

000 001 002 003 004 005 AGENTRACE: BENCHMARKING EFFICIENCY IN LLM 006 AGENT FRAMEWORKS 007 008 009

010 **Anonymous authors**
 011 Paper under double-blind review
 012
 013
 014
 015
 016
 017
 018
 019
 020
 021
 022
 023
 024

ABSTRACT

025 Large Language Model (LLM) agents are rapidly gaining traction across domains
 026 such as intelligent assistants, programming aids, and autonomous decision systems.
 027 While existing benchmarks focus primarily on evaluating the effectiveness of
 028 LLM agents, such as task success rates and reasoning correctness, the efficiency
 029 of agent frameworks remains an underexplored but critical factor for real-world
 030 deployment. In this work, we introduce AgentRace, the first benchmark specifically
 031 designed to systematically evaluate the efficiency of LLM agent frameworks across
 032 representative workloads. AgentRace enables controlled, reproducible comparisons
 033 of runtime performance, scalability, communication overhead, and tool invocation
 034 latency across popular frameworks on diverse task scenarios and workflows. Our in-
 035 depth experiments reveal 9 insights and 12 underlying mechanisms for developing
 036 efficient LLM agents. We believe AgentRace will become a valuable resource
 037 for guiding the design and optimization of next-generation efficient LLM agent
 038 systems. The platform and results are available at the anonymous website <https://agent-race.github.io/>.
 039
 040

1 INTRODUCTION

041 Large Language Models (LLMs) (OpenAI, 2023; Touvron et al., 2023; Liu et al., 2024a; Naveed et al.,
 042 2023; Bai et al., 2023) have rapidly gained widespread popularity due to their exceptional capabilities
 043 in natural language understanding and generation, significantly impacting various applications
 044 including chatbots, content creation, and programming assistants. With these advancements, LLM
 045 agents (Wang et al., 2024; Guo et al., 2024; Zhao et al., 2024; Zhang et al., 2024; Ni & Buehler,
 046 2024), which are autonomous entities powered by LLMs capable of executing complex tasks through
 047 intelligent interactions, have emerged as a promising area of research and practical implementation.

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

053 However, for LLM agents to be widely deployed in real-world scenarios in the future, the efficiency of
 054 their frameworks is critically important. Efficient execution, scalability, and minimal communication
 055 overhead are essential for ensuring timely responses and practical usability, particularly in resource-
 056 constrained and latency-sensitive environments. Despite the proliferation of LLM agent frameworks,
 057 such as LangChain (LangChain, 2025), AutoGen (Wu et al., 2023), and AgentScope (Gao et al.,
 058 2024), a systematic benchmark evaluating these frameworks' performance efficiency remains absent.

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

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

054 3. How to improve the efficiency of agent execution?

055

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

062 In summary, our contributions include:

063

- 064 • 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.
 - 065 • We conduct a comprehensive and in-depth assessment of efficiency across frameworks,
 revealing previously undocumented sources of inefficiency.
 - 066 • We provide actionable insights for both practitioners and researchers to optimize the deploy-
 ment of efficient LLM-based agents.
 - 067 • We release the entire benchmark suite and experimental results, providing a platform to
 identify the efficiency issues of LLM agents.
- 068

069 2 BACKGROUND AND RELATED WORK

070

071 2.1 LLM AGENTS

072

073 LLMs agents (Yao et al., 2023; Zhao et al., 2024) are systems that combine the generative capabilities
 074 of LLMs with additional components such as memory, planning, and tool usage to perform complex
 075 tasks autonomously. These agents can interpret user inputs, plan actions, interact with external tools,
 076 and adapt based on feedback, enabling more dynamic and context-aware behaviors. Many agents
 077 have been developed, where some are generic agents that are designed to execute general tasks and
 078 some are specialized agents for some concrete task. For example, ReAct (Yao et al., 2023) is a typical
 079 general agent workflow, where the agent thinks and take actions interatively. MetaGPT (Hong et al.,
 080 2023) is an agent designed for software development, where each agent plays a different role to
 081 simulate a software company. In this work, we aim to evaluate the efficiency of different LLM agent
 082 frameworks, thus focusing on using the widely used general agent workflows.

083

084 2.2 LLM AGENT FRAMEWORKS

085

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

094

095 2.3 BENCHMARKS FOR LLM AGENTS

096

097 There have been many benchmarks for LLM agents (Andriushchenko et al., 2024; Chang et al.,
 098 2024; Huang et al., 2023; Shen et al., 2024; Liu et al., 2024b). However, most of these benchmarks
 099 usually focus on ability or trustworthiness perspectives, and do not exploit the efficiency part. For
 100 example, AgentBench (Liu et al., 2024b) report *Step Success Rate* as the main metric showing the
 101 independent accuracy of each action step, due to the current struggles for LLMs to ensure overall
 102 task success rates. Beyond benchmarks focusing solely on success rates, AgentBoard (Chang et al.,
 103 2024) proposes a comprehensive evaluation framework for LLM agents. It introduces a fine-grained
 104 *Progress Rate* metric to track incremental advancements during task execution, along with an open-
 105 source toolkit for multi-faceted analysis. WORFBENCH (Huang et al., 2023) introduces a unified
 106

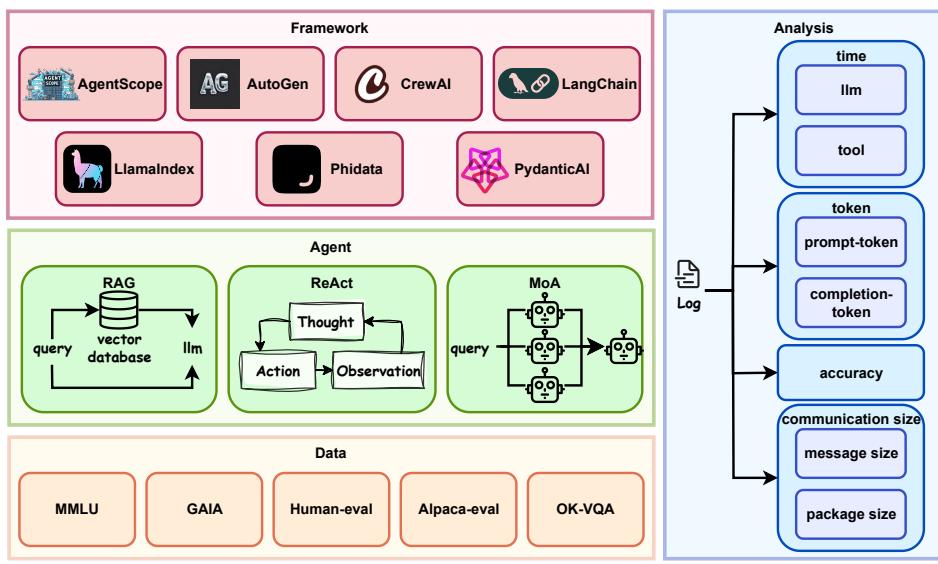


Figure 1: The architecture of AgentRace.

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-**

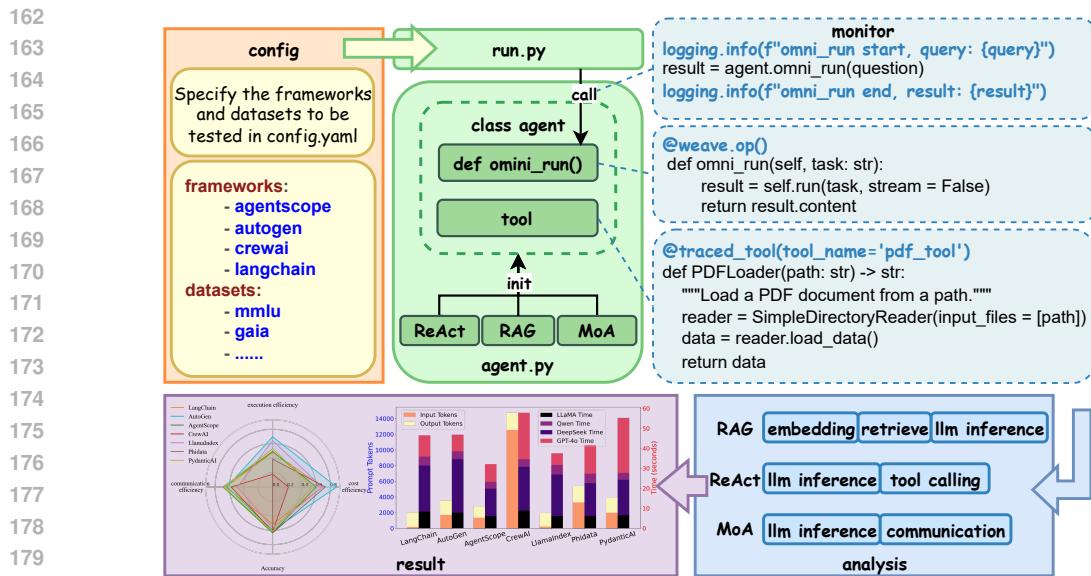


Figure 2: The pipeline of AgentRace.

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** (agnog-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

216 Table 1: The supported functionalities of AgentRace. ✓ denotes that the functionality is implemented
 217 in AgentRace. ○ denotes that the functionality is supported in the original framework.
 218

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

236
 237 the corresponding agent with unified interfaces. During execution, the agent interacts with the chosen
 238 framework and tools under controlled settings, while a monitoring layer is dynamically attached to
 239 capture runtime behavior. Finally, the analysis stage aggregates the collected traces into structured
 240 logs and performance visualizations for reproducibility and cross-framework comparison.
 241

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

252 3.3 FUNCTIONALITIES

253 The core functionalities supported by AgentRace are summarized in Table 1. Our benchmark
 254 currently supports three representative agent workflows executed across seven widely used LLM
 255 agent frameworks, utilizing a unified pool of eleven tools. While some of these capabilities are
 256 natively supported by the frameworks, approximately 50% of the functionalities are implemented by
 257 ourselves to ensure full compatibility and coverage. To maintain a fair comparison across frameworks,
 258 we adopt a standardized implementation for any functionality that is not natively provided. This
 259 ensures that differences in evaluation metrics stem from the underlying framework behavior, rather
 260 than implementation gaps. For more implementation details, please refer to Appendix A.
 261

262 4 EXPERIMENTS AND INSIGHTS

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

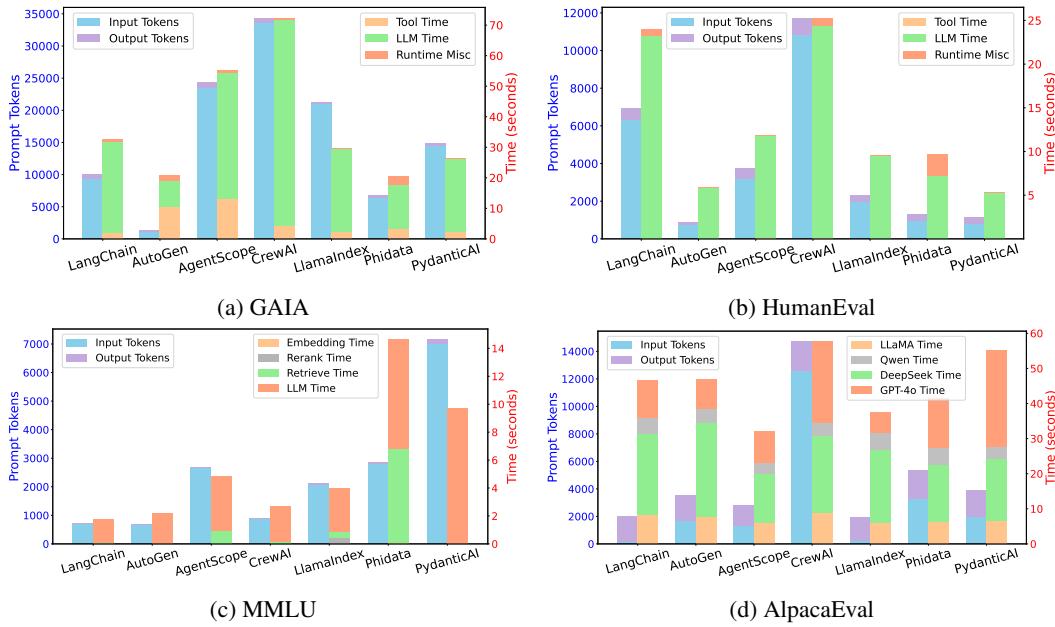


Figure 3: Token consumption and execution time per query of different frameworks.

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-4o 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-4o 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 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.



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

378 **Potential Optimizations** While LLM inference remains the dominant bottleneck in most of our
 379 benchmarks, more complex, tool-heavy scenarios, such as document analysis or multimodal agent
 380 tasks, may shift the performance bottleneck toward tool execution. Frameworks aiming to support
 381 such use cases must pay greater attention to optimizing tool orchestration and external API integration.
 382

383 4.4 RAG
 384

385 *Insight 3: While agents usually involve external databases for information retrieval, the database
 386 performance is overlooked in several frameworks. Vector database is recommended.*
 387

388 **Key Observations** While RAG workflows are increasingly adopted to enhance factual grounding,
 389 our benchmarking reveals that database performance, particularly during embedding and retrieval, is
 390 a critical yet frequently neglected factor. Figure 3c illustrates the variation in retrieval latency across
 391 frameworks, exposing significant performance disparities.
 392

393 **Underlying Mechanism-4: Embedding-Pipeline Design** One notable example is AgentScope,
 394 which demonstrates high vectorization latency. This stems from its design: during the database setup
 395 phase, AgentScope invokes a large embedding model to compute dense vector representations. The
 396 latency of this embedding model, often implemented as a separate LLM call, substantially increases
 397 the overall vectorization time. Similarly, Phidata exhibits elevated vectorization latency due to its use
 398 of a two-step pipeline. First, its built-in `csv_tool` loads documents row-by-row; then, it applies a
 399 `SentenceTransformer` model to compute embeddings. Our benchmark confirms that Phidata's
 400 `csv_tool` itself is a relatively slow component, compounding the overall vectorization time. From
 401 our observation, vector databases such as Faiss (Douze et al., 2024) are good choices.
 402

403 **Potential Optimizations** These observations highlight the need for more attention to retrieval
 404 pipeline design, especially in frameworks that aim to support real-time or large-scale RAG deploy-
 405 ments. Optimization opportunities include batching document embeddings, using faster embedding
 406 models, minimizing redundant file reads, and caching frequent queries.
 407

408 4.5 COMMUNICATION SIZE
 409

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

413 **Key Observations** In multi-agent frameworks, communication between agents is often overlooked
 414 as a source of inefficiency. However, our analysis reveals large discrepancies in communication size
 415 across frameworks, as shown in Table 2. These differences arise not only from framework-specific
 416 message formats but also from architectural design choices.
 417

418 **Underlying Mechanism-5: Inefficient Communication Architecture** Frameworks such as Cre-
 419 wAI, which adopt a centralized communication pattern, exhibit significantly higher communication
 420 costs. In these designs, a central agent coordinates multiple sub-agents by sequentially delegating
 421 subtasks and collecting responses. For example, in CrewAI's MoA implementation, the center agent
 422 queries three sub-agents in sequence and aggregates their outputs. Each LLM invocation by the
 423 center agent accumulates prior messages in memory, causing the prompt size and the communication
 424 payload to grow linearly with the number of sub-agents.
 425

426 **Underlying Mechanism-6: Package Design** In addition to the core message, Phidata returns a
 427 duplicated content field that mirrors the final message. This, combined with additional metadata
 428 fields, results in large communication sizes.
 429

430 **Potential Optimizations** These findings indicate that communication cost is not merely a func-
 431 tion of task complexity but also of framework design. Future agent frameworks should consider
 decentralized communication protocols and agent sampling to reduce unnecessary transfer overhead.
 432

432
433 Table 2: Communication size between agents (Unit: Byte). We report the content size (e.g., the
434 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	165.07/0	209.08/44.01	284.078/0	514.962/0	1180.078/898	354.508/0	96.022/0
	Agent2	165.07/0	209.08/44.01	284.078/0	483.740/0	1171.078/889	341.160/0	95.425/0
	Agent3	165.07/0	209.08/44.01	284.078/0	619.516/0	1164.078/882	343.219/0	97.116/0
To Aggregation Agent	Agent1	1983.02/3	2066.04/52.4	1659.318/0	2497.929/0	2022.417/33.689	6128.259/2639.113	2000.542/0
	Agent2	2011.83/3	2071.24/57.38	1511.311/0	1754.701/0	2054.878/39.118	6131.272/2629.426	1927.093/0
	Agent3	2072.98/3	2156.04/66.81	1889.247/0	2151.097/0	2116.377/48.641	5715.126/2465.817	1892.344/0

440
441 Table 3: Scalability Evaluation of AlpacaEval.
442

	Worker Agents	LangChain	AutoGen	AgentScope	CrewAI	LlamaIndex	Phidata	PydanticAI
Time (Unit: Second)	3	36.50	36.85	32.12	64.00	27.32	50.22	46.45
	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
	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
Total Token	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
	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
Communication Size (Unit: Byte)	3	6563.04	6920.56	5912.11	8021.94	9708.91	19013.54	6108.54
	6	14029.26	14383.36	10506.82	17863.90	19965.41	21684.95	12206.18
	9	20468.68	22325.97	16275.87	24769.81	29936.97	21320.89	16278.34
	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

455
456

4.6 SCALABILITY

457 *Insight 5: MoA scalability is governed by agent-invocation policy.*458
459
460 **Key Observations** We evaluate the scalability of the MoA workflow by increasing the number of
461 worker agents from 3 to 6, 9, 12, and 15, while keeping the additional agents identical in configuration
462 to the original ones. Table 3 reports the results on AlpacaEval. For frameworks such as AgentScope
463 and LangChain, both execution time and token consumption grow almost linearly with the number
464 of worker agents, reflecting sequential scheduling policies. In contrast, frameworks like PydanticAI
465 exhibit a significantly slower growth rate, suggesting a fundamentally different invocation strategy.
466467 **Underlying Mechanism-7: Parallel Execution** In PydanticAI, the observed runtime is shorter than
468 the aggregate of individual tool and LLM invocation times. This efficiency stems from its parallel
469 execution architecture: agent calls and tool invocations are dispatched asynchronously, allowing
470 multiple operations to overlap in time. As a result, the end-to-end latency is effectively bounded by
471 the slowest operation rather than the sum of all operations.
472473 **Potential Optimizations** Our analysis indicates that task-level parallelism remains largely underex-
474 plored in current frameworks. Incorporating asynchronous scheduling and concurrent invocation can
475 substantially improve scalability in multi-agent workflows, especially under real-world conditions
476 where latency and throughput are critical.
477478

5 CONCLUSION

479
480 We introduce AgentRace, a comprehensive benchmark platform for evaluating the efficiency of LLM
481 agent frameworks. AgentRace covers a diverse set of datasets, agent workflows, and frameworks,
482 enabling a fair and reproducible comparison across real-world scenarios. Through extensive and
483 in-depth experiments, we reveal key insights and underlying mechanisms. These findings highlight
484 critical optimization opportunities in the design and deployment of LLM-based agents. We hope
485 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.
486

486 **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
 487 reproducibility verification in Appendix B.6.
 488

489 REFERENCES
 490

- 492 Together.ai. <https://www.together.ai/>, 2024. Accessed: 2024-07-16.
- 493 Masarena, 2025. URL <https://lins-lab.github.io/MASArena/>. Accessed: 2025-09-
 494 23.
- 495 agno-ag. Phidata, 2025. URL <https://docs.phidata.com/introduction>. Accessed:
 496 2025-05-15.
- 497 Maksym Andriushchenko, Alexandra Souly, Mateusz Dziemian, Derek Duenas, Maxwell Lin, Justin
 498 Wang, Dan Hendrycks, Andy Zou, Zico Kolter, Matt Fredrikson, et al. Agentarm: A benchmark
 499 for measuring harmfulness of llm agents. *arXiv preprint arXiv:2410.09024*, 2024.
- 500 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,
 501 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- 502 Ma Chang, Junlei Zhang, Zhihao Zhu, Cheng Yang, Yujiu Yang, Yaohui Jin, Zhenzhong Lan,
 503 Lingpeng Kong, and Junxian He. Agentboard: An analytical evaluation board of multi-turn llm
 504 agents. *Advances in Neural Information Processing Systems*, 37:74325–74362, 2024.
- 505 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared
 506 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
 507 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 508 Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvassy, Pierre-Emmanuel
 509 Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. *arXiv preprint
 510 arXiv:2401.08281*, 2024.
- 511 Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled
 512 alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.
- 513 Dawei Gao, Zitao Li, Xuchen Pan, Weirui Kuang, Zhijian Ma, Bingchen Qian, Fei Wei, Wenhao
 514 Zhang, Yuexiang Xie, Daoyuan Chen, et al. Agentscope: A flexible yet robust multi-agent platform.
 515 *arXiv preprint arXiv:2402.14034*, 2024.
- 516 Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest,
 517 and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and
 518 challenges. *arXiv preprint arXiv:2402.01680*, 2024.
- 519 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
 520 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint
 521 arXiv:2009.03300*, 2020.
- 522 Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang,
 523 Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multi-agent
 524 collaborative framework. *arXiv preprint arXiv:2308.00352*, 3(4):6, 2023.
- 525 Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. Benchmarking large language models as ai
 526 research agents. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023.
- 527 LangChain. Langchain, 2025. URL <https://www.langchain.com/>. Accessed: 2025-05-15.
- 528 Zeping Lee. GB/T 7714-2015 BibTeX Style. <https://github.com/zepinglee/gbt7714-bibtex-style>, 2025. GitHub repository.
- 529 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
 530 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera-
 531 tion for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:
 532 9459–9474, 2020.

- 540 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
 541 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*
 542 *arXiv:2412.19437*, 2024a.
- 543 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,
 544 Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. In *ICLR*, 2024b.
- 545 LlamaIndex. Llamaindex, 2025. URL <https://www.llamaindex.ai/>. Accessed: 2025-05-15.
- 546 Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual
 547 question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf*
 548 *conference on computer vision and pattern recognition*, pp. 3195–3204, 2019.
- 549 Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia:
 550 a benchmark for general ai assistants. In *The Twelfth International Conference on Learning*
 551 *Representations*, 2023.
- 552 Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman,
 553 Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language
 554 models. *arXiv preprint arXiv:2307.06435*, 2023.
- 555 Bo Ni and Markus J Buehler. Mechagents: Large language model multi-agent collaborations can
 556 solve mechanics problems, generate new data, and integrate knowledge. *Extreme Mechanics*
 557 *Letters*, 67:102131, 2024.
- 558 OpenAI. Gpt-4 technical report, 2023.
- 559 PydanticAI. Pydanticai: A python agent framework for generative ai, 2025. URL <https://ai.pydantic.dev/>. Accessed: 2025-05-15.
- 560 Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang, Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng
 561 Li, and Yueting Zhuang. Taskbench: Benchmarking large language models for task automation.
 562 *Advances in Neural Information Processing Systems*, 37:4540–4574, 2024.
- 563 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
 564 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cris-
 565 tian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu,
 566 Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,
 567 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel
 568 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,
 569 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,
 570 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,
 571 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh
 572 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen
 573 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,
 574 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models,
 575 2023.
- 576 Junlin Wang, Jue WANG, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances
 577 large language model capabilities. In *The Thirteenth International Conference on Learning*
 578 *Representations*, 2025. URL <https://openreview.net/forum?id=h0ZfDIRj7T>.
- 579 Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai
 580 Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents.
 581 *Frontiers of Computer Science*, 18(6):186345, 2024.
- 582 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang,
 583 Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent
 584 conversation. *arXiv preprint arXiv:2308.08155*, 2023.
- 585 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
 586 React: Synergizing reasoning and acting in language models. In *International Conference on*
 587 *Learning Representations (ICLR)*, 2023.

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

597 Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. Expel: Llm
598 agents are experiential learners. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
599 volume 38, pp. 19632–19642, 2024.
600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648
649

A EXPERIMENTAL DETAILS

650
651

A.1 DETAILS ABOUT THE DATASETS

652
653

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 general-purpose 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 instruction-following 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.

654
655

A.2 DETAILS ABOUT THE WORKFLOWS

656
657

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.

658
659

A.3 DETAILS ABOUT THE FRAMEWORKS

660

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.

702 **A.4 VERSIONS OF EVALUATED FRAMEWORKS**
 703

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

706 707 Table 4: Versions of the LLM Agent frameworks employed in this paper.
 708

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

714 715 **A.5 HYPERPARAMETERS**
 716

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

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

723 724 **B ADDITIONAL RESULTS**
 725

726 **B.1 ACCURACY**
 727

728 Table 5: Accuracy of each framework on each dataset.
 729

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

730 731 732 733 Table 5 presents the accuracy of each framework. In general, the accuracy differences among
 734 frameworks are relatively small when using the same underlying LLM. However, there are still some
 735 notable exceptions.

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

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

744 **Underlying Mechanism-8: Structured Output Misalignment** Some frameworks, such as Lla-
 745 maIndex, require tool inputs to conform to a strict dictionary format. However, GPT-4o does not
 746 consistently produce structured outputs that align with these expectations, leading to frequent tool
 747 invocation failures. This issue can be partially mitigated if the framework explicitly enforces the
 748 format requirement during the registration phase or input schema definition. In contrast, other
 749 frameworks such as LangChain adopt stricter enforcement mechanisms. ReAct-style agents in these
 750 systems perform rigid output validation and initiate automatic retries when the model’s response
 751 deviates from the expected invocation structure. While such mechanisms increase robustness against
 752 malformed outputs, they may backfire in certain scenarios.

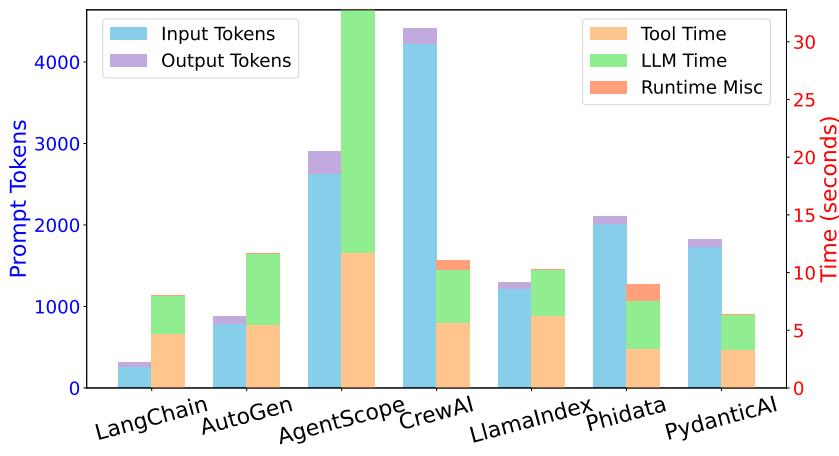


Figure 5: OK-VQA

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 experiment for 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	9358.35	1159.48	23520.479	33621.857	20935.364	6386.667	14459.17
	Output	637.92	180.66	785.891	664.511	304.976	323.558	320.588
	Total	9996.27	1340.15	24306.37	34286.369	21240.339	6710.224	14779.758
Time	llm	29.491	8.464	41.17	67.68	27.244	14.375	23.779
	Search	1.58856	9.4219	7.291	4.031	1.4399	1.83012	1.2275
	PDF loader	0.02423455	0.0009297	0.217	0.00965	0.0001352	0.001147	0.001395
	CSV reader	0.00003333	0.000336	0.000297	0.000196	0.00016616	0.0007207	0.0003148
	XLSX reader	0.06422606	0.002387	0.00405	0.00422	0.004254	0.003858	0.003795
	Text file reader	0.0004194	0.00002909	0.0000193	0.00123	0.000034839	0.0002107	8.6865E-06
	Doc reader	0.00009758	0.0002212	0.00000883	0.000278	0.0001135	0.000073355	0.000056241
	MP3 loader	-	-	0.729	0.000346	0.03341	0.03821	0.02965
	Figure loader	0.5345976	1.05489	4.083	0.03164	0.8767	1.4065	1.2104
	Video loader	-	-	0.0000271	0.000999	-	1.38445E-05	3.1952E-06
	Code executor	0.0152988	0.00005333	0.752	0.09565	0.05782	0.003035	0.0001414
Total tool time		2.22746732	10.4807	13.076	4.18	2.4126	3.2839	2.4732
Total time		32.492	20.76	55.092	72.195	29.795	20.396	26.238

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

Table 7: HumanEval Detailed Results

Framework	Token			Time		
	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
AgentScope	3180.689	561.518	3742.207	11.738	0.131	11.906
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 8: MMLU Detailed Results

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

Table 10: OK-VQA Detailed Results

Framework	token			time		
	prompt	output	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

B.3.1 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.

918
919
Table 11: Scalability Evaluation of AlpacaEval Using AgentScope

Number of Worker Agent			3	6	9	12	15
Token	llama	prompt output total	85.451 382.45 467.901	137.84 796.68 934.52	206.76 1204.91 1411.67	275.68 1641.18 1916.86	344.6 2021.9 2366.5
	qwen	prompt output total	61.815 311.109 372.924	89.92 555.47 645.39	134.88 848.94 983.82	179.84 1139.47 1319.31	224.8 1497.77 1722.57
	deepseek	prompt output total	52.478 416.639 469.117	71.74 841.37 913.11	107.61 1253.25 1360.86	143.48 1704.54 1848.02	179.35 2100.42 2279.77
	gpt	prompt output total	1138.243 352.564 1490.807	2237.83 412.43 2650.26	3351.55 439.44 3790.99	4542.02 442.62 4984.64	5677.57 434.15 6111.72
		llama	6.063	12.76	19.307	25.547	35.311
		qwen	3.415	6.523	10.819	13.866	18.237
Time	deepseek		13.726	32.81	48.833	67.114	84.318
	gpt		8.89	15.468	14.33	14.373	15.813
	total		32.119	67.607	93.357	122.987	153.784
Communication	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
	agent1 to aggregator		1659.318/124	3256.270/250	4960.120/375	6718.820/500	8266.330/625
	agent2 to aggregator		1511.311/122	2375.700/246	4051.120/369	5477.260/492	7080.860/615
	agent3 to aggregator		1889.247/126	3705.450/254	5510.530/381	7497.600/508	9255.700/635

940
941
Table 12: Scalability Evaluation of AlpacaEval Using AutoGen

Number of Worker Agent			3	6	9	12	15
Token	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
	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
		llama	7.812	14.667	25.424	34.833	37.816
		qwen	3.977	12.653	21.064	28.736	35.71
Time	deepseek		26.745	46.011	71.345	71.98	104.207
	gpt		8.274	19.816	22.41	30.398	17.817
	total		36.854	47.339	50.843	55.6	46.428
Communication	prompt to agent1		209.08/44.01	236.48/86.04	359.7/129.06	479.6/172.08	599.5/215.1
	prompt to agent2		209.08/44.01	236.48/86.04	359.7/129.06	479.6/172.08	599.5/215.1
	prompt to agent3		209.08/44.01	236.48/86.04	359.7/129.06	479.6/172.08	599.5/215.1
	agent1 to aggregator		2066.04/52.24	4618.13/103.41	7069.64/156.52	9297.61/208.1	11541.91/258.55
	agent2 to aggregator		2071.24/57.38	4450.9/112.37	6777.28/172.84	8661.68/217.75	10768.69/271.29
	agent3 to aggregator		2156.04/66.81	4604.89/128.62	7399.95/204.09	9384.64/258.11	11497.32/317.86

963
964
B.3.2 THE NUMBER OF TOOLS965
966
Insight 8: Increasing the number of tools has only a minimal impact on execution time across
967
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.*968
969
Key Observations We conduct scalability experiments on the GAIA dataset, examining the effect
970
of varying the number of tools across different frameworks. In addition to each frameworks original
971
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

972
973
974 Table 13: Scalability Evaluation of AlpacaEval Using LangChain
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995

Number of Worker Agent			3	6	9	12	15
Token	llama	prompt	70.49	105.84	158.76	211.68	264.60
		output	428.55	1054.54	1518.52	2037.28	2537.08
		total	499.04	1160.38	1677.28	2248.96	2801.68
	qwen	prompt	64.84	93.92	140.88	187.84	234.80
		output	446.05	1007.95	1446.68	2017.43	2436.53
		total	510.91	1101.87	1587.56	2205.27	2671.33
	deepseek	prompt	38.50	41.74	62.61	83.48	104.35
		output	501.11	1132.22	1677.97	2224.98	2792.75
		total	539.61	1173.96	1740.58	2308.46	2897.10
	gpt	prompt	1522.48	3300.82	4734.44	6353.31	7823.07
		output	444.81	693.66	661.37	685.78	700.94
		total	1967.29	3994.48	5395.81	7039.09	8524.01
Time	llama		8.275	12.061	19.123	23.213	34.437
			4.480	10.838	16.584	24.812	29.335
	qwen		23.084	40.801	66.156	73.476	115.888
			10.699	13.741	17.592	33.688	32.068
	deepseek		36.502	37.958	47.112	59.725	66.075
Communication	gpt						
			165.07/0	153.76/0	230.64/0	307.52/0	384.4/0
	prompt to agent1		165.07/0	153.76/0	230.64/0	307.52/0	384.4/0
			165.07/0	153.76/0	230.64/0	307.52/0	384.4/0
	prompt to agent2		165.07/0	153.76/0	230.64/0	307.52/0	384.4/0
			165.07/0	153.76/0	230.64/0	307.52/0	384.4/0
	prompt to agent3		1983.02/3	4703.67/6	6787.84/9	9117.26/13	11314.20/17
			2011.83/3	4334.61/6	6286.30/9	8621.19/13	10546.46/17
	agent1 to aggregator		2072.98/3	4529.70/6	6702.62/9	8880.47/13	11164.34/17

997
998 Table 14: Scalability Evaluation of AlpacaEval Using PydanticAI
999

Number of Worker Agent			3	6	9	12	15
Token	llama	prompt	61.347	95.5	126.29	139.77	161.71
		output	429.543	938.08	1273.35	1327.71	1559.8
		total	490.889	1033.58	1399.64	1467.48	1721.51
	qwen	prompt	41.217	58.39	76.32	80.85	94.52
		output	433.739	939.44	1213.31	1257.55	1608.87
		total	474.957	997.83	1289.63	1338.4	1703.39
	deepseek	prompt	31.802	41.44	50.5	50.88	58
		output	434.81	931.31	1210.15	1150.95	1311.62
		total	485.612	972.75	1260.65	1201.83	1369.62
	gpt	prompt	1845.724	3531.53	4673.28	4739.51	5684.52
		output	596.876	636.99	633.62	637.09	691.85
		total	2442.6	4168.52	5306.9	5376.6	6376.37
Time	llama		6.503	15.15	16.68	19.71	21.15
			3.441	8.38	11.2	11.59	13.41
	qwen		17.79	33.34	42.14	40.71	47.19
			27.486	22.05	90.94	91.35	41.02
	deepseek		46.45	42.24	110.78	111.4	62.13
Communication	gpt						
			96.022/0	88.12/0	113.86/0	124.09/0	134.34/0
	prompt to agent1		95.425/0	93.84/0	118.13/0	119.19/0	131.11/0
			97.116/0	94.73/0	108.99/0	103.12/0	113.92/0
	agent1 to aggregator		2000.542/0	4154.19/0	5693.77/0	6003.71/0	6851.79/0
			1927.093/0	4002.04/0	5302.46/0	5314.6/0	6682.15/0
	agent2 to aggregator		1892.344/0	3773.26/0	4941.13/0	4729.39/0	5284.75/0

1022
1023
1024 impact on execution time, it leads to a noticeable increase in LLM token usage. In addition, it
1025 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

Table 15: Scalability Evaluation of AlpacaEval Using CrewAI

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

Table 16: Scalability Evaluation of AlpacaEval Using LlamaIndex

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

Table 17: Scalability Evaluation of AlpacaEval Using Phidata

Number of Worker Agent			3	6	9	12	15
Token	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
	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	llama		6.152	6.55	6.55	9.12	10.33
	qwen		4.707	6.75	5.27	6.09	6.56
	deepseek		16.456	15.43	16.6	19.32	22.07
	gpt		14.208	23.13	25.68	31.67	31.7
Communication	prompt to agent1		354.508/0	325.7/0	310.63/0	329.16/0	334.87/0
	prompt to agent2		341.160/0	309.01/0	293.84/0	319.17/0	326.58/0
	prompt to agent3		343.219/0	304.15/0	288.79/0	307.71/0	314.94/0
	agent1 to aggregator		6128.259/2639.113	7105.44/3163.16	6961.87/3105.45	7291.8/3252.42	7582.27/3388.25
	agent2 to aggregator		6131.272/2629.426	7475.54/3354.1	7269.53/3267.58	7792.87/3505.45	8121.06/3656.93
	agent3 to aggregator		5715.126/2465.817	6165.11/2699.97	6196.23/2734.51	6342.37/2791.33	6571.72/2891.08

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.

Potential Optimizations Building on these findings, agent frameworks should emphasize relevance-aware 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 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 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

Table 20: Number of Failed Runs

	LangChain	AutoGen	AgentScope	CrewAI	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

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.

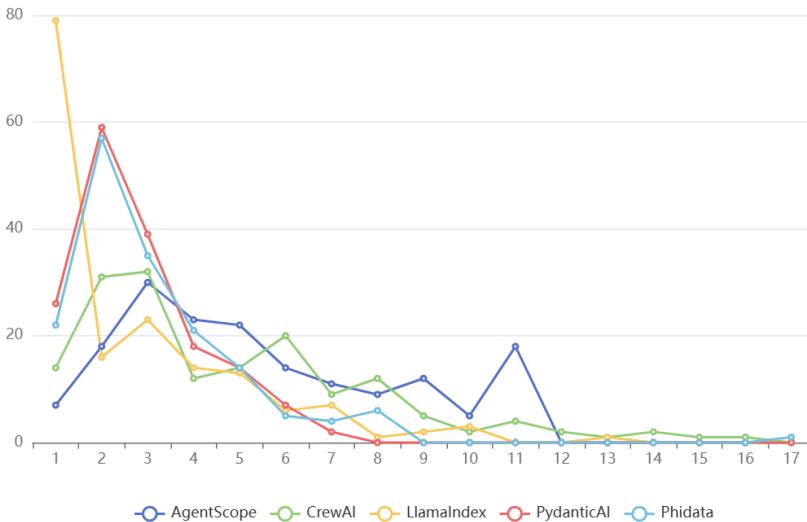


Figure 6: LLM Call Frequency per Query across Different Frameworks

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

- 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,

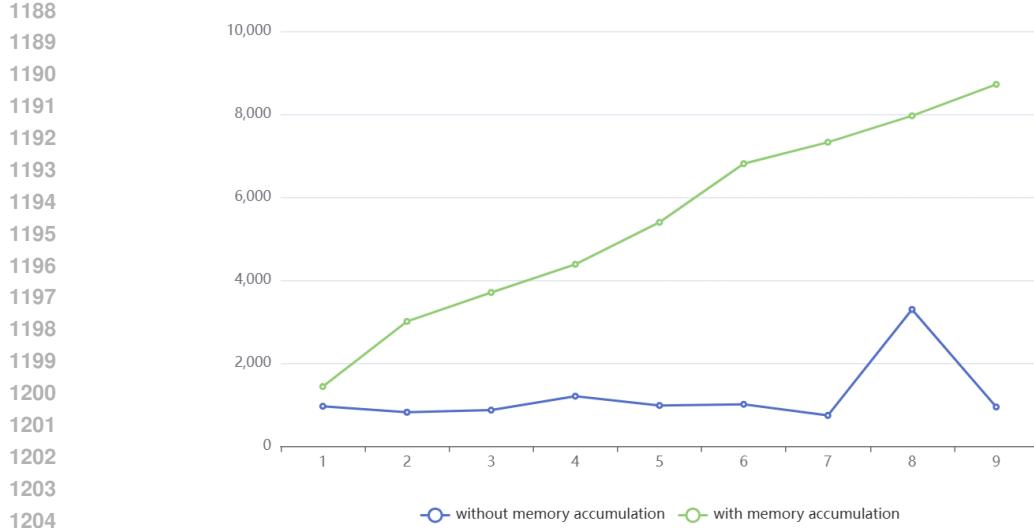


Figure 7: Memory Accumulation Impact

Qwen, DS, and GPT denote GLM-Z1-Rumination-32B-0414, Qwen2.5-7B-Instruct, DeepSeek-V3, and GPT-4o, respectively.

Table 21: The Impact of Agent Execution Order on Tokens

Order	GLMQwenDS			DSQwenGLM			QwenDSGLM		
	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

B.5 CLAUDE-BASED RESULTS

Table 22: Claude-Based HumanEval Results

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

Given that the majority of our experiments are implemented with GPT-4o, 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.

1242 **B.6 REPRODUCIBILITY VERIFICATION**
 1243

1244 Table 23: HumanEval Run 2
 1245

Framework	Token			Time		
	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

1254 Table 24: HumanEval Run 3
 1255

Framework	Token			Time		
	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

1265 Table 25: GAIA Run 1
 1266

Frameworks	Token			Time				
	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	

1285 *Insight 9: Experimental reproducibility is underpinned by the stability of token usage, while variability arises from stochastic tool behaviors and fluctuating LLM invocation dynamics.*

1289 **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 and 9, the token consumption in our experiment is relatively stable. In general, the execution time is usually positively related to the token consumption.

1294 **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

Table 26: GAIA Run 2

Frameworks	Token			Time				
	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

Table 27: GAIA Run 3

Frameworks	Token			Time				
	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

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.

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.

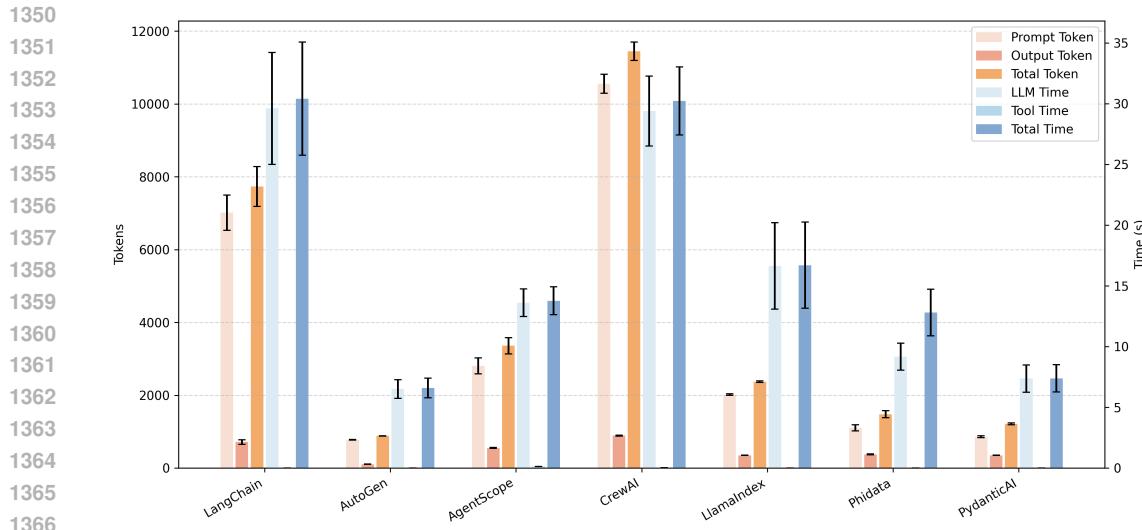


Figure 8: Consistency of Token Consumption and Latency in Repeated Experiments (HumanEval)

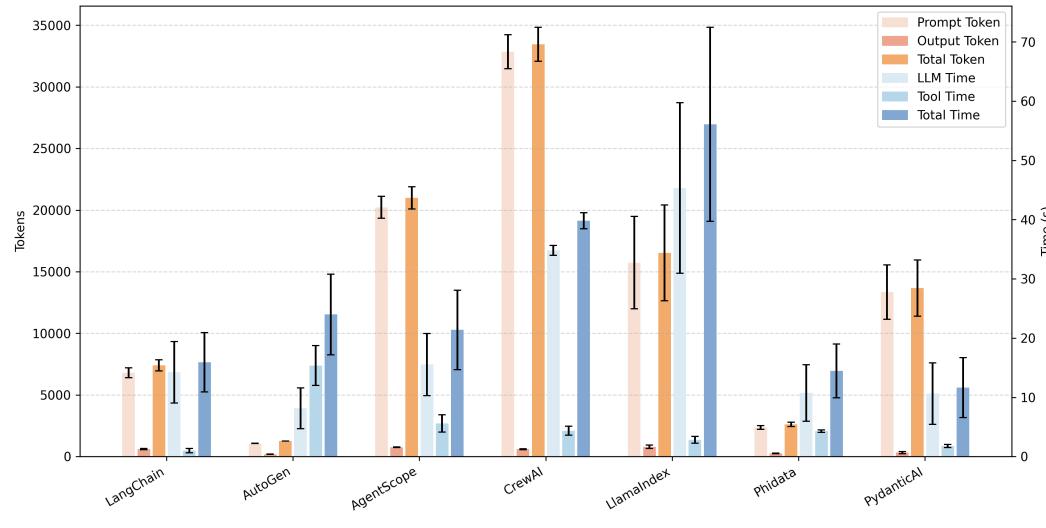


Figure 9: Consistency of Token Consumption and Latency in Repeated Experiments (GAIA)

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:

```

1 You are a ReAct-based assistant.
2 You analyze the question, decide whether to call a tool or directly
3 answer, and then respond accordingly.
4 Use the following format: Question: the input question or request
5 Thought: you should always think about what to do\nAction: the action to
6 take (if any)
7 Action Input: the input to the action (e.g., search query)
8 Observation: the result of the action
9 ... (this process can repeat multiple times)
10 Thought: I now know the final answer
11 Final Answer: the final answer to the original input question or request

```

1404
1405
1406 Table 28: Comparison of Token and Time Ratios Across Hardware Configurations
1407
1408
1409
1410
1411
1412
1413

Framework	Token ratio			Time ratio		
	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

1414

1415

10 Begin!
11 Question: {input}

1416

1417

1418

1419

C.1.1 LANGCHAIN

1420

1421

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.

1422

1423

1424

C.2 RAG

1425

1426

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

1427

1428

1429

C.2.1 AUTOGEN

1430

1431

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

1432

1433

1434

1435

C.2.2 CREWAI

1436

1437

1 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'.

1438

1439

1440

1441

1442

1443

C.2.3 PHIDATA

1444

1445

1446

1 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.

1447

1448

1449

C.2.4 PYDANTICAI

1450

1451

1452

1453

1 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.

1454

1455

1456

1457

C.3 MoA

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

1458
1459

C.3.1 LANGCHAIN

1460 You have been provided with a set of responses from various open-source
 1461 models to the latest user query. Your task is to synthesize these
 1462 responses into a single, high-quality response. It is crucial to
 1463 critically evaluate the information provided in these responses,
 1464 recognizing that some of it may be biased or incorrect. Your response
 1465 should not simply replicate the given answers but should offer a refined,
 1466 accurate, and comprehensive reply to the instruction. Ensure your
 1467 response is well-structured, coherent, and adheres to the highest
 1468 standards of accuracy and reliability.

1468
1469
1470

C.3.2 AGENTSCOPE

1471 You are an assistant called Dave, you should synthesize the answers from
 1472 Alice, Bob and Charles to arrive at the final response.
 1473

1474
1475
1476

For the worker agent, we used the following prompt.

You are an assistant called Alice/Bob/Charles.

1477

1478
1479

C.3.3 CREWAI

1480 You are an agent manager, and You need to assign the questions you
 1481 receive to each of your all agents, and summarize their answers to get a
 1482 more complete answer
 1483 You must give question to all the all agents, and you must summarize
 1484 their answers to get a more complete answer.\nYou need to be the best

1485
1486

For the worker agent, we used the following prompt.

1487
1488
1489
1490

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

1491

1492
1493

C.3.4 PHIDATA

1494 Transfer task to all chat agents (There are 3 agents in your team)", "
 1495 Aggregate responses from all chat agents

1496

1497
1498

C.3.5 PYDANTICAI

1499 Your task is to aggregate all agents results to solve complex tasks.\nYou
 1500 analyze the input, input the task to all tools that can run a single
 1501 agent, and synthesize the results from all agents into a final response.

1502

1503
1504

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:

1505
1506
1507
1508
1509
1510
1511

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?

```

1512
1513 1 question: The attached spreadsheet contains the inventory of a movie and
1514 video game rental store located in Seattle, Washington. What is the title
1515 of the oldest Blu-Ray listed in this spreadsheet? Return it exactly as
1516 it appears., file_name: 32102e3e-d12a-4209-9163-7b3a104efe5d.xlsx,
1517 file_path: path/to/32102e3e-d12a-4209-9163-7b3a104efe5d.xlsx
1518
1519
1520
```

C.5 HUMAN EVAL

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:

```

1525 1 from typing import List
1526 2
1527 3 def has_close_elements(numbers: List[float], threshold: float) -> bool:
1528 4     """ Check if in given list of numbers, are any two numbers closer to
1529 each other than
1530 given threshold.
1531 >>> has_close_elements([1.0, 2.0, 3.0], 0.5)
1532 False
1533 >>> has_close_elements([1.0, 2.8, 3.0, 4.0, 5.0, 2.0], 0.3)
1534 True
1535 """
1536
1537 # Complete the function. Only return code. No explanation, no comments,
1538 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:

```

1548
1549 1 Question:Find the degree for the given field extension Q(sqrt(2), sqrt(3),
1550 sqrt(18)) over Q.
1551 A.0
1552 B.4
1553 C.2
1554 D.6
1555 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.

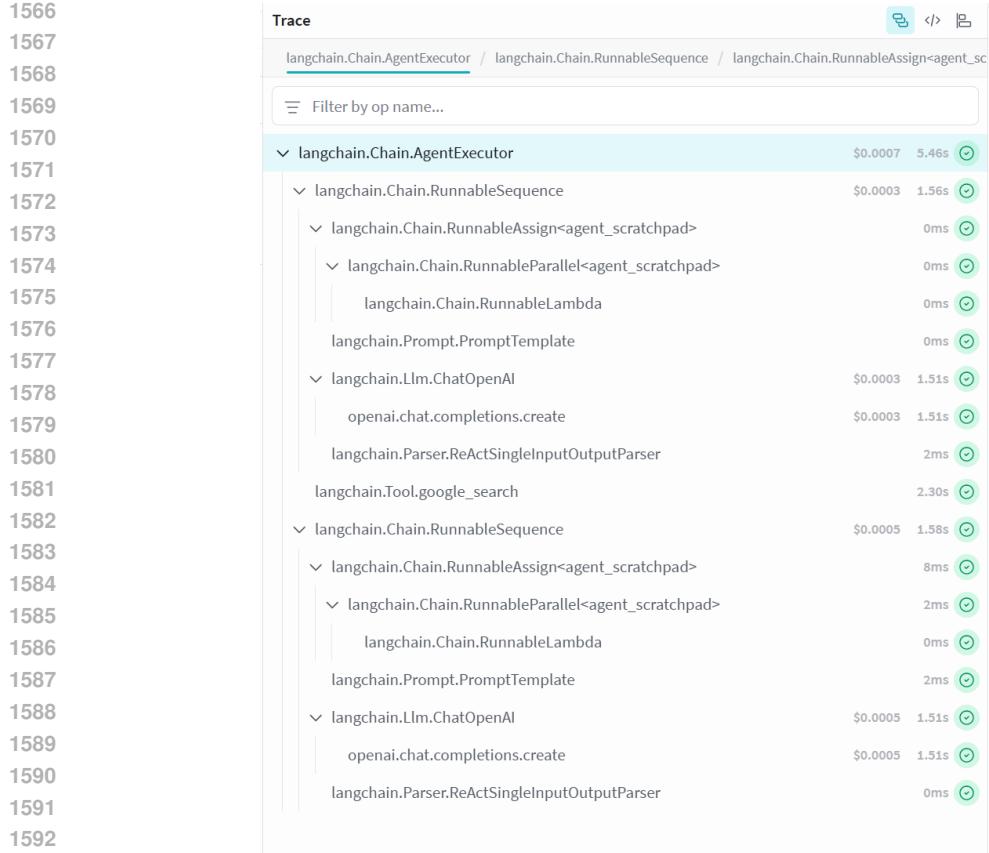


Figure 10: LangChain’s high level of abstraction and encapsulation.

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.

2025-05-12 2-3-4&81172-145 - root - INFO - tool name.pdf tool_time 3.843

2025-05-12 2-3-4&81172-172 - root - INFO - tool end result: [document]metadata=[{"producer": "Project MUSE", "creator": "PyPDF", "creationdate": "", "author": "David Greetham", "language": "English", "title": "Uncoupled OR, How I Lost My Author(s)", "xmp-dc:creator": "David Greetham", "xmp-dc:language": "English", "xmp-dc:title": "Uncoupled OR, How I Lost My Author(s)", "xmp-dc:type": "Text", "xmp-dc:format": "text/html", "xmp-dc:subject": "Uncoupled; OR, How I Lost My Author(s); David Greetham; Textual Cultures: Texts, Contexts, Interpretation, Volume 3, Number 1, Spring 2008, pp. 44-55 (Article)Published by Society for Textual Scholarship in the Vol(D)I: Textual Cultures", "xmp-dc:source": "https://doi.org/10.2970/ies.2008.3.1.4411https://muse.jhu.edu/article/240795/pdf", "total_pages": 13}, {"page_content": "Uncoupled; OR, How I Lost My Author(s); David Greetham; Textual Cultures: Texts, Contexts, Interpretation, Volume 3, Number 1, Spring 2008, pp. 44-55 (Article)Published by Society for Textual Scholarship in the Vol(D)I: Textual Cultures", "xmp-dc:content": "Additional information about this article\nhttps://doi.org/10.2970/ies.2008.3.1.4411https://muse.jhu.edu/article/240795/103224417384 Project MUSE [2025-05-12 16:40 GM]x00Uncoupled; OR, How I Lost My Author(s); David Greetham; Textual Cultures: Texts, Contexts, Interpretation, Volume 3, Number 1, Spring 2008, pp. 44-55 (Article)Published by Society for Textual Scholarship in the Vol(D)I: Textual Cultures", "xmp-dc:link": "https://doi.org/10.2970/ies.2008.3.1.4411https://muse.jhu.edu/article/240795/103224417384"}, {"page_content": "uncoupling, which deviates as the mapping out of a comprehensive range of abstract models for charting the relationships between a translator and a translating text (in the case of Iwao Trivisa) and the construction of an authorial identity, the essay concludes by showing how authoriality (rather than authorship) took over in the development of a scholarly career, resulting in, for example, the founding of the interdisciplinary Society for Textual Scholarship; In ForS Hill's TRIED TO FIND, AND BELOW IT ALTO, A N A U T H O R I R E A L Y D I D B. But somehow "the 'psychic connections'" that my fellow-paneelists identified, never quite took with me. As a very neophyte editor, I was co-pub (subedit) into the Well Clarendon Press edition of John Treviss' Middle English translation of *In Barholmioamus Angelus* = De Proptertatibus Reum. (On the Properties of Vnthings), or DPR, a well-known medieval text. I edited it in 1975. In 2-3-4&81172-172, Trivisa, a diligent scholar and teacher, has written a detailed account of the psychological processes of the screen editor and author (and also contributions) a year earlier. The introduction I dedicate to him in the first volume of *Singer-Scholar* (1935-2007) is also relevant here. Linda Quintana's contribution occurs in vol. 2, pp. 726-724, and I also worked on the notes for vol. 3, pp. 103-224417384. Project MUSE [2025-05-12 16:40 GM]x00David Greetham Uncoupled (45) presumably according to David C. Fowler, the author of *Piers Plowman*, re-named somehow determinedly faceless, basically a faithful translator (and a unsuccessful one) of "other men's flowers". Iwao did try, and, to be honest, quite enjoyed the ultimately failed at-attempt, in part because it gave me the opportunity to delve into some of the 'immore arcane areas of textual research (based on this biographical fact of *Vntrivis* a having devoted himself to translation), and thereby to produce 'whati' now regard as a series of too-neat, too pat, and perhaps also too pre-intentional "models" of translation which Fredson Bowers was generous 'enough to publish in *Studies in Bibliography* (see G e r i t h a m 1984). I don't 'think' whati' has any point now (and besides which I would probably find it too embarrassing a display of my youthful idealism) to try to actually *explain* whati' these fanciful models mean, ranging as they do from whati' called 'In Perfect Linear' (in which all witness in the original text are mapped onto the target text) to 'In Perfect Non-linear' (in which all witness in the original text are mapped onto the target text). I think they 'show' my youthful earnestness, my attempts to situate Trivisa and his text as 'van ideal form' to which I could declare allegiance, and thus embrace as a fa-i-ther figure. The models were quite wonderful in their enthusiasm, and we'll dedicated to the project of gaining a more direct access to my author, ob-s. In while I have obviously 'repressed' this recognition in the oral delivery of this paper, [the recent death of David Fowler has brought home to me more forcefully his role as 'fa-i-ther figure' looming in and other papers in this collection]. It was Fowler who read my first meager publication in a medievalist newsletter mysteri-'ously called 'Rapport' on the strengths and weaknesses of whati' I had presented me as 'my first attempt at a translation'. I am grateful to 'lecture on Trivisa in a textual study' as one of the first to advise me to change my name to 'Schulz' (as did Peter Sheldene, a contributor to E-T and to our conference) and a youngish kindly presence in my early textual life. While this may have been acknowledged that was the path I was to take, and continued to be, wise, cour-ageous, and benevolent mentors till the end of their lives). 4. The models were intended to chart every possible relative relationship between the trans-lated and translating text, from unambiguous to complex, including various types whati' that would cover not only variance in the witnesses of the tra-nslation text and 'unanimity' in the source, but also variance in both source and translation. (vñw46) Textual Cultures 3.1 (2008) inscribed not by a "veil of print" but by the scibl of transcription and mis-transmission. Our/our author lay beyond this confiusion n, namentam as the 'vñdestors' for my arcane models. I think I have to recognise that the project *Vntrivis* ultimately not as effective at 'unveiling' the author as I had hoped and *Vntrivis* thus also not effective at revealing the 'real' author. I think I have to accept that the project *Vntrivis* was a failure, and that the thoughts and the then-domain editorial principles with which we had unconsciously 'Vntrivised' were inappropriate to the reading of Trivisa, and hence to writing a series of essays not setting up an idealised author and an idealised text and its distanciation myself both from these essentially Formalist/New Critical/authorial/philosophical protocols that could be emp loved on *Vntrivis*. This (distantancing (and embarrassment) has continued down to the pres-ent, for in an article recently published in the Italian textual journal *Eccot-Diaria*, I return again to the scene of the crime and use the personal and cultural distanciation as a means of showing just what was, and then what might be 'known' that the author's search an immanent author (mis)represented by the *vñtext*, but inevitably corrupt, witness, has given way to a 'recognition' that the variance shown in these witness'es is valuable evidence for the so-'localization' of the text (G e r i t h a m 2006). In this case, my Ecotida piece, 'vnwhile formally a response to Paul Eggert's 'generous and comprehensive re-vñtext-essay' (2005) of

Figure 11: LangChain occasionally terminated processes prematurely.

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).

1620 D.2 AUTOGEN

1621

1622 Due to the default system prompt being relatively long and containing irrelevant instructions, the
 1623 RAG workflow may consume unnecessary tokens or produce unexpected errors (e.g., attempting to
 1624 invoke non-existent tools). Therefore, it is necessary for users to customize the system prompt.

1625

1626 D.3 AGENTSCOPE

1627

1628 AgentScopes image and audio processing tools internally rely on OpenAI models, causing their
 1629 execution time to partially overlap with that of the LLM itself. This overlap can lead to inflated or
 1630 inaccurate measurements of LLM processing time. Researchers and practitioners should be mindful
 1631 of this issue when conducting time-based evaluations involving AgentScope.

```

1632 1 def openai_image_to_text(
1633 2     image_urls: Union[str, list[str]],
1634 3     api_key: str,
1635 4     prompt: str = "Describe the image",
1636 5     model: Literal["gpt-4o", "gpt-4-turbo"] = "gpt-4o",
1637 6 ) -> ServiceResponse:
1638 7     """
1639 8         Generate descriptive text for given image(s) using a specified model,
1640 9         and
1641 10        return the generated text.
1642 11    Args:
1643 12        image_urls (`Union[str, list[str]]`):
1644 13            The URL or list of URLs pointing to the images that need to
1645 14            be
1646 15            described.
1647 16        api_key (`str`):
1648 17            The API key for the OpenAI API.
1649 18        prompt (`str`, defaults to `"Describe the image"`):
1650 19            The prompt that instructs the model on how to describe
1651 20            the image(s).
1652 21        model (`Literal["gpt-4o", "gpt-4-turbo"]`, defaults to `"gpt-4o"`)
1653 22    :
1654 23        The model to use for generating the text descriptions.
1655 24    Returns:
1656 25        `ServiceResponse`:
1657 26            A dictionary with two variables: `status` and `content`.
1658 27            If `status` is `ServiceExecStatus.SUCCESS`,
1659 28            the `content` contains the generated text description(s).
1660 29    Example:
1661 30        .. code-block:: python
1662 31
1663 32            image_url = "https://example.com/image.jpg"
1664 33            api_key = "YOUR_API_KEY"
1665 34            print(openai_image_to_text(image_url, api_key))
1666 35
1667 36            > {
1668 37            >     'status': 'SUCCESS',
1669 38            >     'content': "A detailed description of the image..."
1670 39            > }
1671 40
1672 41            openai_chat_wrapper = OpenAIChatWrapper(
1673 42                config_name="image_to_text_service_call",
1674 43                model_name=model,
1675 44                api_key=api_key,
1676 45            )
1677 46            messages = Msg(
1678 47                name="service_call",
1679 48

```

```

1674      role="user",
1675      content=prompt,
1676      url=image_urls,
1677  )
1678  openai_messages = openai_chat_wrapper.format(messages)
1679  try:
1680      response = openai_chat_wrapper(openai_messages)
1681  return ServiceResponse(ServiceExecStatus.SUCCESS, response.text)
1682 except Exception as e:
1683     return ServiceResponse(ServiceExecStatus.ERROR, str(e))

1684 def openai_audio_to_text(
1685     audio_file_url: str,
1686     api_key: str,
1687     language: str = "en",
1688     temperature: float = 0.2,
1689 ) -> ServiceResponse:
1690     """
1691         Convert an audio file to text using OpenAI's transcription service.
1692
1693     Args:
1694         audio_file_url (`str`):
1695             The file path or URL to the audio file that needs to be
1696             transcribed.
1697         api_key (`str`):
1698             The API key for the OpenAI API.
1699         language (`str`, defaults to `"en"`):
1700             The language of the input audio. Supplying the input language
1701             in
1702                 [ISO-639-1] (https://en.wikipedia.org/wiki/List\_of\_ISO\_639-1\_codes)
1703             format will improve accuracy and latency.
1704         temperature (`float`, defaults to `0.2`):
1705             The temperature for the transcription, which affects the
1706             randomness of the output.
1707
1708     Returns:
1709         `ServiceResponse`:
1710             A dictionary with two variables: `status` and `content`.
1711             If `status` is `ServiceExecStatus.SUCCESS`,
1712             the `content` contains a dictionary with key 'transcription'
1713             and
1714                 value as the transcribed text.
1715
1716 Example:
1717
1718     .. code-block:: python
1719
1720         audio_file_url = "/path/to/audio.mp3"
1721         api_key = "YOUR_API_KEY"
1722         print(openai_audio_to_text(audio_file_url, api_key))
1723
1724         > {
1725             >     'status': 'SUCCESS',
1726             >     'content': {'transcription': 'This is the transcribed text
1727 from
1728             the audio file.'}
1729         > }
1730
1731     """
1732     try:
1733         import openai
1734     except ImportError as e:
1735         raise ImportError(
1736             "The `openai` library is not installed. Please install it by
1737 "

```

```

1728      "running `pip install openai`.",
1729      ) from e
1730
1731      client = openai.OpenAI(api_key=api_key)
1732      audio_file_url = os.path.abspath(audio_file_url)
1733      with open(audio_file_url, "rb") as audio_file:
1734          try:
1735              transcription = client.audio.transcriptions.create(
1736                  model="whisper-1",
1737                  file=audio_file,
1738                  language=language,
1739                  temperature=temperature,
1740          )
1741          return ServiceResponse(
1742              ServiceExecStatus.SUCCESS,
1743              {"transcription": transcription.text},
1744          )
1745      except Exception as e:
1746          return ServiceResponse(
1747              ServiceExecStatus.ERROR,
1748              f"Error: Failed to transcribe audio {str(e)}",
1749          )

```

Meanwhile, AgentScope’s vector database module, LlamaIndexKnowledge, is implemented based on the BM25Retriever from the llamaindex library. However, the original implementation relies on an outdated version of llamaindex, and recent updates to the library introduced structural changes that break compatibility with the original import statements.

```

1752 from llama_index.retrievers.bm25 import BM25Retriever
1753

```

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.

```

1758 from llama_index.legacy.retrievers.bm25_retriever import BM25Retriever
1759

```

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

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

1782

1783

1784

1785

1786

1787

Response

1788

1789

1790

1791

1792

1793

1794

1795

1796

1797

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

1798

1799

1800

1801

1802

1803

1804

1805

1806

1807

1808

1809

1810

1811

1812

1813

1814

1815

1816

1817

1818

1819

1820

1821



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.

Inputs

Path	Value
self	<openai.resources.chat.completions.completions.Completions object at 0x304c89ab0>
> messages	
model	gpt-4o
temperature	0
tool_choice	auto
> tools	
> 0	
> 1	
> 2	
> 3	
> 4	
> 5	
> 6	
> 7	
> 8	
> 9	
type	function
< function	
name	run_python_code
description	This function runs Python code in the current environment. If successful, returns the value ...
< parameters	
type	object

Instructions

- You are a ReAct-based assistant.
 You analyze the question, decide whether to call a tool or directly answer, and then respond accordingly.
 Use the following format: Question: the input question or request
 Thought: you should always think about what to do
 Action: the action to take (if any)
 Action Input: the input to the action (e.g., search query)
 Observation: the result of the action
 ... (this process can repeat multiple times)
 Thought: I now know the final answer
 Final Answer: the final answer to the original input question or request
 Begin!
 Question: {input}
 - Use markdown to format your answers.

from typing import List

```
def filter_by_substring(strings: List[str], substring: str) -> List[str]:
    """ Filter an input list of strings only for ones that contain given substring
    >>> filter_by_substring([], 'a')
    []
    >>> filter_by_substring(['abc', 'bcd', 'cde', 'array'], 'a')
    ['abc', 'bcd', 'array']
    """

```

Response

The provided function `filter_by_substring` filters an input list of strings to only include strings that contain a specified substring.

Here is the complete implementation for the function:

```
'''python
from typing import List

def filter_by_substring(strings: List[str], substring: str) -> List[str]:
    """ Filter an input list of strings only for ones that contain given substring.
    """

```

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

to inefficiencies. Additionally, similar to Phidata, the code execution tool was seldom triggered (see Figure 15).

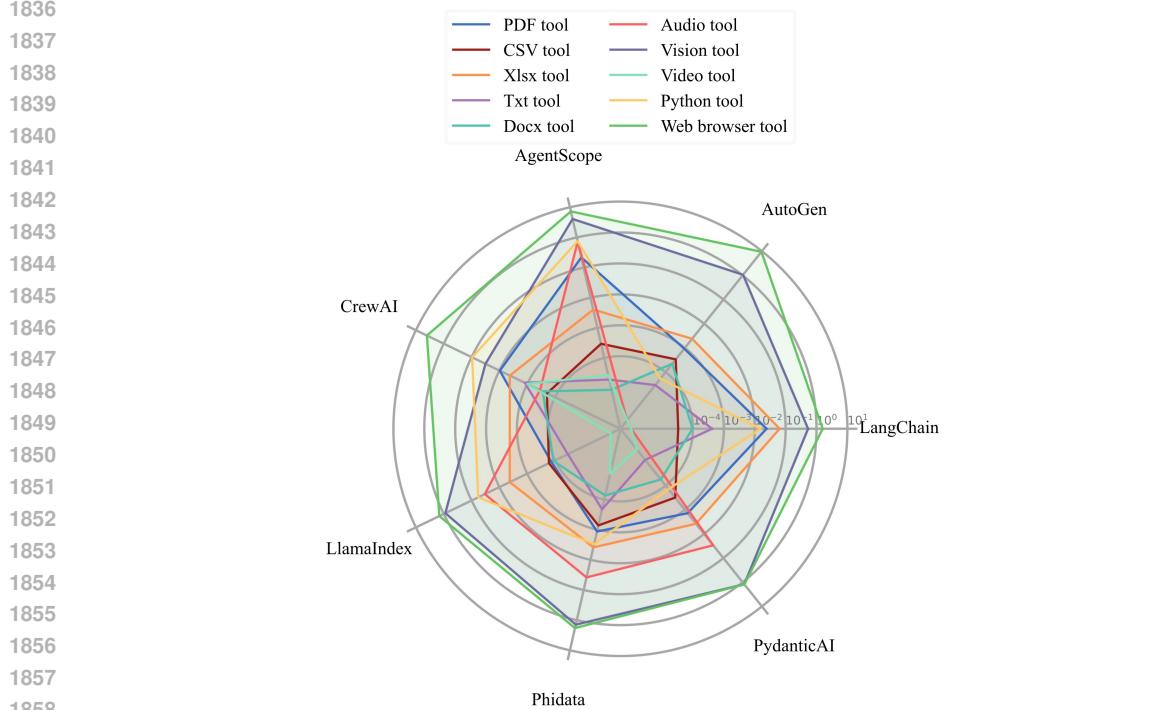
Furthermore, The MoA implementation in the PydanticAI framework is tool-based, and not all three models are invoked for every query. We observe that when the number of sub-agents is 3, 6, 9, 12, and 15, there were 232, 89, 229, 485, and 663 instances, respectively, where sub-agents were not invoked.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

E TOOL IMPLEMENTATION

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

E.1 SEARCH

E.1.1 AUTOGEN



```

189022     raise Exception(f"Error in API request: {response.status_code}")
189123     results = response.json().get("items", [])
189224     def get_page_content(url: str) -> str:
189325         try:
189426             response = requests.get(url, timeout=10)
189527             soup = BeautifulSoup(response.content, "html.parser")
189628             text = soup.get_text(separator=" ", strip=True)
189729             words = text.split()
189830             content = ""
189931             for word in words:
190032                 if len(content) + len(word) + 1 > max_chars:
190133                     break
190234                 content += " " + word
190335             return content.strip()
190436         except Exception as e:
190537             print(f"Error fetching {url}: {str(e)}")
190638             return ""
190739         enriched_results = []
190840         for item in results:
190941             body = get_page_content(item["link"])
191042             enriched_results.append(
191143                 {
191244                     "title": item["title"],
191345                     "link": item["link"],
191446                     "snippet": item["snippet"],
191547                     "body": body
191648                 }
191749             )
191850             time.sleep(1)
191951     return enriched_results

```

E.1.2 PYDANTICAI

```

19191     def google_search(query, num=None):
19202         """
19213             Make a query to the Google search engine to receive a list of results.
19224             Args:
19235                 query (str): The query to be passed to Google search.
19246                 num (int, optional): The number of search results to return.
19257             Defaults to None.
19268
19279             Returns:
192810                str: The JSON response from the Google search API.
192911
193012             Raises:
193113                 ValueError: If the 'num' is not an integer between 1 and 10.
193214             """
193315             try:
193416                 QUERY_URL_TMPL = ("https://www.googleapis.com/customsearch/v1?key"
193517                 ={key}&cx={engine}&q={query}")
193618                 url = QUERY_URL_TMPL.format(
193719                     key=os.environ['GOOGLE_KEY'],
193820                     engine=os.environ['GOOGLE_ENGINE'],
193921                     query=urlllib.parse.quote_plus(str(query))
194022                 )
194123                 if num is not None:
194224                     if not 1 <= num <= 10:
194325                         raise ValueError("num should be an integer between 1 and "
194426                         10, inclusive")
194527                         url += f"&num={num}"
194628                     response = requests.get(url)
194729                     return response.text
194830             except Exception as e:

```

```

1944 28     return f"Error: {e}"
1945

```

1946

1947

1948 E.2 PDF LOADER

```

1949
1950 1 def pdf_load(file_path: str) -> ServiceResponse:
1951 2     try:
1952 3         reader = PdfReader(file_path)
1953 4         text = ""
1954 5         for page in reader.pages:
1955 6             text += page.extract_text() + "\n"
1956 7         return ServiceResponse(status=ServiceExecStatus.SUCCESS, content=
1957 8             text)
1958 9     except Exception as e:
1959 10        return ServiceResponse(ServiceExecStatus.ERROR, str(e))
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

```

E.2 PDF LOADER

E.3 CSV READER

```

import pandas as pd

def csv_load(path:str)->ServiceResponse:
    try:
        df = pd.read_csv(path)
        csv_str = df.to_string(index=False)
        return ServiceResponse(status=ServiceExecStatus.SUCCESS, content=
            csv_str)
    except Exception as e:
        return ServiceResponse(ServiceExecStatus.ERROR, str(e))

```

E.4 XLSX READER

```

def xlsx_load(path:str)->ServiceResponse:
    try:
        excel_file = pd.read_excel(path, sheet_name=None)
        result = ""
        for sheet_name, df in excel_file.items():
            result += f"Sheet: {sheet_name}\n"
            result += df.to_string(index=False) + "\n\n"
        return ServiceResponse(status=ServiceExecStatus.SUCCESS, content=
            result.strip())
    except Exception as e:
        return ServiceResponse(ServiceExecStatus.ERROR, str(e))

```

E.5 TEXT FILE READER

```

import pandas as pd

def txt_load(path:str)->ServiceResponse:
    try:
        with open(path, 'r', encoding='utf-8') as f:
            txt_str = f.read()
        return ServiceResponse(status=ServiceExecStatus.SUCCESS, content=
            txt_str)
    except Exception as e:
        return ServiceResponse(ServiceExecStatus.ERROR, str(e))

```

1998
1999

E.6 DOCX READER

```

2000 1 from docx import Document
2001 2
2002 3 def docs_load(path:str) ->ServiceResponse:
2003 4     try:
2004 5         doc = Document(path)
2005 6         docx_str = "\n".join([para.text for para in doc.paragraphs])
2006 7         return ServiceResponse(status=ServiceExecStatus.SUCCESS, content=
2007 docx_str)
2008 8     except Exception as e:
2009 9         return ServiceResponse(ServiceExecStatus.ERROR, str(e))

```

2009

2010

E.7 MP3 LOADER

```

2012 1 import whisper
2013 2 from typing import cast
2014 3
2015 4 def load_audio(file):
2016 5     model = whisper.load_model(name="base")
2017 6     model = cast(whisper.Whisper, model)
2018 7     result = model.transcribe(str(file))
2019 8     return result["text"]

```

2020

2021

E.8 FIGURE LOADER

```

2022 1 from transformers import DonutProcessor, VisionEncoderDecoderModel
2023 2 import re
2024 3 from PIL import Image
2025 4
2026 5 def load_image(path):
2027 6     image = Image.open(path)
2028 7     processor = DonutProcessor.from_pretrained(
2029 8         "naver-clova-ix/donut-base-finetuned-cord-v2"
2030 9     )
2031 10    model = VisionEncoderDecoderModel.from_pretrained(
2032 11        "naver-clova-ix/donut-base-finetuned-cord-v2"
2033 12    )
2034 13    device = 'cpu'
2035 14    model.to(device)
2036 15    # prepare decoder inputs
2037 16    task_prompt = "<s_cord-v2>"
2038 17    decoder_input_ids = processor.tokenizer(
2039 18        task_prompt, add_special_tokens=False, return_tensors="pt"
2040 19    ).input_ids
2041 20    pixel_values = processor(image, return_tensors="pt").pixel_values
2042 21    outputs = model.generate(
2043 22        pixel_values.to(device),
2044 23        decoder_input_ids=decoder_input_ids.to(device),
2045 24        max_length=model.decoder.config.max_position_embeddings,
2046 25        early_stopping=True,
2047 26        pad_token_id=processor.tokenizer.pad_token_id,
2048 27        eos_token_id=processor.tokenizer.eos_token_id,
2049 28        use_cache=True,
2050 29        num_beams=3,
2051 30        bad_words_ids=[[processor.tokenizer.unk_token_id]],
2052 31        return_dict_in_generate=True,
2053 32    )
2054 33    sequence = processor.batch_decode(outputs.sequences)[0]
2055 34    sequence = sequence.replace(processor.tokenizer.eos_token, "").replace(
2056 35        processor.tokenizer.pad_token, ""

```

```

205236     )
205337     # remove first task start token
205438     text_str = re.sub(r"<.*?>", "", sequence, count=1).strip()
205539     return text_str
2056
2057
2058 E.9 VIDEO LOADER
2059

```

```

2060 1   import whisper
2061 2   from typing import cast
2062 3   from pydub import AudioSegment
2063 4   from pathlib import Path
2064 5
2065 6   def load_video(file):
2066 7       video = AudioSegment.from_file(Path(file), format=file[-3:])
2067 8       audio = video.split_to_mono()[0]
2068 9       file_str = str(file)[:-4] + ".mp3"
2069 10      audio.export(file_str, format="mp3")
2070 11      model = whisper.load_model(name="base")
2071 12      model = cast(whisper.Whisper, model)
2072 13      result = model.transcribe(str(file))
2073 14      return result["text"]
2074
2075
2076 E.10 DATA RETRIEVAL
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

```

```

1   def create_vector_db():
2       import faiss
3       import pickle
4       from sentence_transformers import SentenceTransformer
5       from data.mmlu import merge_csv_files_in_folder
6       dataset=merge_csv_files_in_folder(path to MMLU/dev)
7       docs = []
8       for item in dataset:
9           text = item[0].replace(",please answer A,B,C, or D.",",")+""
10          answer:{item[1]}."
11          docs.append(text)
12          embed_model = SentenceTransformer('all-MiniLM-L6-v2')
13          doc_embeddings = embed_model.encode(docs)
14          dimension = doc_embeddings.shape[1]
15          index = faiss.IndexFlatL2(dimension)
16          index.add(doc_embeddings)
17          faiss.write_index(index, "db/index.faiss")
18          with open("db/index.pkl", "wb") as f:
19              pickle.dump(docs, f)
20
21  def load_vector_db():
22      import faiss
23      import pickle
24      from sentence_transformers import SentenceTransformer
25      class db:
26          def __init__(self):
27              self.index = faiss.read_index("db/index.faiss")
28              with open("db/index.pkl", "rb") as f:
29                  self.docs = pickle.load(f)
30                  self.embed_model = SentenceTransformer('all-MiniLM-L6-v2')
31          def search(self, query, k=5):
32              query_embedding = self.embed_model.encode([query])
33              D, I = self.index.search(query_embedding, k)
34              return [self.docs[i] for i in I[0]]
35
36  return db()

```

2106
2107

E.11 PROBLEM SOLVER

```

2108 1     def twoSum(nums: List[int], target: int) -> List[int]:
2109 2         """
2110 3             Given an array of integers nums and an integer target, return indices
2111 4             of the two numbers such that they add up to target.
2112 5                 Args:
2113 6                     nums (List): an array of integers
2114 7                     target (Int): an integer target
2115 8                 Returns:
2116 9                     List[int]: indices of the two numbers such that they add up to
target.
211710                """
211811            try:
211912                n = len(nums)
212013                for i in range(n):
212114                    for j in range(i + 1, n):
212215                        if nums[i] + nums[j] == target:
212316                            return [i, j]
212417
212518            except Exception as e:
212619                return str(e)

212720
212821            def lengthOfLongestSubstring(s: str) -> int:
212922                """
213023                    Given a string s, find the length of the longest substring without
duplicate characters.
213124                    Arg:
213225                        s (String): a string
213326
213427                    Returns:
213528                        Int: the length of the longest substring without duplicate
characters.
213629                """
213730            try:
213831                left = 0
213932                right = 0
214033                max_len = 0
214134
214235                while right < len(s):
214336                    if s[right] in s[left:right]:
214437                        max_len = max(max_len, right-left)
214538                            left = s.index(s[right], left, right)+1
214639                        max_len = max(max_len, right-left+1)
214740                            right += 1
214841                return max_len
214942            except Exception as e:
215043                return str(e)

215144            def findMedianSortedArrays(nums1: List[int], nums2: List[int]) -> float:
215245                """
215346                    Given two sorted arrays nums1 and nums2 of size m and n respectively,
return the median of the two sorted arrays.
215447                    Args:
215548                        nums1 (List[int]): sorted array 1
215649                        nums2 (List[int]): sorted array 2
215750                    Returns:
215851                        float: the median of the two sorted arrays
215952                """
216053            try:
216154                m, n = len(nums1), len(nums2)

```

```

2160      59
2161      60
2162      61
2163      62
2164      63
2165      64
2166      65
2167      66
2168      67
2169      68
2170      69
2171      70
2172      71
2173      72
2174      73
2175      74
2176      75
2177      76
2178      77
2179      78
2180      79
2181      80
2182      81
2183      82
2184      83
2185      84
2186      85
2187      86
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
2999

```

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.