AGENTRACE: BENCHMARKING EFFICIENCY IN LLM AGENT FRAMEWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Model (LLM) agents are rapidly gaining traction across domains such as intelligent assistants, programming aids, and autonomous decision systems. While existing benchmarks focus primarily on evaluating the effectiveness of LLM agents, such as task success rates and reasoning correctness, the efficiency of agent frameworks remains an underexplored but critical factor for real-world deployment. In this work, we introduce AgentRace, the first benchmark specifically designed to systematically evaluate the efficiency of LLM agent frameworks across representative workloads. AgentRace enables controlled, reproducible comparisons of runtime performance, scalability, communication overhead, and tool invocation latency across popular frameworks on diverse task scenarios and workflows. Our indepth experiments reveal 9 insights and 12 underlying mechanisms for developing efficient LLM agents. We believe AgentRace will become a valuable resource for guiding the design and optimization of next-generation efficient LLM agent systems. The platform and results are available at the anonymous website https://agent-race.github.io/.

1 Introduction

Large Language Models (LLMs) (OpenAI, 2023; Touvron et al., 2023; Liu et al., 2024a; Naveed et al., 2023; Bai et al., 2023) have rapidly gained widespread popularity due to their exceptional capabilities in natural language understanding and generation, significantly impacting various applications including chatbots, content creation, and programming assistants. With these advancements, LLM agents (Wang et al., 2024; Guo et al., 2024; Zhao et al., 2024; Zhang et al., 2024; Ni & Buehler, 2024), which are autonomous entities powered by LLMs capable of executing complex tasks through intelligent interactions, have emerged as a promising area of research and practical implementation.

To accelerate the development of LLM agents, numerous benchmarks and datasets (Andriushchenko et al., 2024; Chang et al., 2024; Huang et al., 2023; Shen et al., 2024) have been proposed to assess LLM agents, primarily focusing on evaluating their effectiveness and reliability in task completion. These benchmarks typically measure task success rates, correctness of generated outputs, overall functional capabilities, and safety of agents.

However, for LLM agents to be widely deployed in real-world scenarios in the future, the efficiency of their frameworks is critically important. Efficient execution, scalability, and minimal communication overhead are essential for ensuring timely responses and practical usability, particularly in resource-constrained and latency-sensitive environments. Despite the proliferation of LLM agent frameworks, such as LangChain (LangChain, 2025), AutoGen (Wu et al., 2023), and AgentScope (Gao et al., 2024), a systematic benchmark evaluating these frameworks' performance efficiency remains absent.

To bridge this significant gap, we introduce **AgentRace**, the first efficiency-focused benchmark platform for LLM agent frameworks, including cost, computational, and communication efficiency. AgentRace enables controlled, reproducible comparisons across frameworks and workflows, aiming to answer the following key research questions:

- 1. What are the primary efficiency bottlenecks in current LLM agent frameworks (e.g., model inference latency, tool calling overhead)?
- 2. What caused the inefficiency of existing LLM agent frameworks?

3. How to improve the efficiency of agent execution?

AgentRace features a modular and extensible design. It supports 7 LLM agent frameworks, 12 types of tools, 3 commonly used workflows, 5 task scenarios, and 4 metrics. The benchmark can be executed with a single command line, facilitating rapid experimentation and reproducibility. We conduct a comprehensive assessment of the efficiency of popular LLM agent frameworks and reveal 9 insights and 12 underlying mechanisms for developing efficient LLM agents. The platform and results are made available through an anonymous website https://agent-race.github.io/.

In summary, our contributions include:

- We introduce AgentRace, the first benchmark platform that systematically evaluates the efficiency of LLM agent frameworks with modular design, filling a critical gap left by existing benchmarks that primarily focus on task success or reasoning correctness.
- We conduct a comprehensive and in-depth assessment of efficiency across frameworks, revealing previously undocumented sources of inefficiency.
- We provide actionable insights for both practitioners and researchers to optimize the deployment of efficient LLM-based agents.
- We release the entire benchmark suite and experimental results, providing a platform to identify the efficiency issues of LLM agents.

2 BACKGROUND AND RELATED WORK

2.1 LLM AGENTS

LLMs agents (Yao et al., 2023; Zhao et al., 2024) are systems that combine the generative capabilities of LLMs with additional components such as memory, planning, and tool usage to perform complex tasks autonomously. These agents can interpret user inputs, plan actions, interact with external tools, and adapt based on feedback, enabling more dynamic and context-aware behaviors. Many agents have been developed, where some are generic agents that are designed to execute general tasks and some are specialized agents for some concrete task. For example, ReAct (Yao et al., 2023) is a typical general agent workflow, where the agent thinks and take actions interatively. MetaGPT (Hong et al., 2023) is an agent designed for software development, where each agent plays a different role to simulate a software company. In this work, we aim to evaluate the efficiency of different LLM agent frameworks, thus focusing on using the widely used general agent workflows.

2.2 LLM AGENT FRAMEWORKS

The development and deployment of LLM agents have been facilitated by various frameworks that provide tools and abstractions for building agentic systems. There have been many LLM agent frameworks. For example, LangChain (LangChain, 2025) offers a modular framework for developing applications with LLMs, supporting integrations with various data sources and tools. It provides a low-level agent orchestration framework, a purpose-built deployment platform, and debugging tools. Besides LangChain, there are also many other popular LLM agent frameworks. In our platform, we select some popular and easy-to-use frameworks for integration. For the detailed introduction of these frameworks, please refer to Section 3.1.

2.3 BENCHMARKS FOR LLM AGENTS

There have been many benchmarks for LLM agents (Andriushchenko et al., 2024; Chang et al., 2024; Huang et al., 2023; Shen et al., 2024; Liu et al., 2024b). However, most of these benchmarks usually focus on ability or trustworthiness perspectives, and do not exploit the efficiency part. For example, AgentBench (Liu et al., 2024b) report *Step Success Rate* as the main metric showing the independent accuracy of each action step, due to the current struggles for LLMs to ensure overall task success rates. Beyond benchmarks focusing solely on success rates, AgentBoard (Chang et al., 2024) proposes a comprehensive evaluation framework for LLM agents. It introduces a fine-grained *Progress Rate* metric to track incremental advancements during task execution, along with an open-source toolkit for multi-faceted analysis. WORFBENCH (Huang et al., 2023) introduces a unified

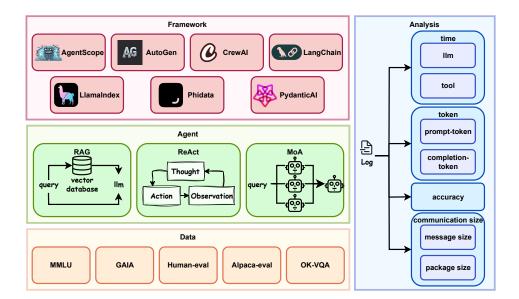


Figure 1: The architecture of AgentRace.

framework for evaluating workflow generation, including both linear and graph-structured workflows. Its evaluation metric, WORFEVAL, quantifies generation performance across these tasks. Although the benchmark measures end-to-end efficiency through *Task Execution Time*, it omits a detailed breakdown of computational costssuch as tool execution latency. This lack of granularity obscures potential bottlenecks in workflow optimization. MASArena (MAS, 2025) provides a convenient multi-dimensional framework for agent evaluation, but it lacks a unified implementation for diverse workflows and heterogeneous tool integrations. Moreover, its evaluation benchmarks are limited to domains such as mathematics, code, and textual reasoning.

3 Design of Agentrace

3.1 Modules

To systematically evaluate the efficiency and scalability of LLM agent frameworks, we introduce a modular benchmark platform AgentRace. As shown in Figure 1, this platform comprises four interconnected modules, including **Data**, **Agent**, **Framework**, and **Analysis**, designed to capture diverse agent frameworks, execution workflows, task complexities, and performance analysis.

Data Module: Diverse Task Coverage The Data module defines the core tasks used in our benchmark and plays a critical role in ensuring that LLM agent frameworks are evaluated across a wide range of real-world scenarios. Our design is guided by two key considerations: (1) task diversity in terms of reasoning complexity, tool usage, and interaction patterns; and (2) alignment with widely adopted benchmarks to enable meaningful and comparable evaluations. We select five representative datasets that reflect varying levels of difficulty, domain coverage, and agent requirements, including GAIA (Mialon et al., 2023), HumanEval (Chen et al., 2021), MMLU (Hendrycks et al., 2020), AlpacaEval (Dubois et al., 2024), and OK-VQA (Marino et al., 2019). The datasets cover tool-intensive, structured reasoning, retrieval-augmented workflows, multi-agent, and multi-modal scenarios. The details about the datasets are available at Appendix A.1. The above coverage enables a holistic evaluation of agent frameworks under varied demands, including tool usage, memory handling, retrieval integration, and inter-agent communication.

Agent Module: Workflow Diversity The Agent module captures the diversity of reasoning patterns exhibited by modern LLM-based agents. In designing this module, our goal is to represent a wide range of real-world task execution strategies while ensuring broad compatibility with existing agent frameworks. We instantiate agents using three widely adopted and conceptually distinct workflow paradigms, including **ReAct (Reasoning and Acting)** (Yao et al., 2023), **RAG (Retrieval-**

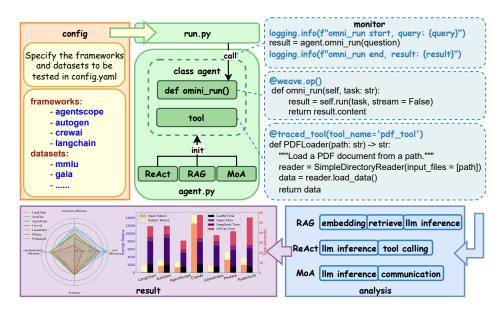


Figure 2: The pipeline of AgentRace.

Augmented Generation), and **MoA** (**Mixture of Agents**) (Wang et al., 2025). These workflows reflect sequential prompting, retrieval-grounded answering, and distributed multi-agent collaboration. By supporting all three within our benchmark, we enable a comprehensive evaluation of agent frameworks under varying reasoning styles and system architectures. The details about the workflows are available at Appendix A.2.

Framework Module: Broad Ecosystem Coverage The Framework module integrates a wide spectrum of open-source LLM agent frameworks including LangChain (LangChain, 2025), AutoGen (Wu et al., 2023), AgentScope (Gao et al., 2024), CrewAI (Lee, 2025), LlamaIndex (LlamaIndex, 2025), Phidata (agno-agi, 2025), and PydanticAI (PydanticAI, 2025), each with distinct design philosophies, runtime environments, and abstraction layers. In selecting the frameworks, we focus on two primary considerations: (1) their popularity and influence in the developer and research communities, and (2) the feasibility of easy deployment and integration within our benchmarking platform. In particular, our implementations are designed to extend functionalities absent from certain frameworks, while leveraging native components whenever available so as not to replace or override existing optimizations.

Analysis Module: Measuring Efficiency The Analysis module defines the core metrics used to evaluate the system-level efficiency of LLM agent frameworks. We focus on three dimensions: computational efficiency, cost efficiency, and communication efficiency. Specifically, we measure the following four key metrics: (1) Execution Time: The total wall-clock time from agent invocation to task completion. This includes the full execution pipeline, including LLM inference, tool calls, etc. (2) Token Consumption: The total number of input and output tokens processed by the LLM during the task. This reflects the computational cost of inference and directly impacts the monetary cost in API-based deployments. (3) Communication Size: The total volume of data exchanged between agents. This metric captures inefficiencies in prompt formatting, serialization, and inter-agent message passing, particularly relevant in multi-agent setting. (4) Accuracy: To ensure correctness is preserved during efficiency evaluation, we also include a task-specific accuracy metric. This ensures that frameworks are functionally correct.

3.2 PIPELINE

The design of the AgentRace benchmark pipeline is illustrated in Figure 2. The pipeline is fully modular and consists of three main stages: (1) configuration, (2) execution and monitoring, and (3) analysis and visualization. In the configuration stage, users specify experimental parameters (e.g., framework, workflow, dataset, and tools) in a YAML file. The executor parses this file and instantiates

Table 1: The supported functionalities of AgentRace. ✓ denotes that the functionality is implemented in AgentRace. ○ denotes that the functionality is supported in the original framework.

		LangChain	AutoGen	AgentScope	CrewAI	LlamaIndex	Phidata	PydanticAI
	ReAct	0	✓	0	✓	0	✓	✓
Workflow	RAG	0	✓	0	√	0	0	✓
	MoA	✓	0	✓	0	0	✓	✓
	Search	0	✓	0	0	0	0	✓
	PDF loader	0	✓	✓	✓	0	/	✓
	CSV reader	0	✓	✓	0	0	0	✓
	XLSX reader	0	✓	✓	/	0	/	✓
	Text file reader	0	✓	✓	0	0	0	✓
Tools	doc reader	0	✓	✓	√	0	1	✓
10018	MP3 loader	0	√	0	1	0	1	✓
	Figure loader	✓	✓	0	0	0	1	✓
	Video loader	✓	✓	✓	✓	1	1	✓
	Code executor	0	0	0	0	0	0	✓
	data retrieval	0	✓	0	✓	0	0	√
	LeetCode solver	✓	✓	1	✓	1	✓	✓

the corresponding agent with unified interfaces. During execution, the agent interacts with the chosen framework and tools under controlled settings, while a monitoring layer is dynamically attached to capture runtime behavior. Finally, the analysis stage aggregates the collected traces into structured logs and performance visualizations for reproducibility and cross-framework comparison.

Tracer and Logger The monitoring layer is designed to provide fine-grained yet low-overhead instrumentation. We implement two complementary components: a *logger* for recording high-level events and a *tracer* for intercepting fine-grained tool calling operations. The logger tracks each LLM inference call, data retrieval request, and inter-agent communication, capturing metadata such as token count, latency, and payload size. To address the scalability challenge of monitoring diverse tool invocations, we introduce a generic tracer wrapper, traced_tool, that instruments tool execution transparently. Developers can annotate a tool with a single wrapper, after which its statistics are automatically recorded. This design allows AgentRace to maintain both extensibilitynew tools can be integrated without modifying the core frameworkand reproducibility, as all traces are stored in a standardized log format compatible with downstream analysis modules.

3.3 Functionalities

The core functionalities supported by AgentRace are summarized in Table 1. Our benchmark currently supports three representative agent workflows executed across seven widely used LLM agent frameworks, utilizing a unified pool of eleven tools. While some of these capabilities are natively supported by the frameworks, approximately 50% of the functionalities are implemented by ourselves to ensure full compatibility and coverage. To maintain a fair comparison across frameworks, we adopt a standardized implementation for any functionality that is not natively provided. This ensures that differences in evaluation metrics stem from the underlying framework behavior, rather than implementation gaps. For more implementation details, please refer to Appendix A.

4 EXPERIMENTS AND INSIGHTS

We conduct in-depth analysis for the efficiency of LLM agent frameworks. Due to the page limit, we present the representative results in the main paper. We present additional experimental details in Appendix A, accuracy comparison in Appendix B.1, results on additional datasets in Appendix B.2, scalability experiments in Appendix B.3, extended analysis between prompt token counts and execution time in Appendix B.4, experiments on additional models in Appendix B.5, reproducibility

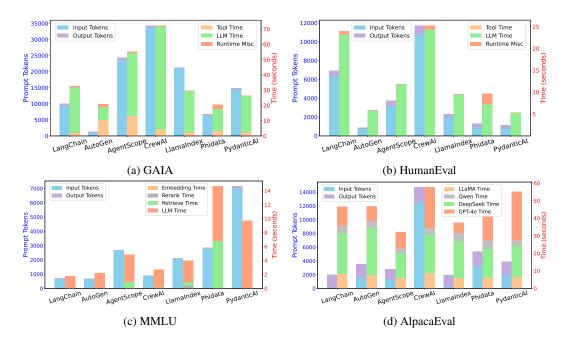


Figure 3: Token consumption and execution time per query of different frameworks.

verification in Appendix B.6, prompts in Appendix C, bugs and features of the investigated frameworks in Appendix D, tool implementation in Appendix E, and usage of LLMs in Appendix F.

4.1 EXPERIMENTAL SETUP

Setting We evaluate 7 LLM agent frameworks using our benchmarking platform, AgentRace, ensuring a standardized and reproducible execution environment. By default, with three repeated runs, experiments are conducted on a Linux server equipped with 12-core Intel(R) Xeon(R) Silver 4214R CPUs and a single NVIDIA RTX 3080 Ti GPU. While most of our metrics and findings are independent of hardware setup, we have also added experiments on additional servers to demonstrate the robustness of our results in Appendix B.6.

Datasets We use five representative datasets across different agent workflows: GAIA, HumanEval and OK-VQA are executed with the ReAct workflow, MMLU is evaluated using RAG, and AlpacaEval is tested under the MoA.

Models Unless otherwise specified, GPT-40 is used as the default LLM. We also conduct experiments using other models in Appendix B.5. For MoA, we instantiate the first-layer agents with a diverse set of open models: LLaMA-3.3-70B-Instruct-Turbo, Qwen2.5-7B-Instruct-Turbo, and DeepSeek-V3. We use TogetherAI (tog, 2024) for querying these models. GPT-40 is used as the aggregation agent to integrate their outputs. In the RAG setting, the MMLU test set is used to construct the retrieval database.

4.2 EXECUTION TIME AND TOKEN CONSUMPTION

Insight 1: LLM inference usually dominates runtime across all agent frameworks, and inefficient prompt engineering, such as appending full histories and using verbose prompts, exacerbates both latency and cost.

Key Observations Figure 3 presents the breakdown of agent execution time across four benchmark scenarios. The results on OK-VQA are available at Appendix B.2. Across all settings, LLM inference consistently dominates runtime. Even in the GAIA scenario, which is explicitly designed to be tool-intensive and involves frequent calls to external APIs, LLM inference accounts for more than 85% of the total execution time in most frameworks. This highlights that LLM inference, due to its

computational demands and frequent invocation, remains the primary bottleneck in agent execution, regardless of the complexity or type of task. Moreover, we observe that the cost of LLM inference is further exacerbated by large variations in token efficiency across frameworks. There is a strong positive correlation between LLM inference time and token consumption.

Underlying Mechanism-1: Appending Full History to Prompts We observe that CrewAI and AgentScope elevate token usage arises from their design choice. In their implementation, the LLM stores all intermediate inputs and outputs in memory and appends this memory to each new prompt. As a result, the prompt length grows with every step of reasoning, causing a high token consumption.

Underlying Mechanism-2: Using Verbose Prompts In the ReAct workflow, LlamaIndex consumes a significant amount of prompts, primarily due to the observation portion returned to the LLM after tool invocation. Additionally, for queries that fail to execute successfully, the number of reasoning and action iterations increases, leading to a corresponding growth in the observation-related prompts. For a more detailed analysis of the underlying causes, please refer to Appendix B.2.

Potential Optimizations These findings underscore the importance of efficient prompt engineering and memory management in agent framework design. Strategies such as selective memory summarization, compact formatting, and prompt compression are crucial for reducing token usage.

4.3 TOOL CALLING

Insight 2: Tool execution efficiency varies widely across frameworks, with search and figure-related tools introducing disproportionately high latency.

Key Observations We analyze the execution cost of various tool types across multiple LLM agent frameworks, as illustrated in Figure 4. The results reveal substantial variation in tool execution efficiency between frameworks, particularly for high-cost operations. Among all tool categories, search and figure-related tools usually incur the highest latency, often dominating total tool execution time within a workflow. For instance, the figure loader takes 2.7 seconds to execute in CrewAI, but exceeds 30 seconds in AgentScope, indicating considerable framework-dependent overhead. In

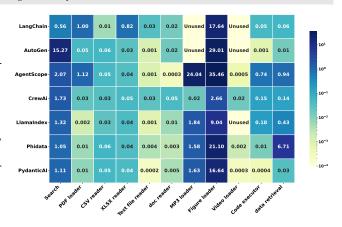


Figure 4: The execution time per call for each tool.

contrast, lightweight tools such as Text file reader and doc reader typically complete in under a millisecond, demonstrating minimal variance.

Additionally, some frameworks (e.g., AgentScope) show disproportionately high total tool processing time, driven primarily by inefficient handling of image processing or multimedia tasks. This highlights the importance of optimizing high-latency tools, particularly in scenarios where tool invocation is frequent or tightly coupled with LLM inference.

Underlying Mechanism-3: Orchestration Depth and I/O Overhead The pronounced disparity in execution times can be attributed to heterogeneous orchestration layers and I/O pathways across frameworks. Heavy operations, especially image-centric routines in figure-related tools, trigger large data transfers and repeated external API calls, amplifying serialization and network overhead. Frameworks with leaner orchestration logic (e.g., CrewAI) perform these steps with fewer intermediate abstractions, thereby reducing latency, whereas frameworks with deeper abstraction stacks (e.g., AgentScope) accumulate additional processing overhead. Consequently, tool latency scales not only with the intrinsic cost of the operation but also with the efficiency of each frameworks data handling, scheduling, and resource management pipelines.

Potential Optimizations While LLM inference remains the dominant bottleneck in most of our benchmarks, more complex, tool-heavy scenarios, such as document analysis or multimodal agent tasks, may shift the performance bottleneck toward tool execution. Frameworks aiming to support such use cases must pay greater attention to optimizing tool orchestration and external API integration.

4.4 RAG

Insight 3: While agents usually involve external databases for information retrieval, the database performance is overlooked in several frameworks. Vector database is recommended.

Key Observations While RAG workflows are increasingly adopted to enhance factual grounding, our benchmarking reveals that database performance, particularly during embedding and retrieval, is a critical yet frequently neglected factor. Figure 3c illustrates the variation in retrieval latency across frameworks, exposing significant performance disparities.

Underlying Mechanism-4: Embedding-Pipeline Design One notable example is AgentScope, which demonstrates high vectorization latency. This stems from its design: during the database setup phase, AgentScope invokes a large embedding model to compute dense vector representations. The latency of this embedding model, often implemented as a separate LLM call, substantially increases the overall vectorization time. Similarly, Phidata exhibits elevated vectorization latency due to its use of a two-step pipeline. First, its built-in csv_tool loads documents row-by-row; then, it applies a SentenceTransformer model to compute embeddings. Our benchmark confirms that Phidata's csv_tool itself is a relatively slow component, compounding the overall vectorization time. From our observation, vector databases such as Faiss (Douze et al., 2024) are good choices.

Potential Optimizations These observations highlight the need for more attention to retrieval pipeline design, especially in frameworks that aim to support real-time or large-scale RAG deployments. Optimization opportunities include batching document embeddings, using faster embedding models, minimizing redundant file reads, and caching frequent queries.

4.5 COMMUNICATION SIZE

Insight 4: Inefficient communication architecture and package design lead to high communication overhead in the multi-agent setting.

Key Observations In multi-agent frameworks, communication between agents is often overlooked as a source of inefficiency. However, our analysis reveals large discrepancies in communication size across frameworks, as shown in Table 2. These differences arise not only from framework-specific message formats but also from architectural design choices.

Underlying Mechanism-5: Inefficient Communication Architecture Frameworks such as CrewAI, which adopt a centralized communication pattern, exhibit significantly higher communication costs. In these designs, a central agent coordinates multiple sub-agents by sequentially delegating subtasks and collecting responses. For example, in CrewAI's MoA implementation, the center agent queries three sub-agents in sequence and aggregates their outputs. Each LLM invocation by the center agent accumulates prior messages in memory, causing the prompt size and the communication payload to grow linearly with the number of sub-agents.

Underlying Mechanism-6: Package Design In addition to the core message, Phidata returns a duplicated content field that mirrors the final message. This, combined with additional metadata fields, results in large communication sizes.

Potential Optimizations These findings indicate that communication cost is not merely a function of task complexity but also of framework design. Future agent frameworks should consider decentralized communication protocols and agent sampling to reduce unnecessary transfer overhead.

Table 2: Communication size between agents (Unit: Byte). We report the content size (e.g., the transferred outputs from the last agent) and overhead size (e.g., header), separated by /.

		LangChain	AutoGen	AgentScope	CrewAI	LlamaIndex	Phidata	PydanticAI
From Global Agent	Agent1 Agent2 Agent3	165.07/0 165.07/0 165.07/0	209.08/44.01 209.08/44.01 209.08/44.01	284.078/0 284.078/0 284.078/0	514.962/0 483.740/0 619.516/0	1180.078/898 1171.078/889 1164.078/882	354.508/0 341.160/0 343.219/0	96.022/0 95.425/0 97.116/0
To Aggregation Agent	Agent1 Agent2 Agent3	1983.02/3 2011.83/3 2072.98/3	2066.04/52.4 2071.24/57.38 2156.04/66.81	1659.318/0 1511.311/0 1889.247/0	2497.929/0 1754.701/0 2151.097/0	2022.417/33.689 2054.878/39.118 2116.377/48.641	6128.259/2639.113 6131.272/2629.426 5715.126/2465.817	2000.542/0 1927.093/0 1892.344/0

Table 3: Scalability Evaluation of AlpacaEval.

	Worker Agents	LangChain	AutoGen	AgentScope	CrewAI	LlamaIndex	Phidata	PydanticAI
	3	36.50	36.85	32.12	64.00	27.32	50.22	46.45
Time (Unit: Second)	6	37.96	47.34	67.61	120.54	36.87	60.42	42.24
	9	47.11	50.84	93.36	212.76	43.85	63.84	110.78
(Unit. Second)	12	59.73	55.60	122.99	218.34	53.77	78.80	111.40
	15	66.08	46.43	153.78	245.26	67.23	83.42	62.13
	3	3516.85	3537.22	2800.75	14732.43	1933.51	5398.71	3894.06
	6	7430.69	7211.57	5143.28	34558.34	3869.52	6940.13	7172.68
Total Token	9	10401.23	10653.76	7547.34	55923.96	5557.50	7785.16	9256.82
	12	13801.78	13692.51	10068.83	61244.79	7190.98	8819.67	9384.31
	15	16894.12	16886.17	12480.56	80200.01	8873.19	9938.26	11170.89
	3	6563.04	6920.56	5912.11	8021.94	9708.91	19013.54	6108.54
Communication Size	6	14029.26	14383.36	10506.82	17863.90	19965.41	21684.95	12206.18
(Unit: Byte)	9	20468.68	22325.97	16275.87	24769.81	29936.97	21320.89	16278.34
(Omt. Dyte)	12	27541.48	28782.73	22032.48	26822.83	39846.67	22383.08	16394.10
	15	34178.20	35606.42	27526.39	30897.88	49926.39	23251.44	19198.06

4.6 SCALABILITY

Insight 5: MoA scalability is governed by agent-invocation policy.

Key Observations We evaluate the scalability of the MoA workflow by increasing the number of worker agents from 3 to 6, 9, 12, and 15, while keeping the additional agents identical in configuration to the original ones. Table 3 reports the results on AlpacaEval. For frameworks such as AgentScope and LangChain, both execution time and token consumption grow almost linearly with the number of worker agents, reflecting sequential scheduling policies. In contrast, frameworks like PydanticAI exhibit a significantly slower growth rate, suggesting a fundamentally different invocation strategy.

Underlying Mechanism-7: Parallel Execution In PydanticAI, the observed runtime is shorter than the aggregate of individual tool and LLM invocation times. This efficiency stems from its parallel execution architecture: agent calls and tool invocations are dispatched asynchronously, allowing multiple operations to overlap in time. As a result, the end-to-end latency is effectively bounded by the slowest operation rather than the sum of all operations.

Potential Optimizations Our analysis indicates that task-level parallelism remains largely underexplored in current frameworks. Incorporating asynchronous scheduling and concurrent invocation can substantially improve scalability in multi-agent workflows, especially under real-world conditions where latency and throughput are critical.

5 CONCLUSION

We introduce AgentRace, a comprehensive benchmark platform for evaluating the efficiency of LLM agent frameworks. AgentRace covers a diverse set of datasets, agent workflows, and frameworks, enabling a fair and reproducible comparison across real-world scenarios. Through extensive and in-depth experiments, we reveal key insights and underlying mechanisms. These findings highlight critical optimization opportunities in the design and deployment of LLM-based agents. We hope AgentRace provides a guideline for future work in developing efficient, scalable, and robust agent systems, and we plan to continuously extend the benchmark as the LLM agent ecosystem evolves.

Reproducibility Statement We have provided our code on an anonymous website https://agent-race.github.io/. We have also provided the experimental details in Appendix A and reproducibility verification in Appendix B.6.

REFERENCES

- Together.ai. https://www.together.ai/, 2024. Accessed: 2024-07-16.
- Masarena, 2025. URL https://lins-lab.github.io/MASArena/. Accessed: 2025-09-23.
 - agno-agi. Phidata, 2025. URL https://docs.phidata.com/introduction. Accessed: 2025-05-15.
 - Maksym Andriushchenko, Alexandra Souly, Mateusz Dziemian, Derek Duenas, Maxwell Lin, Justin Wang, Dan Hendrycks, Andy Zou, Zico Kolter, Matt Fredrikson, et al. Agentharm: A benchmark for measuring harmfulness of llm agents. *arXiv preprint arXiv:2410.09024*, 2024.
 - Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
 - Ma Chang, Junlei Zhang, Zhihao Zhu, Cheng Yang, Yujiu Yang, Yaohui Jin, Zhenzhong Lan, Lingpeng Kong, and Junxian He. Agentboard: An analytical evaluation board of multi-turn llm agents. *Advances in Neural Information Processing Systems*, 37:74325–74362, 2024.
 - Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
 - Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. *arXiv preprint arXiv:2401.08281*, 2024.
 - Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.
 - Dawei Gao, Zitao Li, Xuchen Pan, Weirui Kuang, Zhijian Ma, Bingchen Qian, Fei Wei, Wenhao Zhang, Yuexiang Xie, Daoyuan Chen, et al. Agentscope: A flexible yet robust multi-agent platform. *arXiv preprint arXiv:2402.14034*, 2024.
 - Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges. *arXiv preprint arXiv:2402.01680*, 2024.
 - Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
 - Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 3(4):6, 2023.
 - Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. Benchmarking large language models as ai research agents. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023.
 - LangChain. Langchain, 2025. URL https://www.langchain.com/. Accessed: 2025-05-15.
 - Zeping Lee. GB/T 7714-2015 BibTeX Style. https://github.com/zepinglee/gbt7714-bibtex-style, 2025. GitHub repository.
 - Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33: 9459–9474, 2020.

- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. arXiv preprint
 arXiv:2412.19437, 2024a.
 - Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. In *ICLR*, 2024b.
 - LlamaIndex. Llamaindex, 2025. URL https://www.llamaindex.ai/. Accessed: 2025-05-15.
 - Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pp. 3195–3204, 2019.
 - Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: a benchmark for general ai assistants. In *The Twelfth International Conference on Learning Representations*, 2023.
 - Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language models. *arXiv preprint arXiv:2307.06435*, 2023.
 - Bo Ni and Markus J Buehler. Mechagents: Large language model multi-agent collaborations can solve mechanics problems, generate new data, and integrate knowledge. *Extreme Mechanics Letters*, 67:102131, 2024.
 - OpenAI. Gpt-4 technical report, 2023.
 - PydanticAI. Pydanticai: A python agent framework for generative ai, 2025. URL https://ai.pydantic.dev/. Accessed: 2025-05-15.
 - Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang, Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng Li, and Yueting Zhuang. Taskbench: Benchmarking large language models for task automation. *Advances in Neural Information Processing Systems*, 37:4540–4574, 2024.
 - Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023.
 - Junlin Wang, Jue WANG, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances large language model capabilities. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=h0ZfDIrj7T.
 - Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345, 2024.
 - Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversation. *arXiv preprint arXiv:2308.08155*, 2023.
 - Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.

Yusen Zhang, Ruoxi Sun, Yanfei Chen, Tomas Pfister, Rui Zhang, and Sercan Arik. Chain of agents: Large language models collaborating on long-context tasks. *Advances in Neural Information Processing Systems*, 37:132208–132237, 2024.

Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. Expel: Llm agents are experiential learners. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 19632–19642, 2024.

A EXPERIMENTAL DETAILS

A.1 DETAILS ABOUT THE DATASETS

We select five representative datasets that reflect varying levels of difficulty, domain coverage, and agent requirements: (1) GAIA (Mialon et al., 2023): A comprehensive benchmark for generalpurpose AI assistants. GAIA includes real-world, multi-hop queries that require reasoning over documents, tool invocation, and web interaction. It is the most tool-intensive dataset in our suite, designed to assess the full-stack capabilities of LLM agents. Notably, GPT-4 with plugins achieves only 15% accuracy, while humans reach 92%, indicating significant headroom for improvement. (2) HumanEval (Chen et al., 2021): A code generation benchmark from OpenAI consisting of Python programming problems. Tasks require precise algorithmic reasoning and strict correctness, with deterministic evaluation via unit tests. This dataset helps us evaluate agents capacity for structured reasoning and program synthesis. (3) MMLU (Hendrycks et al., 2020): MMLU spans 57 academic subjects and provides multiple-choice questions across STEM, humanities, and social sciences. We use it to test retrieval-augmented workflows, as it simulates closed-book knowledge challenges and supports grounding in external sources. (4) AlpacaEval (Dubois et al., 2024): An instructionfollowing benchmark that evaluates natural language understanding and response quality. It consists of 805 prompts and uses GPT-4 as a reference evaluator. This dataset is well-suited for multi-agent settings where coordination, aggregation, and language alignment are essential. (5) **OK-VQA** (Marino et al., 2019): A visual question answering benchmark that requires commonsense knowledge beyond images. It contains 14,000 questions over 14,000 images and emphasizes reasoning with external world knowledge. The dataset is for evaluating the efficiency of LLM agent frameworks when handling multimodal tasks.

A.2 DETAILS ABOUT THE WORKFLOWS

AgentRace includes the following workflow paradigms: (1) **ReAct (Reasoning and Acting)** (Yao et al., 2023): This paradigm interleaves natural language reasoning with tool-based actions. By prompting the LLM to first generate intermediate thoughts and then take corresponding actions, ReAct enables agents to dynamically plan and interact with their environment. (2) **RAG (Retrieval-Augmented Generation)** (Lewis et al., 2020): RAG introduces an explicit retrieval step before generation, allowing agents to ground their outputs in relevant external knowledge. In our benchmark, RAG highlights the performance of agent frameworks in integrating retrieval modules, managing memory contexts, and efficiently handling long documents. (3) **MoA (Mixture of Agents)** (Wang et al., 2025): MoA represents a multi-agent architecture where multiple agents collaborate to solve a task. Each agent is often instantiated with a different LLM. An aggregation agent then composes their outputs to form the final answer. This setting captures the growing trend of using multiple LLMs in coordination, and allows us to benchmark frameworks on communication, modularity, and scalability.

A.3 DETAILS ABOUT THE FRAMEWORKS

We integrate the following frameworks: (1) **LangChain** (LangChain, 2025) is a widely adopted framework that offers modular components for building LLM-based applications. It emphasizes tool chaining, prompt templating, memory integration, and external API orchestration. (2) **AutoGen** (Wu et al., 2023), developed by Microsoft, facilitates the creation of advanced LLM agents through multi-agent conversations and automated task planning. (3) **AgentScope** (Gao et al., 2024) supports rapid development of multi-agent systems through a low-code interface. It emphasizes collaboration among agent roles, enabling scalable deployment of agent collectives with minimal boilerplate. (4) **CrewAI** (Lee, 2025) is a lightweight yet expressive Python framework designed for fast iteration. It provides both high-level abstractions and low-level control. (5) **LlamaIndex** (LlamaIndex, 2025) focuses on context-augmented LLM applications by connecting structured and unstructured data sources to LLMs. (6) **Phidata** (agno-agi, 2025) is a framework for building multi-modal AI agents and workflows with memory, knowledge, tools, and reasoning, enabling collaborative problem-solving through teams of agents. (7) **PydanticAI** (PydanticAI, 2025) is an agent framework that is designed for easy development of production-grade applications.

All LLM agent frameworks employed in this study are contemporaneous, with the specific version numbers reported in Table 4.

Table 4: Versions of the LLM Agent frameworks employed in this paper.

Framework	Version	Framework	version
LangChain AutoGen AgentScope CrewAI	0.3.22 0.8.2 0.1.3 0.114.0	LlamaIndex Phidata PydanticAI	0.12.30 2.7.10 0.1.0

A.5 HYPERPARAMETERS

In all experiments, the temperature was set to 0, the top k to 1 (if available), and all other parameters were set to their default values unless otherwise specified.

Except for the cases explicitly noted below, all workflows employ the default prompts provided by their respective frameworks, and the datasets are used without any modification to the original queries.

B ADDITIONAL RESULTS

B.1 ACCURACY

Table 5: Accuracy of each framework on each dataset.

Dataset	LangChain	AutoGen	AgentScope	CrewAI	LlamaIndex	Phidata	PydanticAI
GAIA	0.152 ± 0.012	0.107 ± 0.003	0.212 ± 0.012	0.222 ± 0.009	0.198 ± 0.015	0.191 ± 0.026	0.157 ± 0.012
HumanEval	0.573	0.884	0.884	0.872	0.872	0.902	0.921
MMLU	0.820	0.817	0.827	0.813	0.745	0.792	0.788
OK-VQA	-	0.305	0.436	0.362	0.307	0.331	0.317

Table 5 presents the accuracy of each framework. In general, the accuracy differences among frameworks are relatively small when using the same underlying LLM. However, there are still some notable exceptions.

Insight 6: The complete absence of output constraints in LLMs may lead to tool invocation failures, whereas excessively strict output validation can incur substantial token overhead and decrease the response success rate.

Key Observations In our evaluation, we find that when the model skips tool invocation and instead provides a direct answer (this happens especially with some of the simpler queries in the HumanEval dataset), the framework retries the prompt, often multiple times. Each retry includes previous failed attempts in the context, leading to a rapid increase in prompt length and token consumption as well as a lower likelihood of producing a clean, valid output on later attempts.

Underlying Mechanism-8: Structured Output Misalignment Some frameworks, such as LlamaIndex, require tool inputs to conform to a strict dictionary format. However, GPT-40 does not consistently produce structured outputs that align with these expectations, leading to frequent tool invocation failures. This issue can be partially mitigated if the framework explicitly enforces the format requirement during the registration phase or input schema definition. In contrast, other frameworks such as LangChain adopt stricter enforcement mechanisms. ReAct-style agents in these systems perform rigid output validation and initiate automatic retries when the model's response deviates from the expected invocation structure. While such mechanisms increase robustness against malformed outputs, they may backfire in certain scenarios.

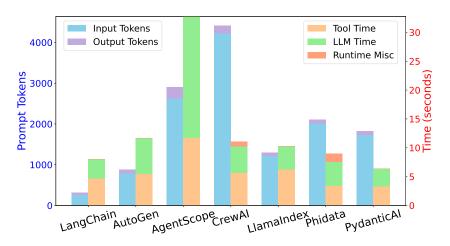


Figure 5: OK-VQA

An additional point to clarify is that the GAIA dataset exhibits relatively low accuracy. This is primarily because GAIA tasks often require complex task planning and the use of multiple tools, posing significant challenges for all evaluated frameworks. It is important to note that the primary focus of this study is not on accuracy, but rather on comparing the performance overhead (e.g., time, token usage) across different frameworks. Therefore, we ensure that the accuracy across frameworks remains broadly comparable, without conducting detailed task-level progress analysis as seen in some related work. By carefully controlling experimental parameters, the fairness of our comparisons remains valid, even in the presence of lower absolute accuracy.

B.2 DETAILED EVALUATION RESULTS

Figure 5 presents the token and time consumption of OK-VQA.

Table 6, 7, 8, 9 and 10 presents the detailed results obtained in this experiment. Unless stated otherwise, the times reported in the table are in seconds per query. The missing data corresponds to instances where the LLM failed to invoke the required tool correctly during the experimentfor example, by not returning outputs in the expected format or by not selecting the appropriate tool for invocation. The following are some noteworthy observations.

Table 6: GAIA Detailed Results

		LangChain	AutoGen	AgentScope	CrewAI	LlamaIndex	Phidata	PydanticAI
Token	Prompt Output Total	9358.35 637.92 9996.27	1159.48 180.66 1340.15	23520.479 785.891 24306.37	33621.857 664.511 34286.369	20935.364 304.976 21240.339	6386.667 323.558 6710.224	14459.17 320.588 14779.758
Time	Ilm Search PDF loader CSV reader XLSX reader Text file reader Doc reader MP3 loader Figure loader Video loader Code executor Total tool time	29.491 1.58856 0.02423455 0.00003333 0.06422606 0.0004194 0.00009758 - 0.5345976 - 0.0152988 2.22746732 32.492	8.464 9.4219 0.0009297 0.000336 0.002387 0.0002909 0.0002212 	41.17 7.291 0.217 0.000297 0.00405 0.0000193 0.00000883 0.729 4.083 0.0000271 0.752 13.076 55.092	67.68 4.031 0.00965 0.000196 0.00422 0.00123 0.000278 0.000346 0.03164 0.000999 0.09565 4.18 72.195	27.244 1.4399 0.0001352 0.00016616 0.004254 0.00034839 0.0001135 0.03341 0.8767 	14.375 1.83012 0.001147 0.0007207 0.003858 0.0002107 0.000073355 0.03821 1.4065 1.38445E-05 0.003035 3.2839 20.396	23.779 1.2275 0.001395 0.0003148 0.003795 8.6865E-06 0.000056241 0.02965 1.2104 3.1952E-06 0.0001414 2.4732 26.238

Insight 7: Token consumption may vary across frameworks even when executing the same workflow, owing to differences in implementation strategies.

810

811

812 813 814 815 816

817 818 819 820

821 822 823 824

827 828 829 830

833 834 835

831 832

836 837 838

839 840 841 842

843 844 845 846

847 848 849

854 855 856 857

858

859

860

861

862

863

Token Time Framework **Prompt** Output **Total** LLM Code executor **Total** LangChain 6326.36 617.13 6943.49 23.221 0.0034 23.968 AutoGen 767.45 106.34 873.79 5.822 0.0002 5.846 3180.689 3742.207 11.906 AgentScope 561.518 11.738 0.131 CrewAI 10817.65 892.798 11710.45 24.22 0.0258 25.24 LlamaIndex 1985.6 342.793 2328.152 9.52 0.003069 9.611 Phidata 967.329 354.427 1321.756 7.181 9.692 **PydanticAI** 812.951 352.543 1165.494 5.258 0.000007158 5.276

Table 7: HumanEval Detailed Results

Table 8: MMLU Detailed Results

		Token		Time				
Framework	Prompt	Output	Total	LLM	Embedding	Retrieve	Total	
LangChain	701.514	4.035	705.55	1.677	11.833	0.055	1.79	
AutoGen	679.788	3.956	683.744	2.171	6.526	0.015	2.182	
AgentScope	2664.315	2.878	2667.193	3.893	92.472	0.935	4.931	
CrewAI	884.536	13.189	897.724	2.51	7.718	0.14	5	
LlamaIndex	2079.702	50.339	2130.042	3.125	4.931	0.4303	3.575	
Phidata	2797.441	37.347	2834.788	7.849	341.611	6.708	17.014	
PydanticAI	6996.242	170.135	7166.378	9.685	5.977	0.03454	9.824	

Key Observations In the results of ReAct workflow, it can be observed that even when using the same ReAct workflow, AgentScope exhibits a significant discrepancy in token usage between the GAIA and HumanEval datasets, with exceptionally high token consumption on GAIA. This is primarily because AgentScope includes the entire memory of the agent in the prompt during every LLM invocation. As the number of reasoning steps increases, the prompt length grows rapidly. While this issue is less apparent in the relatively simple HumanEval dataset, it becomes prominent in the more complex GAIA tasks.

The high token usage observed in CrewAI's ReAct workflow can be attributed to the same reason. In fact, this issue is even more pronounced in CrewAI than in AgentScope, with significantly elevated token consumption observed across both the GAIA and HumanEval datasets.

Underlying Mechanism-9: Overly Detailed Observations In contrast, the majority of token consumption in LlamaIndex and Pydantic arises from the observation segments returned to the LLM after tool invocations. In the GAIA dataset, where tasks are complex and involve frequent tool usage, this results in substantial prompt token overhead.

There are also some issues observed in the MoA workflow. For example, PydanticAI does not require the invocation of all sub-agents during MoA execution, thereby reducing token consumption and runtime overhead. For further details, please refer to the Insight 8 in Appendix B.3.1.

Another example is that in the CrewAI framework, MoA is centrally managed by a global agent, which also plays the role of aggregation agent. The global agent receives the task and sequentially assigns it to sub-agents (e.g., agent1, agent2, agent3). Each sub-agent completes its part and returns the result to the global agent, which then decides the next step. After all agents have responded, the global agent summarizes the results and outputs the final answer. In this setup, the global agent calls the LLM multiple timesonce after each sub-agents response. Because LLMs retain the full context of previous inputs and outputs in a single session, each new call includes all prior interactions. This leads to token accumulation, especially by the third or fourth step, where the prompt becomes much longer. As a result, total token usage becomes higher than in frameworks with different coordination

Table 9: AlpacaEval Detailed Results

			LangChain	AutoGen	AgentScope	CrewAI	LlamaIndex	Phidata	PydanticAI
Token	llama	prompt output total	70.49 428.55 499.04	70.49 431.96 502.45	85.451 382.45 467.901	298.25 518.95 817.201	70.49 430.216 500.707	118.846 438.078 556.924	61.347 429.543 490.889
	qwen	prompt output total	64.84 446.05 510.91	64.85 447.45 512.31	61.815 311.109 372.924	258.083 398.618 656.702	64.81 441.738 506.548	93.899 463.795 557.694	41.217 433.739 474.957
	deepseek	prompt output total	38.5 501.11 539.61	38.5 503.37 541.87	52.478 416.639 469.117	313.01 571.79 884.808	38.485 495.306 533.791	83.391 440.691 524.082	31.802 434.81 485.612
	gpt	prompt output total	1522.48 444.81 1967.29	1529.96 450.63 1980.59	1138.243 352.564 1490.807	11694.576 679.15 12373.72	42.083 350.386 392.47	3003.319 756.689 3760.009	1845.724 596.876 2442.6
Time	llar qw deeps aggre tot	en seek gator	8.275 4.48 23.084 10.699 36.502	7.812 3.977 26.745 8.274 36.854	6.063 3.415 13.726 8.89 32.119	8.835 3.837 21.946 23.114 64	6.069 4.787 20.829 5.849 27.318	6.152 4.707 16.456 14.208 50.217	6.503 3.441 17.79 27.486 46.45
Communication	prompt to prompt to prompt to agent1 to a agent2 to a agent3 to a	agent2 agent3 aggregator aggregator	165.07/0 165.07/0 165.07/0 1983.02/3 2011.83/3 2072.98/3	209.08/44.01 209.08/44.01 209.08/44.01 2066.04/52.24 2071.24/57.38 2156.04/66.81	284.078/118 284.078/118 284.078/118 1659.318/124 1511.311/122 1889.247/126	514.962/0 483.740/0 619.516/0 2497.929/0 1754.701/0 2151.097/0	1180.078/898 1171.078/889 1164.078/882 2022.417/33.689 2054.878/39.118 2116.377/48.641	354.508/0 341.160/0 343.219/0 6128.259/2639.113 6131.272/2629.426 5715.126/2465.817	96.022/0 95.425/0 97.116/0 2000.542/0 1927.093/0 1892.344/0

Table 10: OK-VQA Detailed Results

		token		I	time				
Framework	prompt output total		total	llm	code executor	total			
LangChain	261.033 ± 0.462	52.567 ± 0.473	313.633 ± 0.058	2.948 ± 0.354	4.716 ± 0.115	7.664 ± 0.426			
AutoGen	791.133 ± 0.635	89.467 ± 1.882	880.600 ± 1.609	6.197 ± 0.060	5.171 ± 0.233	11.368 ± 0.240			
AgentScope	2621.367 ± 30.029	283.433 ± 3.355	2902.567 ± 30.346	15.537 ± 5.753	9.043 ± 3.031	24.580 ± 8.779			
CrewAI	4510.933 ± 254.635	269.600 ± 120.951	4780.600 ± 318.263	4.657 ± 0.121	5.578 ± 1.236	10.990 ± 1.536			
LlamaIndex	1219.300 ± 1.682	83.833 ± 0.208	1303.167 ± 1.909	5.548 ± 1.486	5.476 ± 0.627	11.024 ± 0.998			
Phidata	2019.167 ± 2.401	88.500 ± 0.600	2107.500 ± 2.193	4.132 ± 0.054	3.930 ± 0.403	9.039 ± 0.027			
PydanticAI	1728.367 ± 1.674	92.100 ± 0.608	1820.433 ± 2.223	3.034 ± 0.012	3.352 ± 0.057	6.390 ± 0.025			

or memory strategies. This phenomenon will become more apparent in Scalability part as the number of sub agents increases. For further details, please refer to the Insight 4 in Section 4.5.

B.3 SCALABILITY

THE NUMBER OF WORKER AGENTS

To evaluate the scalability of the MoA workflow, we increase the number of worker agents from 3 to 6, 9, 12, and 15, while keeping the newly added agents identical in configuration to the original ones. Metrics from agents using the same LLM are aggregated for reporting. To clearly illustrate how efficiency evolves with increasing numbers of worker agents, we list separate tables (Table 11, 12, 13, 14, 15, 16, 17) for each framework.

B.3.2 THE NUMBER OF TOOLS

Insight 8: Increasing the number of tools has only a minimal impact on execution time across frameworks, but it leads to a noticeable variation in LLM token usage and can cause execution failures when the input exceeds the LLMs maximum context length.

Key Observations We conduct scalability experiments on the GAIA dataset, examining the effect of varying the number of tools across different frameworks. In addition to each frameworks original tool set, we introduce extra LeetCode-solving tools, which are irrelevant for solving the GAIA dataset. The results in Table 18 and 19 show that while increasing the number of tools has only a minimal impact on execution time, it leads to a noticeable increase in LLM token usage. In addition, it can be observed that as the number of tools increases, some test samples encountered execution failures because the input exceed the LLMs maximum context length (see Table 20). Notably, in

918 919 920

Table 11: Scalability Evaluation of AlpacaEval Using AgentScope

938939940941

942 943

937

951

960961962963

964

965966967

968

969

970

971

958 959

Number of Worker Agent 15 6 12 85.451 137.84 206.76 275.68 344.6 prompt 382.45 1204.91 1641.18 2021.9 llama output 796.68 467.901 total 934.52 1411.67 1916.86 2366.5 61.815 89.92 134.88 179.84 224.8 prompt 311.109 555.47 848.94 1139.47 1497.77 qwen output Token 645.39 983.82 1319.31 1722.57 372.924 total 52.478 71.74 107.61 143.48 179.35 prompt deepseek output 416.639 841.37 1253.25 1704.54 2100.42 total 469.117 913.11 1360.86 1848.02 2279.77 1138.243 2237.83 3351.55 4542.02 5677.57 prompt 352.564 412.43 439.44 442.62 434.15 gpt output total 1490.807 2650.26 3790.99 4984.64 6111.72 6.063 12.76 19.307 25.547 35.311 llama 3.415 6.523 10.819 13.866 18.237 gwen Time 13.726 84.318 32.81 48.833 67.114 deepseek 8 89 15 468 14 33 14 373 15 813 gpt total 32.119 67.607 93.357 122,987 153.784 prompt to agent1 284.078/118 389.8/236 584.7/354 779.6/472 974.5/590 prompt to agent2 284.078/118 389.8/236 584.7/354 779.6/472 974.5/590 prompt to agent3 284.078/118 389.8/236 584.7/354 779.6/472 974.5/590 Communication 1659.318/124 3256.270/250 4960.120/375 6718.820/500 8266.330/625 agent1 to aggregator agent2 to aggregator 1511.311/122 2375.700/246 4051.120/369 5477.260/492 7080.860/615 1889.247/126 3705.450/254 5510.530/381 7497.600/508 9255.700/635 agent3 to aggregator

Table 12: Scalability Evaluation of AlpacaEval Using AutoGen

Number of	Worker Age	ent	3	6	9	12	15
	llama	prompt output total	70.49 431.96 502.45	104.14 1004.94 1109.08	158.76 1526.56 1685.32	211.68 2028.21 2239.89	264.6 2529.62 2794.22
Token	qwen	prompt output total	64.85 447.45 512.31	93.18 993.87 1087.05	140.88 1532.12 1673	187.84 1940.46 2128.3	234.8 2419.98 2654.78
	deepseek	prompt output total	38.5 503.37 541.87	40.68 1109.77 1150.45	62.61 1686.42 1749.03	83.48 2249.68 2333.16	104.35 2802.48 2906.83
	gpt	prompt output total	1529.96 450.63 1980.59	3194.7 670.29 3864.99	4830 716.41 5546.41	6290.22 700.94 6991.16	7807.46 722.88 8530.34
Time	llama qwen deepseek gpt total		7.812 3.977 26.745 8.274 36.854	14.667 12.653 46.011 19.816 47.339	25.424 21.064 71.345 22.41 50.843	34.833 28.736 71.98 30.398 55.6	37.816 35.71 104.207 17.817 46.428
Communication	prompt to prompt to prompt to agent1 to a agent2 to a agent3 to a	o agent2 o agent3 oggregator oggregator	209.08/44.01 209.08/44.01 209.08/44.01 2066.04/52.24 2071.24/57.38 2156.04/66.81	236.48/86.04 236.48/86.04 236.48/86.04 4618.13/103.41 4450.9/112.37 4604.89/128.62	359.7/129.06 359.7/129.06 359.7/129.06 7069.64/156.52 6777.28/172.84 7399.95/204.09	479.6/172.08 479.6/172.08 479.6/172.08 9297.61/208.1 8661.68/217.75 9384.64/258.11	599.5/215.1 599.5/215.1 599.5/215.1 11541.91/258.55 10768.69/271.29 11497.32/317.86

the LlamaIndex framework, the addition of the extra LeetCode-solving tools results in a significant decrease in both token consumption and execution time.

Underlying Mechanism-10: Reduced Tool-Call Tendency Increasing the size of the tool inventory paradoxically reduces the agents propensity to invoke tools. On the same test set, adding 10 or 20 LeetCode-solving tools raises the number of queries that make no tool calls from 17 (no extras) to 27 and 25, respectively. Consistent with this shift, the total tool-call counts drop from 630 (0 extra tools) to 454 and 467 (10 and 20 extra tools). These patterns indicate a shallower ReAct trajectory, which in turn reduces LLM token consumption and overall execution time.

Table 13: Scalability Evaluation of AlpacaEval Using LangChain

Number of	Worker Age	nt	3	6	9	12	15
	llama	prompt output total	70.49 428.55 499.04	105.84 1054.54 1160.38	158.76 1518.52 1677.28	211.68 2037.28 2248.96	264.60 2537.08 2801.68
Token	qwen	prompt output total	64.84 446.05 510.91	93.92 1007.95 1101.87	140.88 1446.68 1587.56	187.84 2017.43 2205.27	234.80 2436.53 2671.33
	deepseek prompt output total		38.50 501.11 539.61	41.74 1132.22 1173.96	62.61 1677.97 1740.58	83.48 2224.98 2308.46	104.35 2792.75 2897.10
	gpt	prompt output total	1522.48 444.81 1967.29	3300.82 693.66 3994.48	4734.44 661.37 5395.81	6353.31 685.78 7039.09	7823.07 700.94 8524.01
Time	qwo deeps gp	llama qwen deepseek gpt total		12.061 10.838 40.801 13.741 37.958	19.123 16.584 66.156 17.592 47.112	23.213 24.812 73.476 33.688 59.725	34.437 29.335 115.888 32.068 66.075
Communication	prompt to agent1 prompt to agent2 prompt to agent3 agent1 to aggregator agent2 to aggregator agent3 to aggregator		165.07/0 165.07/0 165.07/0 1983.02/3 2011.83/3 2072.98/3	153.76/0 153.76/0 153.76/0 4703.67/6 4334.61/6 4529.70/6	230.64/0 230.64/0 230.64/0 6787.84/9 6286.30/9 6702.62/9	307.52/0 307.52/0 307.52/0 9117.26/13 8621.19/13 8880.47/13	384.4/0 384.4/0 384.4/0 11314.20/17 10546.46/17 11164.34/17

Table 14: Scalability Evaluation of AlpacaEval Using PydanticAI

Number of	Worker Age	nt	3	6	9	12	15
	llama	prompt output total	61.347 429.543 490.889	95.5 938.08 1033.58	126.29 1273.35 1399.64	139.77 1327.71 1467.48	161.71 1559.8 1721.51
Token	qwen	prompt output total	41.217 433.739 474.957	58.39 939.44 997.83	76.32 1213.31 1289.63	80.85 1257.55 1338.4	94.52 1608.87 1703.39
	deepseek	prompt output total	31.802 434.81 485.612	41.44 931.31 972.75	50.5 1210.15 1260.65	50.88 1150.95 1201.83	58 1311.62 1369.62
	gpt	prompt output total	1845.724 596.876 2442.6	3531.53 636.99 4168.52	4673.28 633.62 5306.9	4739.51 637.09 5376.6	5684.52 691.85 6376.37
Time	qwe deeps gp	llama qwen deepseek gpt total		15.15 8.38 33.34 22.05 42.24	16.68 11.2 42.14 90.94 110.78	19.71 11.59 40.71 91.35 111.4	21.15 13.41 47.19 41.02 62.13
Communication	prompt to prompt to agent1 to a agent2 to a	prompt to agent1 prompt to agent2 prompt to agent3 agent1 to aggregator agent2 to aggregator agent3 to aggregator		88.12/0 93.84/0 94.73/0 4154.19/0 4002.04/0 3773.26/0	113.86/0 118.13/0 108.99/0 5693.77/0 5302.46/0 4941.13/0	124.09/0 119.19/0 103.12/0 6003.71/0 5314.6/0 4729.39/0	134.34/0 131.11/0 113.92/0 6851.79/0 6682.15/0 5284.75/0

Potential Optimizations Building on these findings, agent frameworks should emphasize relevanceaware tool-set curation and dynamic exposure to tools to contain prompt growth and reduce the risk of context-length failures. Regulating ReAct depth and enforcing explicit token budgets can curb

Table 15: Scalability Evaluation of AlpacaEval Using CrewAI

Number of	Worker Age	nt	3	6	9	12	15
	llama	prompt output total	298.25 518.95 817.201	536.95 1186.13 1723.09	706.39 1495.89 2202.26	760.54 1597.78 2358.31	795.95 1741.84 2537.63
Token	qwen	prompt output total	258.083 398.618 656.702	432.87 862.44 1309.12	565.05 1123.05 1688.11	571.23 1088.7 1650.93	589.12 1119.49 1708.61
	deepseek	prompt output total	313.01 571.79 884.808	432.87 1007.86 1440.73	526 1147.04 1673.04	544.75 1181.72 1726.48	668.04 1436.84 2104.88
	gpt	prompt output total	11694.576 679.15 12373.72	28948.53 1136.86 30085.4	49040.19 1320.35 50360.55	54145.65 1363.42 55509.07	72234.23 1614.66 73848.9
Time	llan qwe deeps gp tota	en eek t	8.835 3.837 21.946 23.114 64	20.9 7.7 32.49 53.26 120.54	32.04 16.64 48.37 101.92 212.76	44.25 13.84 50.72 102.374 218.34	27.61 14.61 45.43 159.36 245.26
Communication	prompt to prompt to prompt to agent1 to a agent2 to a agent3 to a	agent2 agent3 ggregator ggregator	514.962/0 483.740/0 619.516/0 2497.929/0 1754.701/0 2151.097/0	925.12/0 912.35/0 900.54.5/0 5921.52/0 4421.22/0 4783.14/0	1425.23/0 1252.74/0 1386.75/0 7929.36/0 6342.21/0 6433.52/0	1724.32/0 1328/0 1327.32/0 8623.56/0 7021.42/0 6798.21/0	1963.23/0 1456.32/0 1587.73/0 9765.36/0 8126.57/0 7998.67/0

Table 16: Scalability Evaluation of AlpacaEval Using LlamaIndex

Number of	Worker Age	ent	3	6	9	12	15
	llama	prompt output total	70.49 430.216 500.707	105.84 1007.91 1113.75	158.76 1502.61 1661.37	211.68 2012.48 2224.16	264.6 2501.66 2766.26
Token	qwen	prompt output total	64.81 441.738 506.548	93.92 972.25 1066.17	140.88 1431.39 1572.27	187.84 1914.73 2102.57	234.8 2420.34 2655.14
	deepseek	prompt output total	38.485 495.306 533.791	41.74 1107.88 1149.62	62.61 1695.19 1757.8	83.48 2216.87 2300.35	104.35 2794.66 2899.01
	gpt	prompt output total	42.083 350.386 392.47	24.68 515.31 539.99	24.68 541.38 566.06	24.68 539.22 563.9	24.68 528.1 552.78
Time	llar qw deep: gr tot	en seek ot	6.069 4.787 20.829 5.849 27.318	12.44 10.69 41.18 9.39 36.87	18.98 14.77 61.97 9.66 43.85	25.65 22.23 81.83 10.4 53.77	35.58 27.49 93.12 16.06 67.23
Communication	prompt to prompt to agent1 to a agent2 to a agent3 to a	o agent2 o agent3 aggregator aggregator	1180.078/898 1171.078/889 1164.078/882 2022.417/33.689 2054.878/39.118 2116.377/48.641	2181.8/1796.0 2163.8/1778.0 2149.8/1764.0 4585.09/67.1 4372.32/72.31 4512.6/90.07	3272.7/2694.0 3245.7/2667.0 3224.7/2646.0 6813.56/99.75 6456.6/106.13 6923.71/137.8	4363.6/3592.0 4327.6/3556.0 4299.6/3528.0 9126.32/133.64 8647.86/143.64 9081.69/180.66	5454.5/4490.0 5409.5/4445.0 5374.5/4410.0 11342.05/169.22 10907.85/181.15 11437.99/227.25

Table 17: Scalability Evaluation of AlpacaEval Using Phidata

Number of	Worker Age	nt	3	6	9	12	15
	llama	prompt output total	118.846 438.078 556.924	114.5 555.91 670.41	110.58 551.62 662.2	116.61 576.04 692.65	118.06 603.61 721.67
Token	qwen	prompt output total	93.899 463.795 557.694	87.57 634.08 721.65	83.21 621.29 704.5	90.21 663.48 753.69	91.97 707.2 799.17
	deepseek	prompt output total	83.391 440.691 524.082	76.54 505.74 582.28	72.94 525.82 598.76	77.18 527.8 604.98	78.63 545.07 623.7
	gpt	prompt output total	3003.319 756.689 3760.009	4180.34 785.45 4965.79	5040.94 778.76 5819.7	5973.25 795.1 6768.35	6991.86 801.86 7793.72
Time	llan qwe deeps gp tots	en seek t	6.152 4.707 16.456 14.208 50.217	6.55 6.75 15.43 23.13 60.42	6.55 5.27 16.6 25.68 63.84	9.12 6.09 19.32 31.67 78.8	10.33 6.56 22.07 31.7 83.42
Communication	prompt to prompt to prompt to agent1 to a agent2 to a agent3 to a	agent2 agent3 ggregator ggregator	354.508/0 341.160/0 343.219/0 6128.259/2639.113 6131.272/2629.426 5715.126/2465.817	325.7/0 309.01/0 304.15/0 7105.44/3163.16 7475.54/3354.1 6165.11/2699.97	310.63/0 293.84/0 288.79/0 6961.87/3105.45 7269.53/3267.58 6196.23/2734.51	329.16/0 319.17/0 307.71/0 7291.8/3252.42 7792.87/3505.45 6342.37/2791.33	334.87/0 326.58/0 314.94/0 7582.27/3388.25 8121.06/3656.93 6571.72/2891.08

unnecessary tool exploration, while compact, standardized tool specifications help decouple token usage from catalog size.

Table 18: Effect of LeetCode-solving tools on execution time (seconds)

	LangChain	AutoGen	AgentScope	CrewAI	LlamaIndex	Phidata	PydanticAI
no LeetCode-solving tools	12.86	8.41	19.57	11.87	24.26	10.23	10.31
10 LeetCode-solving tools	11.79	8.58	22.31	10.35	19.47	10.99	8.33
20 LeetCode-solving tools	10.78	8.36	21.95	11.14	20.89	10.98	9.58

Table 19: Effect of LeetCode-solving tools on Token

	no Lee	no LeetCode-solving tools			10 LeetCode-solving tools			20 LeetCode-solving tools		
	Prompt	Output	Total	Prompt	Output	Total	Prompt	Output	Total	
LangChain	7199.33	553.2	7753	11489.89	586.61	12076.50	12779.90	502.75	13282.65	
AutoGen	1195.98	185.19	1381.18	2200.19	191.82	2392.01	3011.2	182.87	3194.07	
AgentScope	17161.55	828.68	17990.23	31878.31	780.23	32658.54	32464.93	804.56	33269.48	
CrewAI	16475.12	582.82	17057.95	11670.07	552.16	12222.23	17398.34	557.75	17956.09	
LlamaIndex	101042.29	729.57	101771.86	35111.65	348.83	35460.48	32899.47	253.21	33152.68	
Phidata	3293.59	270.75	3564.33	4957.96	295.79	5253.75	6104.55	267.34	6371.88	
PydanticAI	13273.91	373.74	13647.66	12356.90	321.95	12678.85	16682.93	324.13	17025.06	

B.4 EXTENDED ANALYSIS ON INSIGHT 1

Our experiments in Section 5.2 reveal a strong correlation between prompt token counts and execution time across frameworks (see Figure 3). This is primarily due to two factors: 1) the frequency of LLM calls and tool invocations per query; 2) memory accumulation across queries.

Figure 6 shows that CrewAI and AgentScope have significantly higher average LLM call frequencies per query (5.33 and 4.78) compared to LlamaIndex, PydanticAI, and Phidata (2.76, 2.79, and 3.38). This difference explains their greater token consumption and longer runtimes, which stem from more frequent LLM calls and the resulting memory accumulation.

During the experiments, we observed the following patterns, indicating that some frameworks invoke tools more frequently than others:

Table 20: Number of Failed Runs

1	1	36
1	1	37
1	1	38
1	1	39

	LangChain	AutoGen	AgentScope	CrewAl	LlamaIndex	Phidata	PydanticAI
no irrelevant tools	0	0	1	1	1	0	1
10 irrelevant tools	0	0	2	1	1	1	1
20 irrelevant tools	0	0	4	3	1	0	1

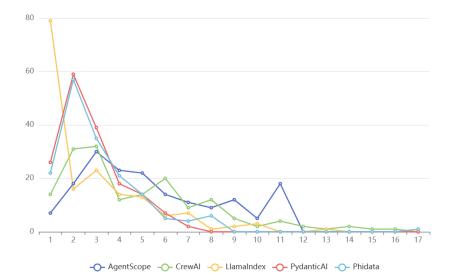


Figure 6: LLM Call Frequency per Query across Different Frameworks

- 1) AgentScope and CrewAI frequently invoke the Web tool to obtain precise results, leading to substantially higher token usage due to lengthy text outputs. In our tests, they called the Web tool 494 and 608 times respectively, far exceeding the maximum of 102 observed in other frameworks.
- 2) AgentScope often writes and executes code to solve problems, which requires returning large code blocks that further increase token usage. It used the code execution tool 122 times, while other frameworks did so no more than 21 times.

Moreover, AgentScope stands out for retaining conversational memory across queries by continuously appending prior interactions to the prompt. Unlike earlier tests that re-instantiated the Agent to avoid memory buildup, running 9 GAIA queries without resets confirmed significant memory accumulation (see Figure 7).4

Meanwhile, in our MoA workflow experiments, we observed that some frameworks invoke worker

agents in parallel, whereas others do so serially. Specifically, we observe that CrewAIs built-in MoA
workflow integrates the previous worker agents output with the initial prompt, performs a secondary
summarization, and then passes the result to the next worker agent. To further explore this behavior,
we varied the order of worker agents in CrewAI and present the results in Table 21. Here, GLM,
Qwen, DS, and GPT denote GLM-Z1-Rumination-32B-0414, Qwen2.5-7B-Instruct, DeepSeek-V3,
and GPT-40, respectively.

Table 21: The Impact of Agent Execution Order on Tokens

	GLMQwenDS			DSQwenGLM			QwenDSGLM		
Order	Prompt	Output	Total	Prompt	Output	Total	Prompt	Output	Total
GLM	1296.82	734.62	2031.44	2953.52	1909.5	4863.02	584.74	324.9	909.64
Qwen	241.86	383.12	624.98	279.96	557.84	837.8	255.92	525.3	781.22
DS	447.0	968.5	1415.5	279.36	568.14	847.5	246.38	556.64	803.02
GPT	36750.26	1119.44	37869.7	36732.26	1129.24	37861.5	17375.24	455.8	17831.04

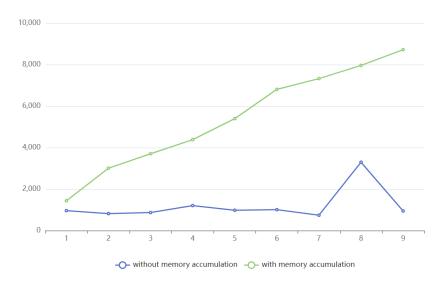


Figure 7: Memory Accumulation Impact

B.5 CLAUDE-BASED RESULTS

Table 22: Claude-Based HumanEval Results

Token				Accuracy			
Framework	Prompt	Output	Total	LLM	Code executor	Total	
LangChain AutoGen	5568.08 920.84	675.88 292.88	6243.96 1213.71	41.644 12.847	0.0140 0.00047	41.932 13.182	0.585 0.823

Given that the majority of our experiments are implemented with GPT-40, and considering the widespread adoption of open-source models, we additionally evaluate the Claude-3-Opus model on the HumanEval dataset within the LangChain and AutoGen frameworks. The results are presented in Table 22.

Notably, AutoGen exhibited slightly lower accuracy compared to GPT-based agents. Upon inspection, we found that Claude did not fabricate test data when invoking the Python execution tool, which rendered the self-checking mechanism ineffective. In LangChain, Claude occasionally emitted tool outputs directly, bypassing the expected format and causing execution failures.

These behaviors suggest that when using Claude-3-Opus as the underlying model for ReAct-style agents, further prompt adaptation may be necessary to ensure compatibility with existing framework toolchains.

B.6 REPRODUCIBILITY VERIFICATION

Table 23: HumanEval Run 2

		Token			Time	
Framework	Prompt	Output	Total	LLM	Code executor	Total
LangChain	6769.16	695.15	7464.31	27.063	0.01267	27.82
AutoGen	790.29	108.26	898.55	5.685	0.000353	5.711
AgentScope	2429.72	530.323	2960.043	13.42	0.121	13.57
CrewAI	10026.98	914.96	10941.95	29.75	0.0432	30.47
LlamaIndex	2052	347.9	2399.9	19.81	0.00381	19.84
Phidata	1083.32	376.46	1459.79	11	8.99E-05	16.3
PydanticAI	903.6	353.48	1257.08	9.13	2.32E-05	9.15

Table 24: HumanEval Run 3

		Token			Time	
Framework	Prompt	Output	Total	LLM	Code executor	Total
LangChain	7953.34	832.63	8785.97	38.562	0.015723	39.471
AutoGen	769.72	105.78	875.5	8.027	0.000279	8.199
AgentScope	2804.341	568.36	3372.701	15.686	0.139	15.858
CrewAI	10822.16	867.08	11689.24	34.19	0.0342	34.98
LlamaIndex	2017.37	362.85	2380.23	20.61	0.00293	20.64
Phidata	1258.7	393.46	1652.16	9.36	0.000227	12.4
PydanticAI	874.49	340.66	1215.15	7.73	2.44E-05	7.74

Table 25: GAIA Run 1

Token				Time						
Frameworks	Prompt	Output	Total	LLM	Search	PDF loader	CSV reader	XLSX reader		
LangChain	6493.9	562.42	7052.33	8.26	0.724	0.000713	2.73E-05	-		
AutoGen	1078.7	183	1261.7	9.65	17.29	0.00347	0.00035	8.91E-05		
AgentScope	19192.78	747.25	19940.02	12.03	1.32	1.48	0.000358	0.00147		
CrewAI	31286.37	612.44	31898.81	34.55	4.66	0.0205	0.000138	0.00272		
LlamaIndex	12370.81	688.83	13059.64	38.4	1.019	0.000618	4.63E-06	0.00196		
Phidata	2387.39	260.78	2648.17	13.16	4.296	0.00257	8.37E-06	8.18E-05		
PydanticAI	15680.58	410.12	16090.7	10.81	0.744	0.461	0.000302	0.000111		

time										
Text file reader	doc reader	MP3 loader	Figure loader	Video loader	Code executor	total tool time	total time			
0.0197	-	-	-	-	0.0176	0.762	10.15			
4.63E-05	5.82E-05	-	-	-	1.15E-05	17.294	27.04			
6.32E-06	2.52E-06	0.125	0.443	2.99E-06	0.996	4.359	16.575			
0.000832	0.00015	0.000375	0.105	-	0.194	4.795	39.86			
0.00113	3.94E-06	3.91E-06	0.839	-	0.387	2.248	47			
4.24E-05	0.000141	0.098	0.075	-	0.000427	4.473	13.16			
0.117	6.33E-05	0.0951	0.141	-	6.39E-05	1.558	11.68			

Table 26: GAIA Run 2

		Token			Time							
Frameworks	Prompt	Output	Total	LLM	Search	PDF loader	CSV reader	XLSX reader				
LangChain	6659.4	598.16	7257.56	17.61	0.78	0.000908	3.82E-05	-				
AutoGen	1063.48	195.52	1259	4.206	11.477	0.000736	0.000223	0.000161				
AgentScope	20787.67	785.02	21572.68	12.997	1.438	2.876	0.000248	0.000841				
CrewAI	33422.3	564.65	33986.94	35.75	4.77	0.0072	0.000146	0.0023				
LlamaIndex	15079.24	731.95	15811.19	35.69	1.196	0.000308	2.19E-06	0.0021				
Phidata	2481.73	279.04	2760.76	5.25	4.055	0.00074	1.37E-05	0.000173				
PydanticAI	11306.87	259.62	11566.48	5.361	1.12	0.535	0.000261	7.93E-05				

Time										
Text file reader	doc reader	MP3 loader	Figure loader	Video loader	Code executor	Total tool time	Total time			
0.0103	_	-	-	_	0.000699	0.797	18.89			
3.39E-05	9.33E-05	-	-	-	1.58E-05	11.478	16.211			
2.60E-06	2.00E-06	0.241	0.406	1.45E-06	0.285	5.248	18.55			
0.000477	0.000147	0.000283	0.0314	-	0.00647	4.82	41.14			
0.00042	9.75E-05	6.96E-06	0.399	-	1.196	2.794	46.28			
0.000166	7.73E-05	0.144	0.108	-	0.000132	4.308	10.69			
0.125	9.10E-05	0.186	0.126	-	1.75E-05	2.091	6.59			

Insight 9: Experimental reproducibility is underpinned by the stability of token usage, while variability arises from stochastic tool behaviors and fluctuating LLM invocation dynamics.

Key Observations To verify the reliability and reproducibility of our results, we conduct repeated experiments on the HumanEval and GAIA datasets. The outcomes are reported in Table 7, 23, 24 for HumanEval and in Table 25, Table 26, Table 27 for GAIA. As illustrated by the error bars in Figure 8

Table 27: GAIA Run 3

Token				Time						
Frameworks	Prompt	Output	Total	LLM	Search	PDF loader	CSV reader	XLSX reader		
LangChain	7262.24	651.28	7913.52	16.86	1.16	0.246	2.55E-05	-		
AutoGen	1067.48	186.24	1253.72	10.59	17.33	0.000685	0.000285	0.000195		
AgentScope	20689.4	761.78	21451.18	21.58	2.446	2.035	0.000199	0.0019		
CrewAI	33866.8	621.44	34488.23	34.15	3.446	0.00617	0.000171	0.00251		
LlamaIndex	19764.47	964	20728.47	61.89	2.395	0.00203	0.000678	0.00631		
Phidata	2187.99	233.53	2421.52	13.81	3.92	0.000728	6.04E-06	0.000103		
PydanticAI	13059.31	296.36	13355.67	15.76	0.783	0.637	3.79E-06	7.88E-05		

	time										
Text file reader	doc reader	MP3 loader	Figure loader	Figure loader	Code executor	Total tool time	Total time				
0.00904	-	-	-	-	0.00125	1.417	18.78				
1.70E-05	2.31E-04	-	-	-	2.00E-05	17.33	28.71				
3.24E-06	4.85E-06	0.164	0.683	4.46E-06	1.88	7.215	29.03				
0.00047	0.000141	0.000283	-	-	0.014	3.47	38.44				
0.000464	0.000239	0.0405	0.69	-	0.307	3.443	74.998				
0.000117	7.83E-05	0.0989	0.0788	-	0.000497	4.1	19.52				
0.0382	5.67E-05	0.0824	0.151	-	5.66E-02	1.75	16.685				

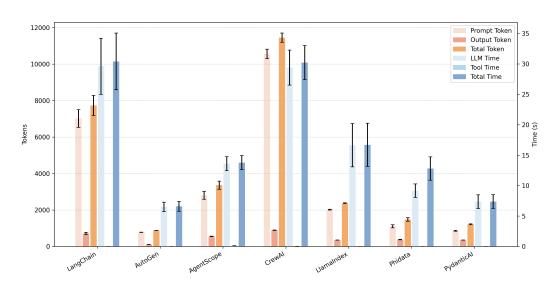


Figure 8: Consistency of Token Consumption and Latency in Repeated Experiments (HumanEval)

and 9, the token consumption in our experiment is relatively stable. In general, the execution time is usually positively related to the token consumption.

Underlying Mechanism-11: Stochastic Tool Behaviors Figure 9 indicates that the LlamaIndex framework yields a relatively high standard deviation on the GAIA dataset. This can be attributed to the stochastic nature of tool invocations and the consequent variations in the number of LLM invocation rounds.

Underlying Mechanism-12: Fluctuating LLM invocation dynamics The inherent randomness of certain LlamaIndex built-in toolssuch as the use of whisper in audio-visual modelsfurther amplifies this effect, resulting in a larger standard deviation in the GAIA test results.

Nevertheless, the overall trend remains reproducible.

In addition, to examine the impact of hardware differences, we rerun the GAIA benchmark on a machine equipped with a 40-series GPU with 48GB of memory. We then compare the results with the average values obtained from the RTX 3080 Ti setup, by computing the ratios of key metrics.

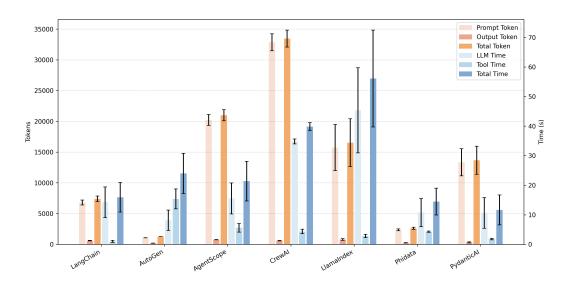


Figure 9: Consistency of Token Consumption and Latency in Repeated Experiments (GAIA)

Table 28: Comparison of Token and Time Ratios Across Hardware Configurations

	T	oken ratio			Time ratio	
Framework	Prompt	Output	Total	LLM	Code executor	Total
LangChain	1.080	1.060	1.083	0.639	0.713	0.638
AutoGen	0.998	1.062	1.008	0.501	1.124	0.923
AgentScope	0.996	1.108	1.000	0.975	0.543	0.858
CrewAI	1.034	0.984	1.033	1.044	1.217	1.076
LlamaIndex	1.354	1.401	1.356	1.014	1.093	1.032
Phidata	0.990	0.963	0.987	1.286	0.961	1.319
PydanticAI	0.912	0.915	0.912	0.629	1.164	0.622

As shown in Table 28, token usage remain largely consistent across most frameworks. Intuitively, the token consumption is independent of hardware setup. In terms of execution time, we observe significant speedup only for LangChain and Pydantic, indicating that these two frameworks benefit more from enhanced GPU capabilities, while others exhibit relatively stable performance regardless of GPU configuration.

C PROMPTS

C.1 REACT

For frameworks that do not have a specific implementation of ReAct, we use the following prompt to build the ReAct workflow:

```
1392
       You are a ReAct-based assistant.
1393
       You analyze the question, decide whether to call a tool or directly
1394
       answer, and then respond accordingly.
1395
       Use the following format: Question: the input question or request
1396 4
       Thought: you should always think about what to do\nAction: the action to
       take (if any)
1397
       Action Input: the input to the action (e.g., search query)
1398
       Observation: the result of the action
1399
       ... (this process can repeat multiple times)
1400 <sub>8</sub>
       Thought: I now know the final answer
1401 9
       Final Answer: the final answer to the original input question or request
      Begin!
1402 10
1403 <sup>11</sup>
       Question: {input}
```

C.1.1 LANGCHAIN

Within the ReAct workflow implemented via LangChain's AgentExecutor, we set the max_iterations parameter to 15 for experiments on the GAIA dataset and to 10 for those on the HumanEval dataset.

1409 C.2 RAG

1415 ₁

1442 1

1420 1

For the following frameworks, we applied specific prompts to improve their token efficiency or to better align with the RAG workflow.

C.2.1 AUTOGEN

You are a helpful assistant. You can answer questions and provide information based on the context provided.

C.2.2 CREWAI

You are a specialized agent for RAG tasks. You just need to give the answer of the question. Don't need any other word. Such as the answer is a number 5 , you need output '5'. Or the answer is A, you need to output 'A'.

C.2.3 PHIDATA

You are a RAG-based assistant. You analyze the question, and call the search_knowledge_base tool to retrieve relevant documents from the knowledge base, and then respond accordingly.

C.2.4 PYDANTICAI

You're a RAG agent. please search information from the given task to build a knowledge base and then retrieve relevant information from the knowledge base.

C.3 MoA

Unless otherwise specified, the following prompt is used for the aggregator agent.

C.3.1 LANGCHAIN

You have been provided with a set of responses from various open-source models to the latest user query. Your task is to synthesize these responses into a single, high-quality response. It is crucial to critically evaluate the information provided in these responses, recognizing that some of it may be biased or incorrect. Your response should not simply replicate the given answers but should offer a refined, accurate, and comprehensive reply to the instruction. Ensure your response is well-structured, coherent, and adheres to the highest standards of accuracy and reliability.

C.3.2 AGENTSCOPE

You are an assistant called Dave, you should synthesize the answers from Alice, Bob and Charles to arrive at the final response.

For the worker agent, we used the following prompt.

1457 | You are an assistant called Alice/Bob/Charles.

C.3.3 CREWAI

1458

1459

1465

1466 1467 1

1468

1469

147014711472

1473

147514761477

1478 1479 ₁

1480

1481 1482 1483

14841485

1486

1487

1488

1489

1490

1491 1492 1493

1494

1495

1496

1497 1498 1499

1501

1502

1503

1474 1

You are an agent manager, and You need to assign the questions you receive to each of your all agents, and summarize their answers to get a more complete answer
You must give question to all the all agents, and you must summarize their answers to get a more complete answer.\nYou need to be the best

For the worker agent, we used the following prompt.

You are one of the agents, you have to make your answers as perfect as possible, there will be a management agent to choose the most perfect answer among the three agents as output, you have to do your best to be selected

C.3.4 PHIDATA

Transfer task to all chat agents (There are 3 agents in your team)", " Aggreagate responses from all chat agents

C.3.5 PYDANTICAI

Your task is to aggregate all agents results to solve complex tasks.\nYou analyze the input, input the task to all tools that can run a single agent, and synthesize the results from all agents into a final response.

C.4 GAIA

In this experiment, we used all levels of questions from the test subset of the GAIA dataset. Below are examples of prompts used in our system, depending on whether a file is attached:

question: A paper about AI regulation originally submitted to arXiv.org in June 2022 features a figure with three axes, each labeled with a pair of opposing terms. Which of these terms is used to describe a type of society in a Physics and Society article submitted to arXiv.org on August 11, 2016?

question: The attached spreadsheet contains the inventory of a movie and video game rental store located in Seattle, Washington. What is the title of the oldest Blu-Ray listed in this spreadsheet? Return it exactly as it appears., file_name: 32102e3e-d12a-4209-9163-7b3a104efe5d.xlsx, file_path: path/to/32102e3e-d12a-4209-9163-7b3a104efe5d.xlsx

C.5 HUMANEVAL

To avoid generating explanatory text or pseudo-code that hinders automated accuracy evaluation, we slightly modify the original HumanEval queries by adding minimal prompts.Below is an example of the prompt used for HumanEval problems:

```
1504
       from typing import List
1505 2
       def has_close_elements(numbers: List[float], threshold: float) -> bool:
1506 3
           """ Check if in given list of numbers, are any two numbers closer to
1507 4
       each other than
1508
           given threshold.
1509
           >>> has_close_elements([1.0, 2.0, 3.0], 0.5)
1510 <sub>7</sub>
           False
           >>> has_close_elements([1.0, 2.8, 3.0, 4.0, 5.0, 2.0], 0.3)
1511 8
           True
```

```
1512

1513 10

1514 12

# Complete the function. Only return code. No explanation, no comments, no markdown.
```

C.6 MMLU

For the MMLU dataset, we constructed the vector database used in the RAG workflow based on the development subset and evaluated the performance of each framework using the test subset. Given the large number of tasks in this dataset, we used only one-quarter of them in our experiments. Considering that tasks from the same domain tend to be spatially adjacent in the dataset, we selected one out of every four tasks in index order. This sampling strategy ensures broader domain coverage and maintains fairness in the evaluation.

Below is an example of the question in MMLU:

```
1527 | Question:Find the degree for the given field extension Q(sqrt(2), sqrt(3), sqrt(18)) over Q.

1528 | A.0 | B.4 | C.2 | D.6 | Answer with A, B, C, or D only
```

C.7 ALPACAEVAL

In this experiment, we used the full set of tasks for the basic MoA experiments, and the first 100 tasks for extended experiments involving more agents. Below is an example of one such task.

D BUGS AND FEATURES

This section summarizes the bugs or features of LLM agent frameworks that we discovered during our evaluation.

D.1 LANGCHAIN

As shown in Figure 10, LangChain's high level of abstraction and encapsulation posed challenges in measuring specific metrics during our experiments.

Additionally, LangChain occasionally terminated processes prematurely after reading files from the GAIA dataset, returning the file content directly rather than proceeding with the expected operations (see Figure 11).

D.2 AUTOGEN

Due to the default system prompt being relatively long and containing irrelevant instructions, the RAG workflow may consume unnecessary tokens or produce unexpected errors (e.g., attempting to invoke non-existent tools). Therefore, it is necessary for users to customize the system prompt.

D.3 AGENTSCOPE

AgentScopes image and audio processing tools internally rely on OpenAI models, causing their execution time to partially overlap with that of the LLM itself. This overlap can lead to inflated or inaccurate measurements of LLM processing time. Researchers and practitioners should be mindful of this issue when conducting time-based evaluations involving AgentScope.

```
def openai_image_to_text(
   image_urls: Union[str, list[str]],
   api_key: str,
```

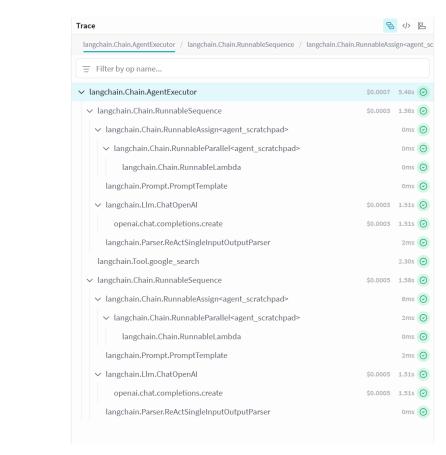


Figure 10: LangChain's high level of abstraction and encapsulation.

2025-05-12 2348:31.145 - root - INFO - tool name: pdf tool, tool time: 3.8433|
2025-05-12 2348:31.175 - root - INFO - tool name: pdf tool, tool time: 3.8433|
2025-05-12 2348:31.175 - root - INFO - tool name: pdf tool, tool time: 3.8433|
2025-05-12 2348:31.175 - root - INFO - tool name: pdf tool, tool time: 3.8433|
2025-05-12 2348:31.175 - root - INFO - tool name: pdf tool, tool time: 3.8433|
2025-05-12 2348:31.175 - root - INFO - tool name: pdf tool, tool time: 3.8433|
2025-05-12 2348:31.175 - root - INFO - tool name: pdf tool, tool time: 3.8433|
2025-05-12 2348:31.175 - root - INFO - tool name: pdf tool, tool time: 3.8433|
2025-05-12 2348:31.175 - root - INFO - tool name: pdf tool, tool time: 3.8433|
2025-05-12 2348:31.175 - root - INFO - tool name: pdf tool time: pdf tool t

Figure 11: LangChain occasionally terminated processes prematurely.

```
prompt: str = "Describe the image",
model: Literal["gpt-40", "gpt-4-turbo"] = "gpt-40",

1614 6
) -> ServiceResponse:
""

Generate descriptive text for given image(s) using a specified model,
and
return the generated text.

Args:
image_urls (`Union[str, list[str]]`):
```

```
1620
                       The URL or list of URLs pointing to the images that need to
1621
       be
1622<sub>14</sub>
                       described.
1623 15
                  api_key (`str`):
                      The API key for the OpenAI API.
1624 16
                  prompt (`str`, defaults to `"Describe the image"`):
1625 <sup>17</sup>
                       The prompt that instructs the model on how to describe
1626 19
                       the image(s).
1627<sub>20</sub>
                  model ('Literal["gpt-40", "gpt-4-turbo"]', defaults to '"gpt-40"')
1628
1629<sup>21</sup>
                       The model to use for generating the text descriptions.
1630 <sup>22</sup>
             Returns:
1631 <sub>24</sub>
                  `ServiceResponse`:
1632<sub>25</sub>
                       A dictionary with two variables: `status` and `content`.
1633 26
                       If `status` is `ServiceExecStatus.SUCCESS`,
                       the `content` contains the generated text description(s).
1634<sup>27</sup>
1635 <sup>28</sup>
             Example:
1636 <sub>30</sub>
1637<sub>31</sub>
                  .. code-block:: python
1638 32
                       image_url = "https://example.com/image.jpg"
1639 33
                       api_key = "YOUR_API_KEY"
1640 34
                       print(openai_image_to_text(image_url, api_key))
1641 36
1642<sub>37</sub>
                  > {
                          'status': 'SUCCESS',
1643 38
                  >
                          'content': "A detailed description of the image..."
                  >
1644<sup>39</sup>
1645 40
             11 11 11
1646<sub>42</sub>
             openai_chat_wrapper = OpenAIChatWrapper(
1647<sub>43</sub>
                  config_name="image_to_text_service_call",
1648 44
                 model_name=model,
                 api_key=api_key,
1649<sup>45</sup>
1650 46
             )
    47
             messages = Msg(
1651<sub>48</sub>
                 name="service_call",
1652<sub>49</sub>
                 role="user",
1653 50
                 content=prompt,
                 url=image_urls,
1654 51
1655 52
             )
    53
             openai_messages = openai_chat_wrapper.format (messages)
1656 54
             try:
1657 55
                 response = openai_chat_wrapper(openai_messages)
                 return ServiceResponse(ServiceExecStatus.SUCCESS, response.text)
1658 56
1659<sup>57</sup>
             except Exception as e:
1660<sup>58</sup>
                 return ServiceResponse(ServiceExecStatus.ERROR, str(e))
1661<sub>60</sub>
        def openai_audio_to_text(
1662<sub>61</sub>
             audio_file_url: str,
1663 62
             api_key: str,
            language: str = "en",
1664 63
1665<sup>64</sup>
            temperature: float = 0.2,
        ) -> ServiceResponse:
    65
1666 <sub>66</sub>
             11 11 11
1667<sub>67</sub>
             Convert an audio file to text using OpenAI's transcription service.
1668 68
1669<sup>69</sup>
1670 70
                  audio_file_url (`str`):
                       The file path or URL to the audio file that needs to be
1671 <sub>72</sub>
                       transcribed.
1672 <sub>73</sub>
                  api_key (`str`):
                       The API key for the OpenAI API.
167374
                  language (`str`, defaults to `"en"`):
    75
```

```
1674 <sub>76</sub>
                        The language of the input audio. Supplying the input language
1675
          in
1676 77
                         [ISO-639-1] (https://en.wikipedia.org/wiki/List_of_ISO_639-1
         codes)
1677
                        format will improve accuracy and latency.
1678<sup>78</sup>
1679 79 80
                   temperature ('float', defaults to '0.2'):
                        The temperature for the transcription, which affects the
1680<sub>81</sub>
                        randomness of the output.
1681 82
1682 83
             Returns:
                   `ServiceResponse`:
1683<sup>84</sup>
                        A dictionary with two variables: `status` and `content`.
1684 85
                        If `status` is `ServiceExecStatus.SUCCESS`,
1685<sub>87</sub>
                        the `content` contains a dictionary with key 'transcription'
1686
        and
                        value as the transcribed text.
168788
1688<sup>89</sup>
1689<sub>91</sub>
             Example:
1690<sub>92</sub>
                   .. code-block:: python
1691 93
                        audio_file_url = "/path/to/audio.mp3"
169294
                        api_key = "YOUR_API_KEY"
1693<sup>95</sup>
1694 96
97
                        print(openai_audio_to_text(audio_file_url, api_key))
1695<sub>98</sub>
                   > {
1696 99
                           'status': 'SUCCESS',
                   >
                           'content': {'transcription': 'This is the transcribed text
1697100
        from
1698
1699<sub>102</sub>
                   the audio file.' }
                   > }
1700<sub>103</sub>
1701<sub>104</sub>
             try:
1702105
                   import openai
              except ImportError as e:
1703<sup>106</sup>
1704<sup>107</sup><sub>108</sub>
                   raise ImportError(
                        "The `openai` library is not installed. Please install it by
1705
1706109
                         "running `pip install openai`.",
                   ) from e
1707110
1708<sup>111</sup>
              client = openai.OpenAI(api_key=api_key)
1709<sup>112</sup>
             audio_file_url = os.path.abspath(audio_file_url)
1710<sub>114</sub>
              with open(audio_file_url, "rb") as audio_file:
1711<sub>115</sub>
                   try:
                        transcription = client.audio.transcriptions.create(
171216
                             model="whisper-1",
1713<sup>117</sup>
1714<sup>118</sup><sub>119</sub>
                              file=audio_file,
                              language=language,
1715<sub>120</sub>
                              temperature=temperature,
1716<sub>121</sub>
1717122
                        return ServiceResponse (
1718<sup>123</sup>
                              ServiceExecStatus.SUCCESS,
1719<sup>124</sup><sub>125</sub>
                              {"transcription": transcription.text},
1720<sub>126</sub>
                   except Exception as e:
1721<sub>127</sub>
                        return ServiceResponse (
                              ServiceExecStatus.ERROR,
1722128
1723<sup>129</sup>
                              f"Error: Failed to transcribe audio {str(e)}",
1724<sup>130</sup>
```

Meanwhile, AgentScope's vector database module, LlamaIndexKnowledge, is implemented based on the BM25Retriever from the llamaindex library. However, the original implementation relies on an

1725 1726

1727

outdated version of llamaindex, and recent updates to the library introduced structural changes that break compatibility with the original import statements.

```
from llama_index.retrievers.bm25 import BM25Retriever
```

To ensure a consistent environment without modifying the frameworks built-in vector database logic, we resolved the issue by duplicating the LlamaIndexKnowledge module and updating the import paths to match the newer llamaindex version.

```
from llama_index.legacy.retrievers.bm25_retriever import BM25Retriever
```

D.4 CREWAI

When our MOA invokes a large number of agents (>=12), CrewAI system occasionally fails to call all agents completely during execution as intended. For example, when we request 12 sub-agents to be activated, some queries may only trigger 9 or fewer agents.

D.5 LLAMAINDEX

What integer-rounded percentage of the total length of the harlequin shrimp recorded in Omar Valencfia-Mendez 2017 paper was the sea star fed to the same type of shrimp in G. Curt Fiedler's 2002 paper?

Response



Thought: The question requires specific data from two scientific papers: Omar Valencia-Mendez's 2017 paper and G. Curt Fiedler's 2002 paper. I need to find the total length of the harlequin shrimp recorded in the 2017 paper and the size of the sea star fed to the shrimp in the 2002 paper. Then, I can calculate the percentage and round it to the nearest integer.

```
Action: google_search
Action Input: "Omar Valencia-Mendez 2017 harlequin shrimp total length"
```

Figure 12: LlamaIndex frequently fails to invoke tools correctly.

As shown in Figure 12, LlamaIndex frequently fails to invoke tools correctly, primarily due to the lack of prompt constraints and insufficient post-processing checks on LLM outputs. Without explicit guidance and validation mechanisms, the LLM often produces outputs that do not conform to the expected dictionary format, resulting in tool invocation failures.

D.6 PHIDATA

As shown in Figure 13, in the ReAct workflow, Phidata passes the available tools to the LLM via the "tools" field. Unlike Llamaindex, which emphasizes the functionality and usage of tools in the system prompt, Phidata rarely invokes the code execution tool when processing queries from humaneval.

D.7 PYDANTICAI

By further visualizing the experimental data (see Figure 14), we found that within the PydanticAI ReAct framework, the same tool was often invoked simultaneously multiple times, potentially leading to inefficiencies. Additionally, similar to Phidata, the code execution tool was seldom triggered (see Figure 15).

Furthermore, The MoA implementation in the PydanticAI framework is tool-based, and not all three models are invoked for every query. We observe that when the number of sub-agents is 3, 6, 9, 12, and 15, there were 232, 89, 229, 485, and 663 instances, respectively, where sub-agents were not invoked. These skipped invocations are randomly distributed across different queries, resulting in lower token consumption than expected.

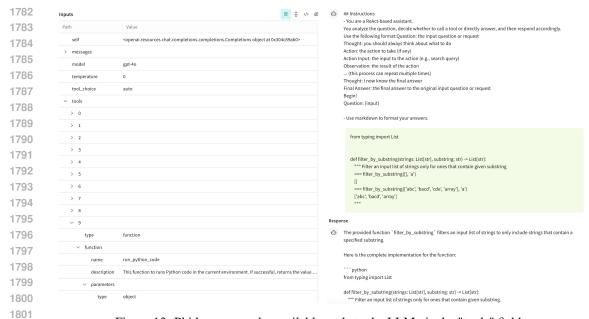


Figure 13: Phidata passes the available tools to the LLM via the "tools" field.

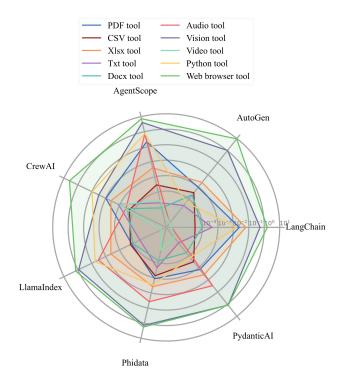


Figure 14: Visualization of the average execution time per run of different tools across different frameworks.

E TOOL IMPLEMENTATION

For frameworks that do not include the required tools, we adopted a unified implementation as follows.

```
2025—65-12 0129723,091 — root — 1MFO — omni_rum start, query; question: What mainsls that were mentioned in both Illis Lagbowardos's and Olga Tapia's papers on the alvei species of the genus named for Copenhag 2025—65-12 0129724,334 — root — 1MFO — LUM name: gpt—do, prompt Losens: 200, point Losens: 1025, 1d: 01506646-6566—7313-3662-0000660000000, limitestump: 1746984563,595017 2025—65-12 0129724,334 — root — 1MFO — LUM name: gpt—do, prompt Losens: 80, completion, tokens: 203, total_tokens: 203, total_
```

Figure 15: PydanticAI's simultaneous invocations of the same tool.

E.1 SEARCH

1836 1837 1838

1843 1844 1845

1847

E.1.1 AUTOGGEN

```
1849 1
        def google_search(query: str, num_results: int = 2, max_chars: int = 500)
         -> list: # type: ignore[type-arg]
1850
             import os
1851 <sup>2</sup>
             import time
1852
             import requests
1853 <sub>5</sub>
             from bs4 import BeautifulSoup
1854 6
             from dotenv import load_dotenv
             load_dotenv()
1855 <sup>7</sup>
             google_api_key = os.environ['GOOGLE_KEY']
1856 <sup>8</sup>
             search_engine_id = os.environ['GOOGLE_ENGINE']
1857 10
             if not search_engine_id or not search_engine_id:
1858<sub>11</sub>
                  raise ValueError ("API key or Search Engine ID not found")
1859 12
             url = "https://www.googleapis.com/customsearch/v1"
             params = {
1860 <sup>13</sup>
                  "key": google_api_key,
1861 <sup>14</sup>
                  "cx": search_engine_id,
1862 16
                  "q": query,
1863<sub>17</sub>
                  "num": num_results
1864 18
             }
             response = requests.get(url, params=params) # type: ignore[arg-type]
1865 19
             if response.status_code != 200:
1866 <sup>20</sup>
1867<sup>21</sup><sub>22</sub>
                  print (response.json())
                  raise Exception(f"Error in API request: {response.status_code}")
1868<sub>23</sub>
             results = response.json().get("items", [])
             def get_page_content(url: str) -> str:
1869 24
                  try:
1870<sup>25</sup>
1871 <sup>26</sup>
                       response = requests.get(url, timeout=10)
                       soup = BeautifulSoup(response.content, "html.parser")
1872 <sub>28</sub>
                       text = soup.get_text(separator=" ", strip=True)
1873 <sub>29</sub>
                       words = text.split()
                       content = ""
1874 30
1875 31
                       for word in words:
1876 32
                            if len(content) + len(word) + 1 > max_chars:
1877 33
34
                                break
                           content += " " + word
1878 35
                       return content.strip()
1879 36
                  except Exception as e:
                       print(f"Error fetching {url}: {str(e)}")
1880<sup>37</sup>
1881 <sup>38</sup>
                       return ""
1882<sub>40</sub><sup>39</sup>
             enriched_results = []
             for item in results:
1883<sub>41</sub>
                  body = get_page_content(item["link"])
1884 42
                  enriched_results.append(
1885<sup>43</sup>
1886 <sup>44</sup>
                            "title": item["title"],
                            "link": item["link"],
1887<sub>46</sub>
                            "snippet": item["snippet"],
1888<sub>47</sub>
                            "body": body
1889 48
                       }
```

```
1890 <sub>50</sub>
                      time.sleep(1)
1891<sub>51</sub>
                return enriched_results
1892
```

E.1.2 PYDANTICAI

1893

```
1894
1895
        def google_search(query, num=None):
1896
1897
             Make a query to the Google search engine to receive a list of results.
1898
1899 4
             Aras:
                  query (str): The query to be passed to Google search.
1900 <sup>5</sup>
                  num (int, optional): The number of search results to return.
1901 <sup>6</sup>
        Defaults to None.
1902 <sub>7</sub>
1903<sub>8</sub>
             Returns:
                 str: The JSON response from the Google search API.
1904 9
1905^{\,10}
1906 . 11
             Raises:
    12
                 ValueError: If the 'num' is not an integer between 1 and 10.
1907
1908<sub>14</sub>
             try:
                  QUERY_URL_TMPL = ("https://www.googleapis.com/customsearch/v1?key
1909 15
        ={key}&cx={engine}&q={query}")
1910
                 url = QUERY_URL_TMPL.format(
1911 16
                      key=os.environ['GOOGLE_KEY'],
1912<sub>18</sub>
                      engine=os.environ['GOOGLE_ENGINE'],
1913<sub>19</sub>
                      query=urllib.parse.quote_plus(str(query))
191420
                  if num is not None:
1915<sup>21</sup>
                      if not 1 <= num <= 10:
1916<sup>22</sup>
1917<sup>23</sup>
                           raise ValueError("num should be an integer between 1 and
        10, inclusive")
1918<sub>24</sub>
                      url += f"&num={num}"
1919 25
                  response = requests.get(url)
                  return response.text
1920<sup>26</sup>
1921 <sup>27</sup>
             except Exception as e:
                  return f"Error: {e}"
1922
```

E.2 PDF LOADER

1923 1924

1935 1936

```
1925
1926 1
       def pdf_load(file_path: str) -> ServiceResponse:
1927 <sup>2</sup>
            try:
                reader = PdfReader(file_path)
1928 <sup>3</sup>
                text = ""
1929
                for page in reader.pages:
1930 <sub>6</sub>
                     text += page.extract_text() + "\n"
                return ServiceResponse(status=ServiceExecStatus.SUCCESS,content=
1931 7
       text)
1932
            except Exception as e:
1933
                return ServiceResponse(ServiceExecStatus.ERROR, str(e))
1934
```

E.3 CSV READER

```
1937
1938 1
       import pandas as pd
1939 <sup>2</sup>
1940 <sup>3</sup>
       def csv_load(path:str)->ServiceResponse:
1941
                 df = pd.read_csv(path)
1942 <sub>6</sub>
                 csv_str = df.to_string(index=False)
1943 7
                 return ServiceResponse(status=ServiceExecStatus.SUCCESS,content=
       csv_str)
```

```
1944
            except Exception as e:
1945
                return ServiceResponse(ServiceExecStatus.ERROR, str(e))
1946
1947
1948
       E.4 XLSX READER
1949
1950
       def xlsx_load(path:str) ->ServiceResponse:
1951 <sub>2</sub>
            try:
1952 3
                excel_file = pd.read_excel(path, sheet_name=None)
                result = ""
1953 4
                for sheet_name, df in excel_file.items():
1954 <sup>5</sup>
                     result += f"Sheet: {sheet_name}\n"
1955
                     result += df.to_string(index=False) + "\n\n"
1956<sub>8</sub>
                return ServiceResponse(status=ServiceExecStatus.SUCCESS,content=
       result.strip())
1957
            except Exception as e:
1958 <sup>9</sup>
                return ServiceResponse(ServiceExecStatus.ERROR, str(e))
1959 <sup>10</sup>
1960
1961
       E.5 TEXT FILE READER
1962
1963
       import pandas as pd
1964
1965 <sub>3</sub>
       def txt_load(path:str) ->ServiceResponse:
1966<sub>4</sub>
            try:
                with open(path, 'r', encoding='utf-8') as f:
1967 5
                     txt_str = f.read()
1968 <sup>6</sup>
                return ServiceResponse(status=ServiceExecStatus.SUCCESS,content=
1969
       txt_str)
1970<sub>8</sub>
            except Exception as e:
1971 9
                return ServiceResponse(ServiceExecStatus.ERROR, str(e))
1972
1973
1974
       E.6 DOCX READER
1975
1976<sub>1</sub>
       from docx import Document
1977 2
       def docs_load(path:str) ->ServiceResponse:
1978 <sup>3</sup>
            try:
1979
                doc = Document(path)
1980
                docx_str = "\n".join([para.text for para in doc.paragraphs])
1981 7
                return ServiceResponse(status=ServiceExecStatus.SUCCESS,content=
1982
       docx_str)
            except Exception as e:
1983 8
                return ServiceResponse(ServiceExecStatus.ERROR, str(e))
1984
1985
       E.7 MP3 LOADER
1987
1988
       import whisper
1989
       from typing import cast
```

E.8 FIGURE LOADER

def load_audio(file):

return result["text"]

model = whisper.load_model(name="base")

model = cast(whisper.Whisper, model)

result = model.transcribe(str(file))

1990 ₃ 1991 ₄

1992 5

1993 ⁶

1994

199519961997

```
1998
1999 1
       from transformers import DonutProcessor, VisionEncoderDecoderModel
       import re
2000 <sup>2</sup>
       from PIL import Image
2001 <sup>3</sup>
2002 4
       def load_image(path):
2003 <sub>6</sub>
           image = Image.open(path)
2004 7
           processor = DonutProcessor.from pretrained(
                                   "naver-clova-ix/donut-base-finetuned-cord-v2"
2005 8
2006 9
2007 10
           model = VisionEncoderDecoderModel.from_pretrained(
                "naver-clova-ix/donut-base-finetuned-cord-v2"
2008 <sub>12</sub>
2009 13
           device = 'cpu'
2010 14
           model.to(device)
            # prepare decoder inputs
2011 15
           task_prompt = "<s_cord-v2>"
2012 16 17
            decoder_input_ids = processor.tokenizer(
2013 18
                task_prompt, add_special_tokens=False, return_tensors="pt"
2014 19
           ).input_ids
           pixel_values = processor(image, return_tensors="pt").pixel_values
2015 20
           outputs = model.generate(
2016 <sup>21</sup>
2017 22 23
                pixel_values.to(device),
                decoder_input_ids=decoder_input_ids.to(device),
2018 24
                max_length=model.decoder.config.max_position_embeddings,
2019 25
                early_stopping=True,
                pad_token_id=processor.tokenizer.pad_token_id,
2020 26
2021 <sup>27</sup>
                eos_token_id=processor.tokenizer.eos_token_id,
2022 28 29
                use_cache=True,
                num_beams=3,
2023 <sub>30</sub>
                bad_words_ids=[[processor.tokenizer.unk_token_id]],
2024 31
                return_dict_in_generate=True,
2025 32
           sequence = processor.batch_decode(outputs.sequences)[0]
2026 33
           sequence = sequence.replace(processor.tokenizer.eos_token, "").
2027 34
       replace(
2028 35
                processor.tokenizer.pad_token, ""
2029 36
            # remove first task start token
2030 37
           text_str = re.sub(r"<.*?>", "", sequence, count=1).strip()
2031 38
2032 39
           return text_str
```

E.9 VIDEO LOADER

2033 2034

2035

20492050

2051

```
2036 1
      import whisper
2037 <sub>2</sub>
       from typing import cast
2038 3
       from pydub import AudioSegment
       from pathlib import Path
2039 4
2040 <sup>5</sup>
       def load_video(file):
2041 7
           video = AudioSegment.from_file(Path(file), format=file[-3:])
2042 8
           audio = video.split_to_mono()[0]
           file_str = str(file)[:-4] + ".mp3"
2043 9
           audio.export(file_str, format="mp3")
2044 10
2045 11
           model = whisper.load_model(name="base")
           model = cast(whisper.Whisper, model)
2046 <sub>13</sub>
           result = model.transcribe(str(file))
2047 14
           return result["text"]
2048
```

E.10 DATA RETRIEVAL

```
def create_vector_db():
```

```
2052
            import faiss
2053
            import pickle
2054 4
           from sentence_transformers import SentenceTransformer
2055 5
           from data.mmlu import merge_csv_files_in_folder
           dataset=merge_csv_files_in_folder(path to MMLU/dev)
2056 6
           docs = []
2057 <sup>7</sup>
           for item in dataset:
2058 9
                    text = item[0].replace(",please answer A,B,C,or D.",",")+f"
2059
       answer:{item[1]}."
2060 10
                    docs.append(text)
2061 11
            embed_model = SentenceTransformer('all-MiniLM-L6-v2')
           doc_embeddings = embed_model.encode(docs)
2062 12
           dimension = doc_embeddings.shape[1]
2063 14
            index = faiss.IndexFlatL2(dimension)
2064 15
            index.add(doc_embeddings)
           faiss.write_index(index, "db/index.faiss")
2065 16
           with open("db/index.pkl", "wb") as f:
2066 17
                pickle.dump(docs, f)
2067 <sup>18</sup>
2068 20
       def load_vector_db():
2069 21
           import faiss
2070 22
            import pickle
2071 <sup>23</sup>
            from sentence_transformers import SentenceTransformer
           class db:
2072 <sup>24</sup>
    25
                def __init__(self):
2073 26
                     self.index = faiss.read_index("db/index.faiss")
2074 <sub>27</sub>
                     with open("db/index.pkl", "rb") as f:
2075 28
                         self.docs = pickle.load(f)
                     self.embed_model = SentenceTransformer('all-MiniLM-L6-v2')
2076<sup>29</sup>
                def search(self, query, k=5):
2077 30
                    query_embedding = self.embed_model.encode([query])
2078 32
                    D, I = self.index.search(query_embedding, k)
2079 33
                    return [self.docs[i] for i in I[0]]
2080 34
           return db()
2081
```

E.11 PROBLEM SOLVER

2082

```
2083
2084 1
        def twoSum(nums: List[int], target: int) -> List[int]:
2085 2
            Given an array of integers nums and an integer target, return indices
2086 <sup>3</sup>
         of the two numbers such that they add up to target.
2087
            Args:
2088
                 nums (List): an array of integers
2089 <sub>6</sub>
                 target (Int): an integer target
2090 7
            Returns:
                 List[int]: indices of the two numbers such that they add up to
2091 8
        target.
2092
            11 11 11
2093 10
            try:
2094 11
                 n = len(nums)
2095 12
                 for i in range(n):
                      for j in range (i + 1, n):
2096 13
2097 <sup>14</sup>
                           if nums[i] + nums[j] == target:
2098 15 16
                                return [i, j]
2099 17
                 return []
2100 18
            except Exception as e:
                 return str(e)
2101 19
2102 <sup>20</sup>
2103<sup>21</sup><sub>22</sub>
        def lengthOfLongestSubstring(s: str) -> int:
2104 <sub>23</sub>
2105 24
            Given a string s, find the length of the longest substring without
        duplicate characters.
```

```
2106 <sub>25</sub>
             Arg:
2107 26
                 s (String): a string
2108<sub>27</sub>
2109 28
             Returns:
                Int: the length of the longest substring without duplicate
2110^{29}
2111 30
        characters.
2112<sub>31</sub>
             try:
2113 32
                  left = 0
211433
                  right = 0
2115 34
                  max_len = 0
2116 <sup>35</sup>
                  while right < len(s):</pre>
2117<sub>37</sub>
                       if s[right] in s[left:right]:
2118<sub>38</sub>
                            max_len = max(max_len, right-left)
                            left = s.index(s[right], left, right)+1
2119 39
                       max_len = max(max_len, right-left+1)
2120<sup>40</sup>
2121 41 42
                       right += 1
                  return max_len
2122<sub>43</sub>
             except Exception as e:
                 return str(e)
2124 45
2125 <sup>46</sup>
2126 47
        def findMedianSortedArrays(nums1: List[int], nums2: List[int]) -> float:
2127 49
             Given two sorted arrays nums1 and nums2 of size m and n respectively,
2128
         return the median of the two sorted arrays.
2129 50
             Args:
2130 <sup>51</sup>
                 nums1 (List[int]): sorted array 1
2131 52 53
                 nums2 (List[int]): sorted array 2
             Returns:
2132 <sub>54</sub>
                 float: the median of the two sorted arrays
2133 55
2134 56
             try:
                  m, n = len(nums1), len(nums2)
2135 <sup>57</sup>
2136 58 59
                  def kth_small(k):
2137<sub>60</sub>
                       i = j = 0
2138 61
                       while True:
                            if i == m:
2139 62
                                return nums2[j + k - 1]
2140^{63}
2141 <sup>64</sup>
                            if j == n:
                                 return nums1[i + k - 1]
    65
2142 66
                            if k == 1:
2143 67
                                 return min(nums1[i], nums2[j])
                            pivot_i = min(i + (k >> 1) - 1, m - 1)
214468
                            pivot_j = min(j + (k >> 1) - 1, n - 1)
2145 <sup>69</sup>
                            if nums1[pivot_i] < nums2[pivot_j]:</pre>
2146 70 71
                                 k = pivot_i + 1 - i
2147 72
                                 i = pivot_i + 1
2148 73
                            else:
214974
                                 k = pivot_j + 1 - j
2150 75
                                 j = pivot_j + 1
2151 76 77
                  return (
2152<sub>78</sub>
                       kth\_small((m + n + 1 >> 1))
2153 79
                       if m + n & 1
                       else (kth_small((m + n >> 1) + 1) + kth_small((m + n >> 1)))
215480
2155 81
                       * 0.5
2156 82 83
                  )
             except Exception as e:
2157<sub>84</sub>
                 return str(e)
2158 85
2159 86
```

F USAGE OF LARGE LANGUAGE MODELS

In the preparation of this paper, we employed large language models to assist with language refinement and stylistic improvements. Typical prompts included instructions such as "please polish the following academic text while preserving its technical meaning", "improve clarity and conciseness without altering the content", or "translate the following text into fluent academic English."

The LLMs were not used for generating research ideas, designing experiments, conducting analyses, or interpreting results. All technical content, methodology, and conclusions are the sole work of the authors, who take full responsibility for the accuracy and validity of the presented material.