3.95	1.01×
Movie Rec.	5.64	5.47	1.03×
ParallelQA	21.72	16.69	1.30×

than four function calls in ReAct, which often lead to divergent behavior, show less than 10% accuracy in ReAct. On the other hand, when these examples are processed with `LLMCompiler`, they achieve around 50% accuracy by circumventing repetitive calls. It is worth noting that there are instances with three function calls in ReAct, where an extra search can lead to improved accuracy by retrying with an alternate entity name when the initial search fails, yielding a better accuracy than `LLMCompiler`. While this shows a potential adaptability advantage of ReAct, such instances represent less than 3% of cases.

A.2 Failure Case Analysis of `LLMCompiler`

This section delves into a qualitative analysis of `LLMCompiler`’s failure cases on the ParallelQA benchmark, which can be broadly attributed to failures in the Planner, Executor, or the final output process. Failures in the final output process refer to cases when LLMs are unable to use the observations collected from tool execution (which are incorporated into the context) to deliver the correct answer to the user. Among the 10.6% (36 examples) of `LLMCompiler`’s total failures reported in Tab. 1, we have noted that the Planner, Executor, and final output process contributed to 8%, 64%, and 28% of the failures, respectively. The Planner’s 8% failure rate is exclusive to `LLMCompiler`. For instance, the Planner would incorrectly map inputs and outputs by assigning a wrong identifier as an input to a subsequent task, thereby forming an incorrect DAG. However, with adequate tool definitions and in-context examples, Planner errors are significantly reduced (only 3 instances in total throughout our experiment), underscoring the LLM’s capability to decompose problems into complex multi-task dependencies.

The remaining 92% of the total failures are attributed to the Executor and the final output process. The Executor accounts for most of these failures (64%), with common issues like the `math` tool choosing wrong attributes or mishandling unit conversions. For the final output process (28% of failures), errors include incorrect conclusions from the gathered observations, such as failing to pick the smallest attribute from the collected data. It’s worth noting that these problems are not exclusive to `LLMCompiler` but also occur in ReAct. Nevertheless, `LLMCompiler` tends to have slightly fewer failures in these areas than ReAct as it provides only relevant contexts to each tool, aiding in more accurate information extraction. We believe that optimizing the structure of the agent scratchpad, rather than simply appending observations, could further reduce failures in the final output process.

A.3 `LLMCompiler` Details

A.3.1 Streamed Planner

The Planner may incur a non-trivial overhead for user queries that involve a lot of tasks as it blocks the Task Fetching Unit and the Executor, which must wait for the Planner output before initiating their processes. However, analogous to instruction pipelining in modern computer systems, this can be mitigated by enabling the Planner to asynchronously stream the dependency graph, thereby allowing each task to be immediately processed by the Executor as soon as its dependencies are all resolved. In Table A.1, we present a latency comparison of `LLMCompiler` with and without the streaming mechanism across different benchmarks. The results demonstrate consistent latency improvements with streaming. Particularly, in the ParallelQA benchmark, the streaming feature leads to a latency gain of up to 1.3×. This is attributed to the `math` tool’s longer execution time for ParallelQA, which can effectively hide the Planner’s latency in generating subsequent tasks, unlike the shorter execution times of the `search` tool used in HotpotQA and Movie Recommendation.

A.3.2 User-Supplied Information

LLMCompiler requires the following two inputs from the user:

1. **Tool Definitions:** Users need to specify the tools that LLMs can use, including their descriptions and argument specifications. Optionally, users can also provide in-context examples demonstrating the usage of these tools. This is essentially the same requirement as other frameworks like ReAct and OpenAI function calling.
2. **In-context Examples for the Planner:** Optionally, users can provide LLMCompiler with examples of how the Planner should behave. For instance, in the case of Figure 2, users may provide examples illustrating expected inter-task dependencies for certain queries. Such examples can aid the Planner LLM in generating the appropriate dependency graph in the correct format for incoming inputs. In Appendix A.6, we include the examples that we used in our evaluations.

A.4 Experiment Details

Our experiments evaluate two different common scenarios: (1) using API-based closed-source models; and (2) using open-source models with an in-house serving framework. We use OpenAI’s GPT models as closed-source models, in particular, gpt-3.5-turbo (1106 release) for HotpotQA and Movie Recommendation, gpt-4-turbo (1106 release) for ParallelQA, and gpt-4 (0613 release) for Game of 24. Experiments on HotpotQA, Movie Recommendation, and ParallelQA are all conducted in November 2023 after the 1106 release. The Game of 24 experiments are conducted over a two-month period from September to October 2023. For an open-source model, we use LLaMA-2 [48], which was hosted on 2 A100-80GB GPUs using the vLLM [25] framework. All the runs have been carried out with zero temperature, except for `thought_proposer` and `state_evaluator` for the Game of 24 evaluation, where the temperature is set to 0.7. Since OpenAI has randomness in outputs even with temperature 0, we have conducted 3 runs and reported the average accuracy. Across ReAct, OpenAI parallel function calling, and LLMCompiler, we perform 3, 1, and 5-shot learning for HotpotQA, Movie Recommendation, and ParallelQA, respectively; the same examples across different methods were used to ensure a fair comparison. For the Game of 24, we use 2 in-context examples for the Planner. We use the same instruction prompts across different methods for a fair comparison, except for ReAct[†] in Sec. 4.1 with additional ReAct-specific prompts. For WebShop experiment, we use gpt-4-0613 with 8k context window and gpt-3.5-turbo model with 16k context window.

A.5 Analysis

A.5.1 Parallel Speedup Modeling

While LLMCompiler shows noticeable latency gain in various workloads, it is not achieving the $N \times$ latency speedup for N -way parallel workloads. This is mostly due to the overhead associated with LLMCompiler’s Planner and final answering process that cannot be parallelized. In our Movie Recommendation experiment, LLMCompiler’s Planner and the answering process have an overhead of 1.88 and 1.62 seconds on average, respectively, whose combined overhead already comprises more than half of LLMCompiler’s overall latency in Tab 1. Another source of overhead is the straggler effect among the parallel tasks when they need to join together. We observe the average latency of the slowest `search` to be 1.13 seconds which is nearly $2 \times$ the average latency of all tasks, which is 0.61 seconds. Below, we provide an analytical latency modeling of ReAct, LLMCompiler, and LLMCompiler with streaming, and we provide an analysis of achievable latency speedup.

In this section, our focus is on *embarrassingly parallelizable* workload (pattern Figure 3(a)), as this allows for a clearer understanding of the impact of each component on potential latency gains. For the precise latency analysis, we consider three key components: the Planner, the Task Fetching Unit, and the Executor, in Figure 2. Assume that the Planner generates N different tasks to be done. We define P_i as the Planner’s output corresponding to the i -th atomic task. Each P_i is a blueprint for a specific atomic task, which we refer to as E_i . The execution of E_i involves a specific function call using the appropriate tool. The latency function of each unit in the system is defined to quantify the time taken for specific operations. For the Planner, the latency is denoted as $T_P(P_i)$, representing the time taken by the Planner to generate the plan P_i . Similarly, for the Executor, the latency, $T_E(E_i)$, corresponds to the time required to complete the task E_i . We ignore the latency of Task Formulation Unit as it is negligible in this section. Our focus here is on comparing the latency models of ReAct [58], and LLMCompiler.

To begin our analysis of ReAct’s latency, we express its total latency as:

$$T^R = \sum_{i=1}^N \left(T_P^R(P_i) + T_E(E_i) \right). \quad (1)$$

Here, the superscript R refers to ReAct. In the ReAct agent system, the process typically involves initial thought generation, followed by action generation and the acquisition of observations through function calls associated with the tool. The creation of both thought and action are collectively considered as part of generating P_i . It is important to note that while the Planner’s latency is denoted with a superscript (indicating ReAct), the Executor’s latency does not have such a superscript. This is because the function calling and the tools execution remain the same between ReAct and LLMCompiler.

For LLMCompiler, where all parallelizable tasks are processed concurrently, the total latency is determined by the slowest task among these tasks. Hence, the latency model for LLMCompiler can be represented as:

$$T^C = \sum_{i=1}^N T_P^C(P_i) + \max_{k \in \{1, \dots, N\}} T_E(E_k). \quad (2)$$

This expression captures the sum of all planning times plus the execution time of the longest task, reflecting the system’s focus on parallel execution.

Further, if the Planner employs streaming of the dependency graph, the latency model undergoes a modification and can be expressed as:

$$T^{SC} = \sum_{i=1}^N T_P^C(P_i) + T_E(E_N). \quad (3)$$

It is important to note that $T^{SC} \leq T^C$. This implies that the streaming mechanism allows for a more efficient handling of task dependencies, potentially reducing overall latency.

In evaluating the potential speedup achievable with the LLMCompiler framework compared to ReAct, the speedup metric, denoted as γ , is defined as follows:

$$\gamma = \frac{T^R}{T^C} = \frac{\sum_{i=1}^N (T_P^R(P_i) + T_E(E_i))}{\sum_{i=1}^N T_P^C(P_i) + \max_{k \in \{1, \dots, N\}} T_E(E_k)}. \quad (4)$$

This ratio represents the comparative efficiency of LLMCompiler over ReAct, considering both planning and execution latencies.

To estimate the upper bound of this speedup, γ_{\max} , we assume that the executor latency $T_E(E_i)$ is dominant over the planning latency $T_P(P_i)$ and all the latencies of executing tasks remain the same. Under this assumption, the upper bound is calculated as:

$$\gamma_{\max} \approx \frac{\sum_{i=1}^N T_E(E_i)}{\max_{k \in \{1, \dots, N\}} T_E(E_k)} = N, \quad (5)$$

indicating the theoretical maximum speedup, γ_{\max} , is equal to the number of tasks, N .

On the other hand, the lower bound of the speedup, γ , is observed when the planning latency is the predominant factor. Given that the planning latencies of both ReAct and LLMCompiler are generally similar, the minimum speedup is approximated as:

$$\gamma_{\min} \approx \frac{\sum_{i=1}^N T_P^R(P_i)}{\sum_{i=1}^N T_P^C(P_i)} \approx 1. \quad (6)$$

From these observations, we can conclude that to achieve significant latency gains with LLMCompiler, it is crucial to (i) reduce the planner overhead and (ii) minimize the occurrence of stragglers.

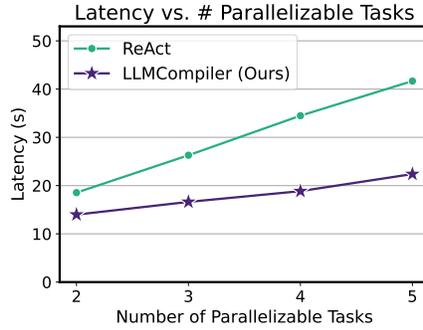


Figure A.5: Latency on the ParallelQA benchmark grouped by the number of maximum parallelizable tasks.

A.5.2 Latency versus Number of Parallelizable Tasks

In Figure A.5, we also report a more detailed latency breakdown on ParallelQA where we show the end-to-end latency as a function of the number of parallel tasks. This is often referred to as weak-scaling in high-performance computing, where the ideal behavior is to have a constant latency as the number of tasks is increased. We can see that ReAct’s latency increases proportionally to the number of tasks which is expected as it executes the tasks sequentially. In contrast, the latency of LLMCompiler increases at a much smaller rate as it can perform multiple function calls in parallel when possible. The reason the end-to-end latency increases slightly with LLMCompiler is due to the overhead of the Planner, which needs to generate plans initially, and which cannot be parallelized. We provide a further analysis of this in Appendix A.5.1.

A.6 User-Supplied Examples for LLMCompiler Configuration

LLMCompiler provides a simple interface that allows for tailoring the framework to different use cases by providing tool definitions as well as optional in-context examples for the Planner. Below, we provide the Planner example prompts that are used to set up the framework for the Movie Recommendation and Game of 24 benchmarks with only a few lines of prompts.

A.6.1 Movie Recommendation Example Prompts

```

Question: Find a movie similar to Mission Impossible, The Silence of the
Lambs, American Beauty, Star Wars Episode IV - A New Hope
Options:
Austin Powers International Man of Mystery
Alesha Popovich and Tugarin the Dragon
In Cold Blood
Rosetta

1. search("Mission Impossible")
2. search("The Silence of the Lambs")
3. search("American Beauty")
4. search("Star Wars Episode IV - A New Hope")
5. search("Austin Powers International Man of Mystery")
6. search("Alesha Popovich and Tugarin the Dragon")
7. search("In Cold Blood")
8. search("Rosetta")
Thought: I can answer the question now.
9. join()
###

```

A.6.2 Game of 24 Example Prompts

```
Question: "1 2 3 4", state_list: [""]
$1 = thought_proposer("1 2 3 4", "")
$2 = state_evaluator("1 2 3 4", "$1")
$3 = top_k_select("1 2 3 4", ["$1"], ["$2"])
$4 = join()
###
Question: "1 2 3 4", state_list: ["1+2=3(left:3 3 4)", "2-1=1(left:1 3
4)", "3-1=2(left:2 2 4)", "4-1=3(left:2 3 3)", "2*1=2(left:2 3 4)"]
$1 = thought_proposer("1 2 3 4", "1+2=3(left:3 3 4)")
$2 = thought_proposer("1 2 3 4", "2-1=1(left:1 3 4)")
$3 = thought_proposer("1 2 3 4", "3-1=2(left:2 2 4)")
$4 = thought_proposer("1 2 3 4", "4-1=3(left:2 3 3)")
$5 = thought_proposer("1 2 3 4", "2*1=2(left:2 3 4)")
$6 = state_evaluator("1 2 3 4", "$1")
$7 = state_evaluator("1 2 3 4", "$2")
$8 = state_evaluator("1 2 3 4", "$3")
$9 = state_evaluator("1 2 3 4", "$4")
$10 = state_evaluator("1 2 3 4", "$5")
$11 = top_k_select("1 2 3 4", ["$1", "$2", "$3", "$4", "$5"], ["$6", "$7",
"$8", "$9", "$10"])
$12 = join()
###
```

A.7 Pre-defined LLMCompiler Planner Prompt

The pre-defined LLMCompiler Planner prompt provides it with specific instructions on how to break down tasks and generate dependency graphs while ensuring that the associated syntax is formatted correctly. This prompt contains specific rules such as assigning each task to a new line, beginning each task with a numerical identifier, and using the \$ sign to denote intermediate variables.

- Each action described above contains input/output types and descriptions.
- You must strictly adhere to the input and output types for each action.
- The action descriptions contain the guidelines. You MUST strictly follow those guidelines when you use the actions.
- Each action in the plan should strictly be one of the above types. Follow the Python conventions for each action.
- Each action MUST have a unique ID, which is strictly increasing.
- Inputs for actions can either be constants or outputs from preceding actions. In the latter case, use the format \$id to denote the ID of the previous action whose output will be the input.
- Ensure the plan maximizes parallelizability.
- Only use the provided action types. If a query cannot be addressed using these, invoke the join action for the next steps.
- Never explain the plan with comments (e.g. #).
- Never introduce new actions other than the ones provided.

A.8 ParallelQA Benchmark Generation

Inspired by the IfQA benchmark [60], our custom benchmark ParallelQA contains 113 examples that are designed to use mathematical questions on factual details of different entities to answer questions, thus requiring a mix of search

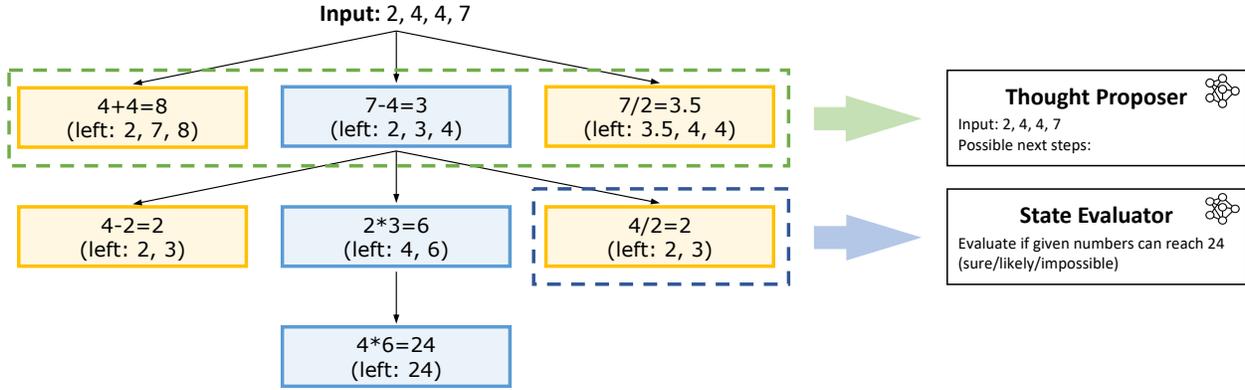


Figure A.6: Visualization of the Tree of Thoughts (ToT) in the Game of 24. Each node represents a distinct proposal, beginning with the root node and branching out through the application of single operations by the thought proposer. Subsequent states are evaluated by the state evaluator for their potential to reach the target number 24. The ToT retains the top-5 states according to their values.

and mathematical operations that are interdependent in various ways. For instance, the benchmark includes examples like “If Texas and Florida were to merge and become one state, as well as California and Michigan, what would be the largest population density among these 2 new states?” requires four parallel search tasks, followed by math tasks dependent on the search outcomes, that can be executed in parallel.

The main objective of the benchmark is to quantify the framework’s ability to decompose an input into multiple tasks to derive an answer. Therefore, we have meticulously selected 56 distinct entities across various domains whose attributes can be accessible from Wikipedia search. By minimizing tool execution (i.e., Wikipedia search) failures, we have aimed our benchmark to effectively assess the frameworks’ abilities to decompose questions into multiple tasks, plan them out, and derive final answers based on observations. Furthermore, to incorporate diverse execution patterns, we crafted various dependency patterns that perform unary and binary math operations after searching for additional information about entities in a given question. We have also curated different questions that accommodate different numbers of maximally parallelizable tasks, ranging from 2 to 5, and we have included varying numbers of joins between parallel function calls as well to increase problem complexity. For instance, we have 2 and 3 joins in Figure 3 (b) and (c), respectively. The benchmark contains 113 different examples, that were populated by GPT-4 based on the aforementioned criteria and labeled by humans afterward.

A.9 Details of the Game of 24 and the Tree-of-Thoughts Approach

The Game of 24 is a mathematical reasoning game that challenges players to manipulate a given set of four numbers, using the basic arithmetic operations of addition, subtraction, multiplication, and division, to arrive at the number 24. The rule of this game is that the given numbers must be used only once. For instance, given the numbers 2, 4, 4, and 7, one possible solution is $4 \times (7 - 4) \times 2 = 24$. This is a non-trivial reasoning benchmark for LLMs, highlighted by the fact that even advanced models like GPT-4 exhibit only a 4% success rate, even when using chain-of-thought prompting [57].

In ToT, the problem is solved in several steps. At each step, the LLM, referred to as the thought proposer, generates thoughts. Each thought is a partial solution that consists of two numbers and an arithmetic operation between them. Then, these thoughts are fed into the state evaluator which assigns a label for each of them. These labels are ‘sure’, ‘likely’, and ‘impossible’, which are given to thoughts to denote how likely they could produce 24 with additional arithmetic operations between the result and the remaining numbers. Only the thoughts that are likely to produce 24 continue onto the next step. This process is illustrated in Figure A.6.

A.10 Details of WebShop Environment

The WebShop environment simulates an online shopping platform. Tasks are designed for the agent to find the item that best matches the given instruction. For instance, if the instruction specifies, “I am looking for a queen-sized bed that is black, and priced lower than 140.00 dollars,” the agent’s task is to pinpoint the bed that precisely fits these criteria: “queen-sized,” “black,” and “priced under 140.00 dollars”. For each item, there is an associated reward

measuring how well this item matches the instruction based on price, item options, and other details contained in the item page. The evaluation metrics are the success rate—the proportion of episodes where the selected product satisfies all requirements—and the average score—the mean reward across episodes.