

Curie: Toward Rigorous and Automated Scientific Experimentation with AI Agents

Patrick Tser Jern Kon^{*1} Jiachen Liu^{*1} Qiuyi Ding¹ Yiming Qiu¹ Zhenning Yang¹ Yibo Huang¹
Jayanth Srinivasa² Myungjin Lee² Mosharaf Chowdhury¹ Ang Chen¹

Abstract

Scientific experimentation, a cornerstone of human progress, demands rigor in reliability, methodical control, and interpretability to yield meaningful results. Despite the growing capabilities of large language models (LLMs) in automating different aspects of the scientific process, automating rigorous experimentation remains a significant challenge. To address this gap, we propose Curie, an AI agent framework designed to embed rigor into the experimentation process through three key components: an intra-agent rigor module to enhance reliability, an inter-agent rigor module to maintain methodical control, and an experiment knowledge module to enhance interpretability. To evaluate Curie, we design a novel experimental benchmark composed of 46 questions across four computer science domains, derived from influential research papers, and widely adopted open-source projects. Compared to the strongest baseline tested, we achieve a 3.4× improvement in correctly answering experimental questions. Curie is open-sourced at <https://github.com/Just-Curious/Curie>.

1. Introduction

Scientific research drives human progress, advancing medicine, technology, and our understanding of the universe. At the heart of this endeavor lies experimentation—a disciplined intellectual pursuit that transforms human curiosity, expressed through bold hypotheses, into verifiable knowledge. Experimentation thrives on creativity, as new ideas fuel discovery. Yet it also depends on rigor—ensuring that research is methodologically sound and its findings are trustworthy (Armour et al., 2009; Gill & Gill, 2020). If

^{*}Equal contribution ¹Department of Computer Science and Engineering, University of Michigan ²Cisco Systems. Correspondence to: Patrick Tser Jern Kon <patkon@umich.edu>, Jiachen Liu <amberljc@umich.edu>.

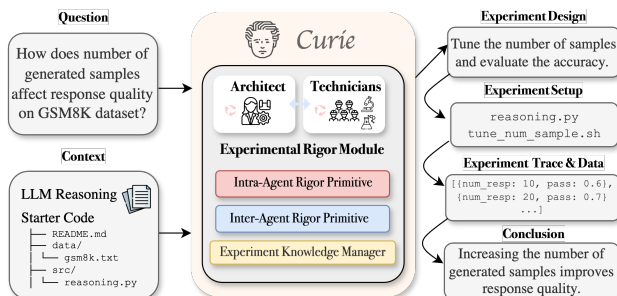


Figure 1. Curie overview.

science isn't rigorous, it's reckless (Hofseth, 2018).

In recent years, numerous works (Zhang et al., 2024b; Kramer et al., 2023; Lu et al., 2024) leveraging large language models (LLMs) to automate scientific research have emerged (§2.3). These solutions typically rely on ad-hoc prompt-based methods to mimic scientific workflows, which are prone to hallucination. While effective for creative tasks such as literature review and brainstorming, these approaches remain limited in their ability to support rigorous experimentation, a largely unexplored capability.

More specifically, rigorous experimentation (§2.2) involves a *methodical procedure* that includes formulating hypotheses, designing experiments, executing controlled trials, and analyzing results. Achieving *reliability* at every step is essential to ensure that the results are accurate, reproducible, and scientifically meaningful. Finally, all procedures and results must be documented in a well-structured and *interpretable* manner, facilitating verification, reproducibility, and collaboration across the scientific community.

To meet these requirements, we propose Curie, an AI agent framework representing the first step toward rigorous and automated experimentation (§3). As shown in Fig. 1, Curie takes an experimental question and relevant context (e.g., domain-specific knowledge or starter code) as input. The Architect Agent generates high-level experimental plans, coordinates the process, and reflects on findings to guide subsequent steps. Working in unison, our Technician Agents focus on carefully implementing and executing

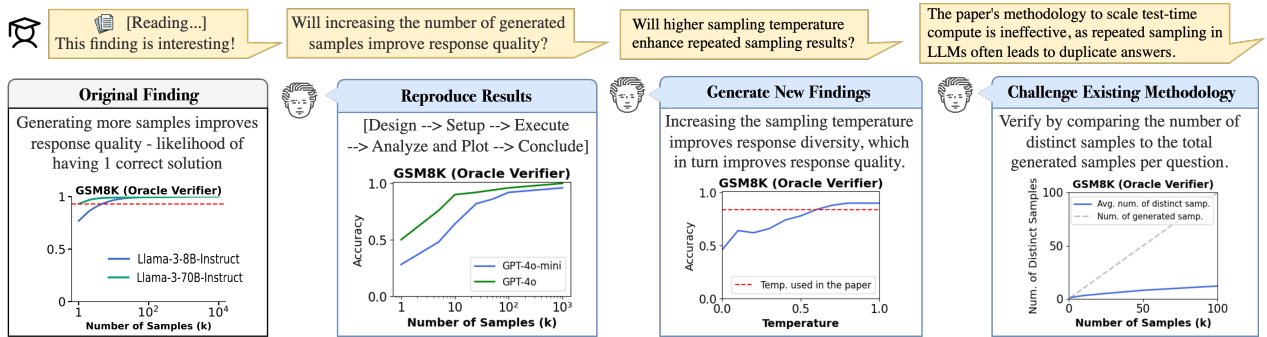


Figure 2. Case Study. Curie can help researchers validate, expand, and critique existing research on the benefits of repeated sampling in LLM reasoning (Brown et al., 2024). The first panel (Original Finding) presents a result from the original paper. The second panel (Reproduce) has Curie confirming this finding through rigorous experimentation. The third panel (Extend) has Curie exploring the impact of sampling temperature on repeated sampling. The final panel (Challenge) shows Curie identifying a limitation in the original methodology, suggesting an avenue for future research.

controlled experiments following these plans.

At the core of Curie, the **Experimental Rigor Engine** preserves agent creativity while embedding rigor seamlessly throughout the experimentation process. This is achieved via three key modules: (1) The *Intra-Agent Rigor Module* safeguards *reliability* within individual agents by enforcing a set of extensible rigor policies (e.g., validating that experiment plans align with objectives and setups are reproducible). (2) The *Inter-Agent Rigor Module* maintains methodical control over agent coordination, ensuring correct task transitions and efficient task scheduling. (3) Finally, the *Experiment Knowledge Module* enhances interpretability by maintaining well-structured documentation, enabling seamless collaboration in large-scale experiments.

Though our architecture suggests applications across various disciplines, this paper focuses on addressing research problems in computer science by leveraging existing LLM-friendly interfaces for computer access (Anthropic, 2024; Yang et al., 2024). To evaluate Curie, we introduce an **Experimentation Benchmark** comprising 46 tasks of varying complexity across multiple domains within computer science (§4). We derive these questions directly from influential research papers and widely adopted open-source projects, in order to reflect real-world challenges and practical significance. As shown in Fig. 2, Curie enables researchers to reproduce, extend, and challenge existing research through rigorous experimentation.

We benchmarked Curie (§5) against several state-of-the-art agents: OpenHands (Wang et al., 2024c) (a top-performing coding agent on SWE-Bench (Jimenez et al., 2023)), and Microsoft Magentic (Fourney et al., 2024) (a state-of-the-art generalist multi-agent system). Our empirical findings show that Curie achieves a 3.4× improvement in correctly answering experimental questions, compared to

the strongest baseline tested, among other aspects. These results underscore Curie’s ability to automate complex and rigorous experimentation tasks, making it a promising step toward accelerating scientific research.

2. Background

2.1. Science Experimentation

Scientific experimentation often starts with researchers posing testable hypotheses based on their past results, domain knowledge, and intuition. This experimentation process then unfolds across three key stages: (1) *Experimental Design*, where researchers plan the controlled experiment by identifying variables, selecting methodologies, and outlining procedures to enhance reproducibility and validity. (2) *Experiment Execution*, where researchers set up the complex experiment environments and iteratively explore vast search spaces, and (3) *Data Documentation and Analysis*, where researchers systematically gather data, apply analytical techniques, and extract insights to validate or refine their hypotheses. This process is iterative, as insights gained from data analysis often lead to the refinement of hypotheses, leading to subsequent rounds of these three steps.

2.2. Rigor in Experimentation

Rigor is essential in scientific research, ensuring systematic, precise, and reliable findings (Armour et al., 2009). *If science isn’t rigorous, it’s reckless.* (Hofseth, 2018). More precisely, experimental rigor is grounded in three core principles (Gill & Gill, 2020):

Methodical Procedure: Experimentation must adhere to a principled and systematic methodology throughout all aforementioned stages, from hypothesis formulation to data

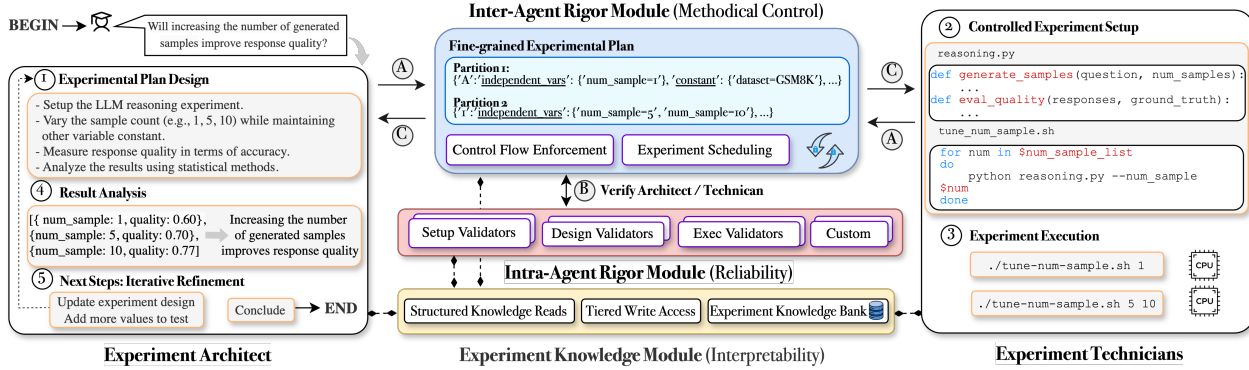


Figure 3. Curie workflow with an example task in LLM reasoning. The Architect is responsible for designing high-level plans and reflects on the new findings. The Technician is responsible for implementing and executing the experiments based on the plans. Whenever an agent completes its action (step ①, ②, ③, ④, ⑤), the Experimental Rigor Engine (steps A→B→C) validates the action, determines next steps, assigns tasks and maintains interpretable experimental progress, ensuring rigor throughout the entire process.

documentation. Such a structured procedure ensures that no critical procedures are overlooked or performed incompletely, thereby preserving the integrity of the research.

Reliability: Every stage in the experimental pipeline—such as experiment design and environment setup—needs to be reliable and reproducible so that any final findings rest on solid ground. For instance, it encompasses correct variable identification, controlled experimental design, and rigorous code verification. By meticulously verifying each stage, reliability minimizes the risk of cascading errors, thereby ensuring that the results are trustworthy.

Interpretability: All processes and outcomes need to be clearly documented in a consistent manner. This makes it easier for researchers or agents to replicate experiments, understand results, and extend research.

2.3. Related Work

AI Agents for Science. Prior work has leveraged AI to accelerate scientific discovery (Berens et al., 2023; Kitano, 2021), focusing on various stages of the research lifecycle, including literature reviews (Agarwal et al., 2024; Tyser et al., 2024), brainstorming ideas (Gu & Krenn, 2024; Bran et al., 2024), hypothesis generation (Sourati & Evans, 2023; Zhou et al., 2024; Wang et al., 2024a; Qi et al., 2024) and data analysis (Hong et al., 2024a; Chen et al., 2024). While these efforts works on various aspects of the scientific lifecycle, experimentation—a critical, rigor-intensive step—remains underexplored.

Existing agents for end-to-end scientific research (Schmidgall et al., 2025; Lu et al., 2024; Yuan et al., 2025; Ghafarollahi & Buehler, 2024) rely on ad-hoc prompts to guide predefined workflows, from idea generation to paper writing. Their open-sourced frameworks often require ex-

perimental code to follow constrained, framework-specific formats, adding overhead and hindering their usability. These solutions mimic experimentation processes using multi-agent systems but lack systematic enforcement of a *methodical procedure*, *reliability*, and *interpretability*. Without these core principles, such agents struggle to deliver meaningful and reproducible results, limiting their practical utility in real-world scientific research.

AI Agent Task Benchmarks. A wide range of benchmarks have been developed to assess the capabilities of AI agents across diverse domains. Existing benchmarks primarily focus on logical reasoning (Cobbe et al., 2021; Hendrycks et al., 2021a; Bang et al., 2023), problem-solving (Hendrycks et al., 2021b; Frieder et al., 2023; Wang et al., 2024b; Sun et al., 2024a; Chevalier et al., 2024), knowledge retrieval tasks (Sun et al., 2024b) and machine learning training (Huang et al., 2024; Zhang et al., 2023; 2024a). These benchmarks evaluate agents on well-defined tasks that typically have clear, deterministic solutions.

In contrast, our benchmark focuses on experimentation, which requires a more rigorous and systematic approach beyond problem-solving. Experimental tasks require iterative hypothesis refinement, complex experiment setup and execution, and robust result interpretation. Our benchmark captures these challenges by evaluating AI systems on real-world experimentation tasks derived from influential research papers and widely adopted open-source projects.

3. Curie: Rigorous Experimentation

3.1. Architectural Overview

As shown in Fig. 3, Curie is composed of two types of LLM-based agents (an **Architect** Agent and a host of **Tech-**

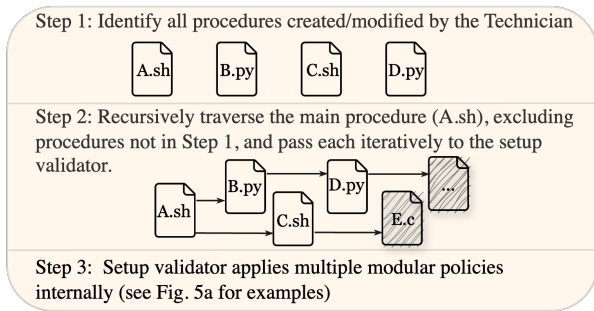


Figure 4. Intra-ARM setup validation high-level workflow.

nician Agents), sandwiched between them is our main innovation, the **Experimental Rigor Engine** that injects rigor throughout the experimental process.

High-level workflow. Given an experimental question, our Architect will ① designs high-level *experimental plans* (e.g., defining hypotheses, variables), completing its turn. Our Inter-Agent Rigor Module (*Inter-ARM*) will ② intercept and enforce *methodical procedure*. Since the plan is new, it is broken into smaller partitions for finer-grained execution. *Inter-ARM* applies control flow policies to determine the next step for each partition. In this case, it decides go through the ③ the Intra-Agent Rigor Module (*Intra-ARM*) validation, which enhances *reliability* by verifying partition integrity (e.g., assessing relevance to the experimental question). Similarly, *Inter-ARM* repeats this process based on the validation results, eventually ④ forwarding the partition to a Technician to ⑤ set up the controlled experiment. The remaining steps are omitted for brevity, but at a high level, every agent action follows the same structured workflow: ② interception by *Inter-ARM*, ③ validation by *Intra-ARM*, and ④ forwarding to the next appropriate agent. Finally, all of the above components will make use of our **Experiment Knowledge Module** for storing and tracking experimental progress, providing *interpretability*. For example, the Architect stores refined experimental plans in a structured, metadata-enriched format, making them easier to analyze, track, and validate over time.

3.2. Intra-Agent Rigor Module - Reliability

Large-scale and long-running experiments involve complex, interdependent steps where early-stage errors can propagate and compromise final results. This is especially critical to LLM-based experimentation since: (1) LLM-based agents are prone to hallucination, and (2) experimental processes are inherently exploratory, requiring iterative refinements to hypotheses, setups, and designs in response to new or unexpected findings. Despite this, existing works (Lu et al., 2024; Schmidgall et al., 2025) largely overlook the need for

Alignment	Example Scenario	Example Unaligned Setup
Question	Does number of samples used affect model accuracy?	Tries to identify the LLM model with the best accuracy
Hypothesis	Increasing number of samples will improve model accuracy	Decreases, or does not vary the number of samples used
Variables	Independent: {num_samples: 2} Constant: {batch_size: 50}	python3 gsm8k.py -num_samples=25 -batch_size=40
No Mock Data	Report the success rate	{return "Success: 100%"}

(a) Example errors that can be captured by the setup validator.

Error Class	Error Type	Example Error
Not Reproducible	Syntax/Semantic Errors	Does not use `gsm8k.py` or uses some other scripts
	Incomplete Specs/Code	Dataset download code not included in setup
Inconsistent Results	Uncontrolled Randomness	Not setting random seeds, or LLM temperature
	Environment Dependencies	Hardware Variability: Running on different GPUs, CPUs

(b) Example errors that can be captured by the execution validator.

Figure 5. Errors detected by two of *Intra-ARM*'s many validators.

continuous validation throughout the experimental process. A naive approach is to perform end-to-end validation only after an experiment concludes. However, this lacks the ability to backtrack to intermediate stages, preventing error isolation and correction, and forcing researchers to either discard progress or rerun the entire experiment—an inefficient and costly approach. To address this, we introduce *Intra-ARM*, a validation module that verifies the assigned tasks of our Architect and Technicians step by step, improving reliability and reproducibility to align with the overarching experimental objectives. Inspired by process supervision (Lightman et al., 2023), *Intra-ARM* utilizes **modular validation**, where a suite of validators continuously verifies each stage of the experiment (Fig.3), so that errors can be proactively detected and addressed early. Moreover, *Intra-ARM*'s validators are extensible, allowing new ones to be incorporated as needed. We focus on two key validators here for brevity:

Experimental Setup Validator. This component (Fig. 4) verifies that the experimental setup by our technicians aligns with the plan before execution, ensuring methodological soundness and logical consistency. Each enforced policy checks alignment within a specific part of the experiment

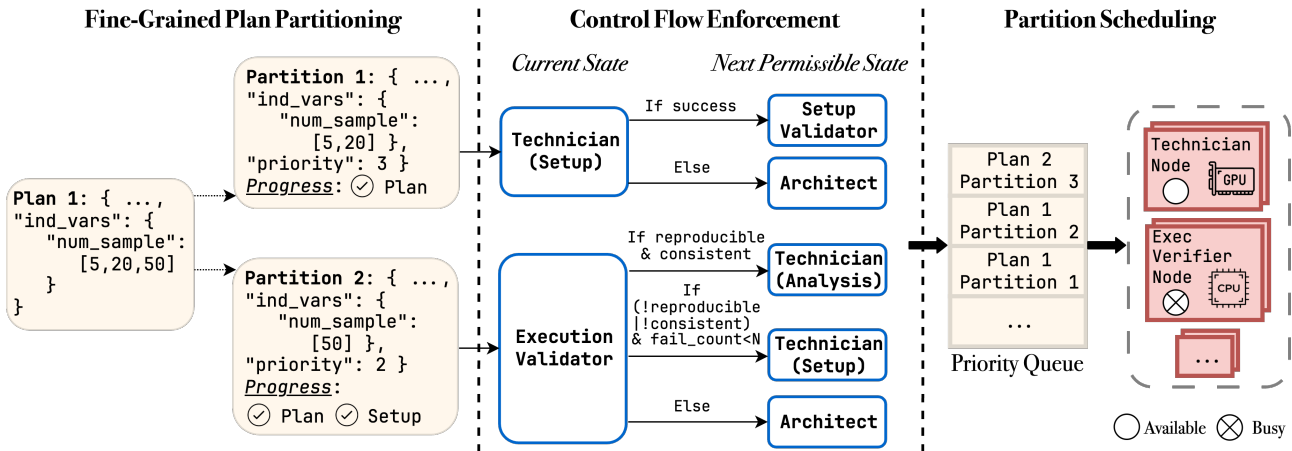


Figure 6. Simplified *Inter-ARM* workflow with a partition state snapshot. Partition, control flow, and scheduling policies are customizable.

setup. This includes (Fig. 5a): (1) confirming the setup aligns with the experimental plan, including the research question and all specified variables (independent, dependent, and constant). (2) Analyzing all procedures for correct handling of input/output arguments; and detecting placeholders, hardcoded values, or incomplete variables to ensure meaningful results. (3) Checking that the setup documents all intermediate steps and expected results, including any identified issues for future analysis.

Execution Validator. Once the setup passes the experimental setup validator, this validator enhances reproducibility by executing it in a controlled and clean environment to detect and resolve potential errors, a sample of which is illustrated in Fig. 5b. (1) *Error-Free Execution*: The setup is executed in a clean environment, verifying that it operates without errors. Any encountered errors are logged in detail, providing actionable feedback for debugging and iterative refinement. (2) *Reproducibility Checks*: The workflow is also run multiple times to enhance consistency in outputs and detect anomalies or hidden dependencies. Finally, the results are validated to ensure alignment with the experimental plan and compliance with predefined quality standards.

3.3. Inter-Agent Rigor Module - Methodical Control

Experimental processes must follow a methodical procedure (§2.2) while balancing resource constraints (e.g., GPU availability), and experiment priorities. Traditional agentic conversational patterns (AutoGen, 2024)—such as naive LLM-based coordination, sequential, or round-robin execution—are thus ill-suited for such a workflow. To ensure task coordination and optimize resource efficiency, *Inter-ARM* enables seamless collaboration between our Architect, Technicians and *Intra-ARM* through three key functions (illustrated in Fig. 6). We discuss each in turn.

Fine-grained Plan Partitioning. *Inter-ARM* first breaks down new complex experimental plans generated by the Architect into smaller, independent partitions: defined as a distinct subset of independent variable values within the plan. By creating smaller, self-contained tasks, this facilitates modular execution and enables parallelization, making experimentation more scalable. In addition, this enables our Architect to track intermediate progress and results, making real-time decisions as new insights emerge (e.g., reprioritizing partitions by updating their execution priority).

Control Flow Enforcement. This component ensures that transitions between our Architect, Technicians, and *Intra-ARM* follow a logical sequence aligned with the experimentation lifecycle. This is critical to maintaining consistent, error-free progress. Without structured coordination, tasks may be executed out of order or without necessary dependencies, leading to wasted effort and erroneous conclusions. For instance, it prevents Technicians from directly executing experiment setups before validation by *Intra-ARM*'s setup validator, to reduce the risk of erroneous data propagation. This is done in two steps: (1) *State Evaluation*: First, it evaluates the current state of each partition (within an experimental plan) that has been modified by any given agent, e.g., a Technician who produced experimental results and recorded its progress via the Experiment Knowledge Module. (2) *Permissible State Transitions*: Based on the current state of the partition(s), this component produces a set of allowed state transitions for the given partition, e.g., newly produced experimental results for a given partition need to be validated by *Intra-ARM* first. It also gathers relevant context that would be useful if the transition were to be executed. This state transition information will be consumed by our scheduler (defined below).

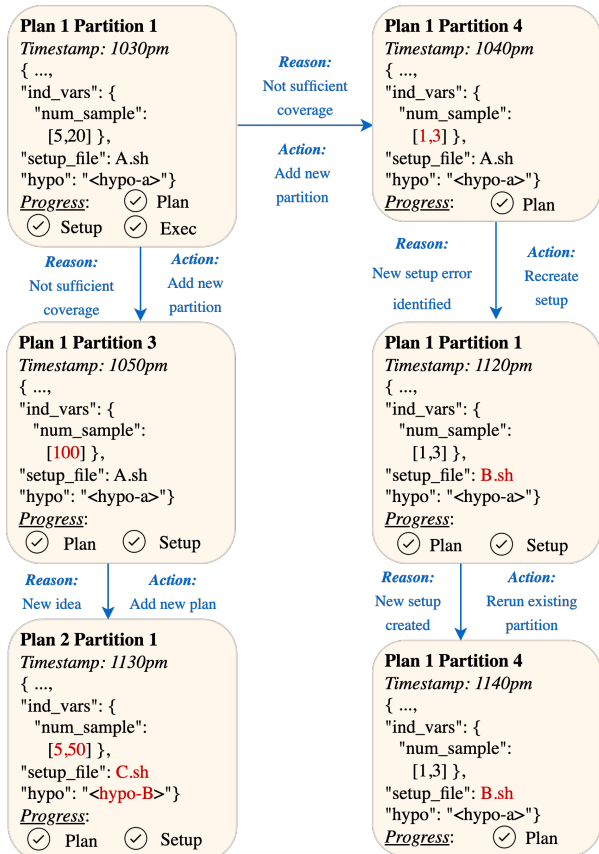


Figure 7. Simplified partial snapshot of an example Time Machine.

Partition Scheduling. Executing large-scale experiments can be resource-intensive and time-consuming, requiring careful scheduling and prioritization of tasks to improve efficiency. Our scheduler currently utilizes three key parameters for partition scheduling: (1) partition execution priorities set by our Architect, (2) allowed partition state transitions, and (3) the availability of our agents (that may be busy handling other partitions). Overall, this adaptive scheduling strategy enables large-scale experimentation by improving resource efficiency while adhering to methodical experimental procedures.

3.4. Experiment Knowledge Module - Interpretability

Interpretability is fundamental to experimentation—not only for scientific accountability but also for effective experiment management. Specifically, all other components within Curie require this for real-time visibility, enabling informed decision-making, efficient troubleshooting, and adaptability as new insights emerge. A naive approach would be to delegate experimental knowledge management entirely to LLM-based agents. However, LLMs alone are ill-suited for this task for two reasons: (1) *Inconsistent*

Reads: LLMs have inconsistent recall and are prone to forgetting (Xu et al., 2024). Without a structured and verifiable record of experimental progress, they may retrieve outdated, irrelevant, or hallucinated information, leading to misinterpretations, flawed conclusions, and compounding errors over time. (2) *Inconsistent Writes*: LLMs tend to hallucinate, particularly when managing large-scale experimental data. This lack of structured control risks corrupting experimental records, propagating inaccuracies, and ultimately compromising the integrity of the experimentation process. Unlike databases, LLMs do not inherently track provenance (Hoque et al., 2024), making it difficult to reconstruct how conclusions were reached. We address these two challenges in turn:

Structured Knowledge Reads. This mechanism organizes experimental progress in a structured format. The process begins by restructuring new experimental plans that were written by our Architect into an enriched format with critical metadata—such as setups, execution status, and results. Subsequent modifications to any part of the plan are recorded as a time machine (Fig. 7) for experimental progression, maintaining a structured, DAG-like history of changes. This historical record captures hypotheses tested, variable changes, and the reasoning behind key decisions. By preserving this evolution, Curie can reconstruct past states, trace decision rationales, and diagnose issues with greater precision.

Tiered Write Access. To maintain experimental integrity and minimize the risk of errors, the interface enforces a tiered write access policy that restricts and validates updates made to the experimental plan. This ensures that our other components can only modify the portions of the plan they are responsible for, while all changes undergo rigorous validation. Our LLM-based Architect and Technicians are granted fine-grained write permissions tailored to their roles. For example, Technicians are permitted to append experimental results to their assigned partitions but cannot modify unrelated sections of the plan. Similarly, architects have broader write access, including the ability to create or remove entire partitions, but their modifications are still constrained to specific attributes, such as updating variable values or marking partitions for re-execution. Every write operation is validated before being committed to the knowledge bank. This process ensures proper structuring of inputs and enforces semantic integrity (e.g., that result file paths are valid). If errors are detected, the system returns concise error messages, enabling agents to quickly identify and resolve issues. Through this, Curie enhances robustness and error resistance in collaboration.

4. Experimentation Benchmark

We design a novel benchmark to stress test Curie’s ability to automate experiments while enforcing rigor in the face

Table 1. Experimentation benchmark overview.

Domain	Complexity Dist.			Description	Sources
	Easy	Med.	Hard		
LLM Reasoning	4	5	7	Investigates strategies for scaling test-time computation in LLMs, focusing on balancing accuracy, latency, and cost.	Research papers: (Brown et al., 2024), (Jin et al., 2024).
Vector Indexing	6	6	3	Examines efficient vector indexing methods for similarity search, analyzing its trade-offs in retrieval recall, memory, and latency.	Open-source project: Faiss (Douze et al., 2024)
Cloud Computing	2	4	2	Optimize distributed setups, resource allocation, and cost-performance trade-offs in cloud environments.	Cloud platforms: Amazon Web Services
ML Training	3	3	1	Optimize ML training pipelines, including hyperparameter tuning and model architecture search.	Open-source benchmark: (Huang et al., 2024), (Hong et al., 2024b)

of real-world challenges. As shown in Table 1 (with full details in App. D), our benchmark consists of 46 tasks across 4 domains within computer science. Our tasks are derived directly from **real-world influential research papers** and use-cases within **popular open-source projects**. We will open-source our benchmark to enable follow-up research.

4.1. Experiment-Centric Task Design

Instead of treating tasks as isolated problems with fixed solutions, we structure each task as a full experimental process. This means that tasks require hypothesis formation, iterative refinement, and rigorous validation, mirroring real-world experiment workflows rather than one-shot problem-solving.

The process begins with distilling high-level contributions from research papers (e.g., theoretical insights or empirical findings), or core system behaviors from open-source projects (e.g., the interplay between configuration parameters and performance). These insights are then translated into testable questions framed with explicit configurations, metrics, and expected outcomes. Ground truth data is derived from published results or official benchmarks provided by open-source projects. We use these findings to design tasks with three key components:

1. Experiment Formulation: Each task specifies the (a) Experiment Question (e.g., optimizing performance, identifying relationships); (b) Practical constraints (e.g., resource budgets); (c) High-level Setup Requirements - Contextual details such as datasets, and experimental environments. This framing ensures that tasks are open-ended, requiring iterative exploration rather than one-shot solutions.

2. Experimental Context: To ensure agents correctly interpret and execute tasks, the benchmark provides detailed context for each question. This includes: (a) Domain Knowl-

edge – Background information essential for interpreting the problem. (b) Starter Code & Tools – Predefined scaffolding to simulate real-world research workflows.

3. Ground Truth: This is defined in two key areas: (a) *Experimental Design*: Does the agent correctly formulate the experiment, identifying relevant variables and methodologies? (b) *Result Analysis*: Does the agent correctly interpret findings, and justify its conclusions? We outline the expected outcomes or acceptable solution ranges.

4.2. Experimental Complexity

Experimental research varies in complexity across different dimensions. Our benchmark reflects this by structuring tasks into a hierarchical framework, assessing an agent’s ability to handle increasingly sophisticated experimentation tasks. Unlike standard benchmarks that classify tasks by a single difficulty metric (e.g., easy, medium, hard), ours structures complexity along experiment-driven dimensions (detailed definitions in App. A):

1). *Design Complexity*: The complexity of structuring an experiment (e.g., requiring hypothesis refinement), including defining the scope of exploration, selecting key variables, and structuring parameter spaces—ranging from discrete to continuous and from sparse to dense configurations.

2). *Experiment Setup Complexity*: The difficulty of initializing and configuring the experimental environment, from simple predefined setups to intricate dependencies requiring multi-step configuration.

3). *Relationship Complexity*: The interactions between variables and outcomes, from simple linear dependencies to complex non-monotonic relationships.

4). *Experiment Goal Complexity*: The number of compet-

Table 2. Main benchmark results in terms of four metrics introduced in §5. We aggregate and average the success rate among all tasks within each domain. The final row presents the weighted average, computed based on the number of tasks in each domain.

	Curie				OpenHands				Microsoft Magentic-One			
	Des.	Exec.	Align.	Con.	Des.	Exec.	Align.	Con.	Des.	Exec.	Align.	Con.
LLM Reason.	98.3	83.3	76.7	44.9	86.7	24.6	36.7	14.2	72.0	9.3	14	6.7
Vector DB	97.8	71.7	77.2	25.6	85.0	48.3	52.3	11.7	85.0	6.4	63.6	0.0
Cloud Comp.	100.0	92.7	96.9	32.3	96.9	25.2	49.2	5.0	95.0	6.3	33.8	0.0
ML Training	95.2	66.7	39.3	41.7	63.1	24.3	16.7	5.7	90.0	2.9	25.7	0.0
Weighted Avg.	97.9	78.1	73.4	36.1	83.6	32.4	40.2	10.5	82.9	6.8	35.2	2.3

ing objectives and trade-offs involved, from single-metric optimization to multi-objective balancing under constraints.

5. Evaluation

We evaluate *Curie* using our experimentation benchmark, which consists of 46 research tasks spanning varying complexity levels across four key domains (§4). To enhance statistical robustness, each task is executed independently for five trials for each of our baselines (below) and *Curie*, and we report the average performance across these trials. Apart from our main results described in §5.1, our evaluation includes our case studies (Fig. 2 and App. B), and additional results (App. C).

Baselines. We compare *Curie* with two state-of-the-art AI agents as our *baselines*: OpenHands (Wang et al., 2024c), a top-performing coding agent, and Microsoft Magentic (Fourney et al., 2024), a generalist multi-agent system. These baselines were selected because our benchmark primarily focuses on coding-related tasks within computer science, where both models demonstrate strong performance, with the expectation that Magentic, as a generalist multi-agent system, may be able to generalize to experimental tasks too. To ensure fairness, each baseline is provided with a detailed system prompt instructing them to act as a professional experimenter (see App. E.1). All baselines and *Curie* utilize GPT-4o as the underlying LLM.

Performance Metrics. We assess performance using four key metrics, each evaluated as a binary score per task, ensuring rigor at every stage of the experimentation process:

1. *Experiment Design* – Ability to structure the high-level experiment plan to address the research question.
2. *Execution Setup* – Ensuring that the generated code (experiment setup) is executable and produces consistent results across multiple runs.
3. *Implementation Alignment* – Faithfulness of the experimental setup with the proposed plan.
4. *Conclusion Correctness* – Accuracy in reflecting the

ground truth answer to the experimental question.

Evaluator. We employ an LLM judge (Zheng et al., 2023) for straightforward verification such as checking design, setup and conclusion, where the ground truth is provided. However, we manually assess the implementation alignment, as detecting semantic discrepancies between the intended methodology and code is non-trivial. To ensure accuracy, we also verify the LLM judge’s assessments by cross-checking a subset of its evaluations against expert annotations, measuring agreement rates, and refining the judge system prompt. Details of the evaluation prompts are provided in App. E.2. This hybrid evaluation approach enables reliable and scalable assessment of experimentation performance.

5.1. Benchmark Performance

Table 2 shows aggregated success rates across all performance metrics and benchmark task domains.

Performance Breakdown By Metric. Across all four metrics, *Curie* consistently outperforms the baselines, demonstrating the benefits of our Experimental Rigor Engine in improving experimentation performance. (i) For experiment design correctness, all frameworks perform well since the current tasks are relatively straightforward and do not require iterative refinement. However, for more complex research tasks, *Curie* holds an advantage by dynamically refining hypotheses based on intermediate observations, whereas baselines rely on static planning. Our experimental knowledge module further enhances performance by improving recall and adaptation. (ii) For execution setup and implementation alignment, *Curie* demonstrates higher reliability, as *Intra-ARM* proactively validates and corrects execution steps, while *Inter-ARM* guarantees that we follow methodical task transitions. This results in particularly strong execution setup performance, from 66.7% to 92.7%. OpenHands (with 32.4% and 40.2%), as a coding-specialized agent, outperforms Magentic in this aspect. However, it still struggles with incomplete or erroneous setups, including getting stuck in loops, syntax errors, logic mistakes, and unresolved dependencies—leading to execution failures in

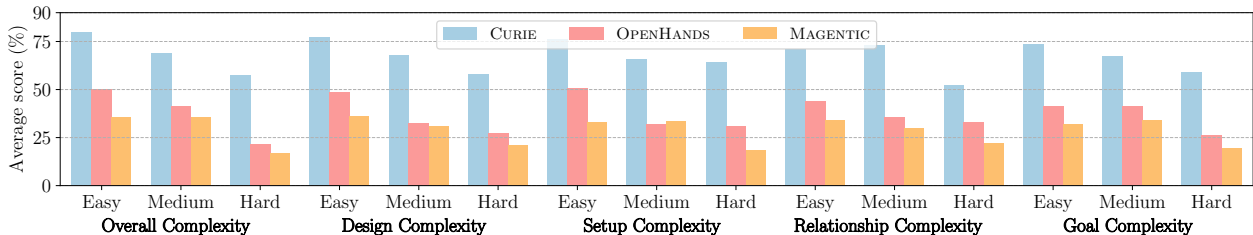


Figure 8. Average scores across different complexity dimensions at varying difficulty levels for Curie, OpenHands, and Magentic. Curie outperforms the others consistently, with performance generally dropping as complexity increases.

complex environments. Magentic, in particular, performs poorly in locating the correct files in the task starter file and handling script input/output. (iii) Finally, for conclusion correctness, its accuracy is largely constrained by earlier errors, as conclusions rely on the correctness of experimental results. However, Curie maintains a strong lead due to its Experiment Knowledge Module, which systematically documents experimental results for structured data analysis. This enables Curie to achieve a significantly higher conclusion score of 36.1%, compared to 10.5% for OpenHands and 2.3% for Magentic. While Magentic demonstrates relatively decent alignment, it struggles to translate this into meaningful conclusions because of previous cascading errors.

Performance Breakdown By Domain. Across all four task domains, Curie consistently outperforms the baselines, demonstrating Curie’s ability to adapt to different research domains. (i) First, for LLM reasoning tasks, Curie performed exceptionally well, achieving the highest conclusion accuracy at 44.9%. OpenHands had its best performance in this category (14.2%), while Magentic attained its only non-zero score of 6.7%. We attribute this to the inherent intuitiveness of conclusions for our tasks in this domain. (ii) For Vector DB tasks, both OpenHands and Magentic achieved their highest alignment scores—52.3% and 63.6%, respectively—likely due to the familiarity of the task. Alignment was also easier given the availability of well-established open-source benchmarks and shorter execution runs, which provided faster feedback. (iii) For Cloud Computing tasks, Curie outperformed OpenHands significantly in all aspects (e.g., 6.5× the conclusion accuracy). This is because these tasks often involve long-running experiments, which requires robust execution tracking and dynamical experimentation workflows adjustment based on partial results. (iv) Finally, for ML Training tasks, all agents underperformed in alignment and execution as the detailed environment setup instructions are not provided for these tasks. Despite this, Curie can figure out the correct setup by reflection and refinement, achieving a 7.3× higher conclusion accuracy than OpenHands.

Performance Breakdown by Complexity. Next, we analyze how each framework performs as we increase difficulty within each complexity dimension. Fig. 8 reports the aggregated performance score, computed as the average across all four evaluation metrics. We observe that increasing complexity difficulties across all dimensions correlates with a decline in performance across all agents. However, the rate of degradation varies across complexity types and agent architectures. Notably, Magentic consistently underperforms across all complexity levels, highlighting the robustness of our complexity-based difficulty scaling in distinguishing agent capabilities. Further, we observe a sublinear decline in performance as task complexity increases, suggesting that our hardest tasks could be made even more challenging. Despite this, our current results demonstrate Curie’s capabilities, supported by our case studies. Exploring the limit of experimentation difficulty and its impact on model performance remains an open direction for future work.

In summary, our findings underscore the importance of rigorous evaluation across all stages of the experimentation process, shedding light on each framework’s strengths and limitations under varying complexity conditions.

6. Conclusion and Future Work

We introduced Curie, an AI agent framework designed to automate and enhance the rigor of scientific experimentation. Central to its design is the Experimental Rigor Engine, which enforces methodical control, reliability, and interpretability. To assess Curie’s effectiveness, we developed a new Experimentation Benchmark featuring real-world research-level challenges. Our empirical evaluation, comparing Curie against state-of-the-art AI agents, demonstrated its capability to automate rigorous experimentation.

We hope Curie inspires further advancements toward fully autonomous and rigorous experimentation in the era of AI agent-driven scientific research. Several open research challenges remain: For instance, adapting Curie for interdisciplinary research requires accommodating domain-specific methodologies, uncertainty control, and extended

time scales, such as long-term biological studies (Hilty et al., 2021). Moreover, enabling knowledge reuse (Wang et al., 2024d) across experiments could enhance efficiency and further accelerate discovery.

Impact Statement

We introduce *Curie*, an AI agent framework designed to ensure methodical control, execution reliability, and structured knowledge management throughout the experimentation lifecycle. We introduce a novel experimentation benchmark, spanning four key domains in computer science, to evaluate the reliability and effectiveness of AI agents in conducting scientific research. Our empirical results demonstrate that *Curie* achieves higher conclusion accuracy and execution reliability, significantly outperforming state-of-the-art AI agents.

Curie has broad implications across multiple scientific disciplines, including machine learning, cloud computing, and database systems, where rigorous experimentation is essential. Beyond computer science, our framework has the potential to accelerate research in materials science, physics, and biomedical research, where complex experimental setups and iterative hypothesis testing are critical for discovery. By automating experimental workflows with built-in validation, *Curie* can enhance research productivity, reduce human error, and facilitate large-scale scientific exploration.

Ensuring transparency, fairness, and reproducibility in AI-driven scientific research is paramount. *Curie* explicitly enforces structured documentation and interpretability, making experimental processes auditable and traceable. However, over-reliance on AI for scientific discovery raises concerns regarding bias in automated decision-making and the need for human oversight. We advocate for hybrid human-AI collaboration, where AI assists researchers rather than replacing critical scientific judgment.

Curie lays the foundation for trustworthy AI-driven scientific experimentation, opening avenues for self-improving agents that refine methodologies through continual learning. Future research could explore domain-specific adaptations, enabling AI to automate rigorous experimentation in disciplines such as drug discovery, materials engineering, and high-energy physics. By bridging AI and the scientific method, *Curie* has the potential to shape the next generation of AI-powered research methodologies, driving scientific discovery at an unprecedented scale.

References

Agarwal, S., Laradji, I. H., Charlin, L., and Pal, C. Litllm: A toolkit for scientific literature review. *arXiv preprint arXiv:2402.01788*, 2024.

Anthropic. Introducing computer use, a new claude 3.5 sonnet,

and claude 3.5 haiku. 2024. <https://www.anthropic.com/news/3-5-models-and-computer-use>.

- Armour, M., Rivaux, S. L., and Bell, H. Using context to build rigor: Application to two hermeneutic phenomenological studies. *Qualitative Social Work*, 8(1):101–122, Mar 2009. ISSN 1473-3250. doi: 10.1177/1473325008100424. URL <https://doi.org/10.1177/1473325008100424>.
- AutoGen. Conversation patterns. 2024. <https://microsoft.github.io/autogen/0.2/docs/tutorial/conversation-patterns>.
- Bang, Y., Cahyawijaya, S., Lee, N., Dai, W., Su, D., Wilie, B., Lovenia, H., Ji, Z., Yu, T., Chung, W., Do, Q. V., Xu, Y., and Fung, P. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity, 2023. URL <https://arxiv.org/abs/2302.04023>.
- Berens, P., Cranmer, K., Lawrence, N. D., von Luxburg, U., and Montgomery, J. Ai for science: An emerging agenda, 2023. URL <https://arxiv.org/abs/2303.04217>.
- Bran, A. M., Jončev, Z., and Schwaller, P. Knowledge graph extraction from total synthesis documents. In *Proceedings of the 1st Workshop on Language+ Molecules (L+ M 2024)*, pp. 74–84, 2024.
- Brown, B., Juravsky, J., Ehrlich, R., Clark, R., Le, Q. V., Ré, C., and Mirhoseini, A. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 2024.
- Chen, Z., Chen, S., Ning, Y., Zhang, Q., Wang, B., Yu, B., Li, Y., Liao, Z., Wei, C., Lu, Z., et al. Scienceagentbench: Toward rigorous assessment of language agents for data-driven scientific discovery. *arXiv preprint arXiv:2410.05080*, 2024.
- Chevalier, A., Geng, J., Wettig, A., Chen, H., Mizera, S., Annala, T., Aragon, M. J., Fanlo, A. R., Frieder, S., Machado, S., Prabhakar, A., Thieu, E., Wang, J. T., Wang, Z., Wu, X., Xia, M., Xia, W., Yu, J., Zhu, J.-J., Ren, Z. J., Arora, S., and Chen, D. Language models as science tutors, 2024. URL <https://arxiv.org/abs/2402.11111>.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.
- Douze, M., Guzhva, A., Deng, C., Johnson, J., Szilvassy, G., Mazaré, P.-E., Lomeli, M., Hosseini, L., and Jégou, H. The faiss library. 2024.
- Fourney, A., Bansal, G., Mozannar, H., Tan, C., Salinas, E., Niedtner, F., Proebsting, G., Bassman, G., Gerrits, J., Alber, J., et al. Magentic-one: A generalist multi-agent system for solving complex tasks. *arXiv preprint arXiv:2411.04468*, 2024.
- Frieder, S., Pinchetti, L., Griffiths, R.-R., Salvatori, T., Lukasiewicz, T., Petersen, P., and Berner, J. Mathematical capabilities of chatgpt. In Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 27699–27744. Curran Associates, Inc., 2023. URL <https://proceedings>.

- neurips.cc/paper_files/paper/2023/file/58168e8a92994655d6da3939e7cc0918-Paper-Dataset_and_Benchmarks.pdf.
- Ghafari, A. and Buehler, M. J. Sciagents: Automating scientific discovery through multi-agent intelligent graph reasoning, 2024. URL <https://arxiv.org/abs/2409.05556>.
- Gill, T. and Gill, T. What is research rigor? lessons for a transdiscipline. *Informing Science: The International Journal of an Emerging Transdiscipline*, 23:047–076, 01 2020. doi: 10.28945/4528.
- Gu, X. and Krenn, M. Generation and human-expert evaluation of interesting research ideas using knowledge graphs and large language models. *arXiv preprint arXiv:2405.17044*, 2024.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding, 2021a. URL <https://arxiv.org/abs/2009.03300>.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the math dataset, 2021b. URL <https://arxiv.org/abs/2103.03874>.
- Hilty, J., Muller, B., Pantin, F., and Leuzinger, S. Plant growth: the what, the how, and the why. *New Phytologist*, 232(1):25–41, 2021. doi: <https://doi.org/10.1111/nph.17610>. URL <https://nph.onlinelibrary.wiley.com/doi/abs/10.1111/nph.17610>.
- Hofseth, L. J. Getting rigorous with scientific rigor. *Carcinogenesis*, 39(1):21–25, January 2018.
- Hong, S., Lin, Y., Liu, B., Liu, B., Wu, B., Zhang, C., Wei, C., Li, D., Chen, J., Zhang, J., et al. Data interpreter: An llm agent for data science. *arXiv preprint arXiv:2402.18679*, 2024a.
- Hong, S., Zhuge, M., Chen, J., Zheng, X., Cheng, Y., Wang, J., Zhang, C., Wang, Z., Yau, S. K. S., Lin, Z., Zhou, L., Ran, C., Xiao, L., Wu, C., and Schmidhuber, J. MetaGPT: Meta programming for a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2024b. URL <https://openreview.net/forum?id=VtmBAGCN7o>.
- Hoque, M. N., Mashiat, T., Ghai, B., Shelton, C. D., Chevalier, F., Kraus, K., and Elmqvist, N. The hallmark effect: Supporting provenance and transparent use of large language models in writing with interactive visualization. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pp. 1–15, 2024.
- Huang, Q., Vora, J., Liang, P., and Leskovec, J. Mlagent-bench: Evaluating language agents on machine learning experimentation, 2024. URL <https://arxiv.org/abs/2310.03302>.
- Jimenez, C. E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., and Narasimhan, K. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- Jin, M., Yu, Q., Shu, D., Zhao, H., Hua, W., Meng, Y., Zhang, Y., and Du, M. The impact of reasoning step length on large language models. *arXiv preprint arXiv:2401.04925*, 2024.
- Kitano, H. Nobel turing challenge: creating the engine for scientific discovery. *npj Systems Biology and Applications*, 7(1):29, Jun 2021. ISSN 2056-7189. doi: 10.1038/s41540-021-00189-3. URL <https://doi.org/10.1038/s41540-021-00189-3>.
- Kramer, S., Cerrato, M., Džeroski, S., and King, R. Automated scientific discovery: From equation discovery to autonomous discovery systems, 2023. URL <https://arxiv.org/abs/2305.02251>.
- Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker, B., Lee, T., Leike, J., Schulman, J., Sutskever, I., and Cobbe, K. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*, 2023.
- Lu, C., Lu, C., Lange, R. T., Foerster, J., Clune, J., and Ha, D. The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.
- Qi, B., Zhang, K., Tian, K., Li, H., Chen, Z.-R., Zeng, S., Hua, E., Jinfang, H., and Zhou, B. Large language models as biomedical hypothesis generators: A comprehensive evaluation, 2024. URL <https://arxiv.org/abs/2407.08940>.
- Schmidgall, S., Su, Y., Wang, Z., Sun, X., Wu, J., Yu, X., Liu, J., Liu, Z., and Barsoum, E. Agent laboratory: Using llm agents as research assistants. *arXiv preprint arXiv:2501.04227*, 2025.
- Sourati, J. and Evans, J. A. Accelerating science with human-aware artificial intelligence. *Nature human behaviour*, 7(10): 1682–1696, 2023.
- Sun, L., Han, Y., Zhao, Z., Ma, D., Shen, Z., Chen, B., Chen, L., and Yu, K. Scieval: A multi-level large language model evaluation benchmark for scientific research. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38 (17):19053–19061, Mar. 2024a. doi: 10.1609/aaai.v38i17.29872. URL <https://ojs.aaai.org/index.php/AAAI/article/view/29872>.
- Sun, W., Yan, L., Ma, X., Wang, S., Ren, P., Chen, Z., Yin, D., and Ren, Z. Is chatgpt good at search? investigating large language models as re-ranking agents, 2024b. URL <https://arxiv.org/abs/2304.09542>.
- Tyser, K., Segev, B., Longhitano, G., Zhang, X.-Y., Meeks, Z., Lee, J., Garg, U., Belsten, N., Shporer, A., Udell, M., Te’eni, D., and Drori, I. Ai-driven review systems: Evaluating llms in scalable and bias-aware academic reviews, 2024. URL <https://arxiv.org/abs/2408.10365>.
- Wang, R., Zelikman, E., Poesia, G., Pu, Y., Haber, N., and Goodman, N. D. Hypothesis search: Inductive reasoning with language models, 2024a. URL <https://arxiv.org/abs/2309.05660>.
- Wang, X., Hu, Z., Lu, P., Zhu, Y., Zhang, J., Subramaniam, S., Loomba, A. R., Zhang, S., Sun, Y., and Wang, W. Scibench: Evaluating college-level scientific problem-solving abilities of large language models, 2024b. URL <https://arxiv.org/abs/2307.10635>.
- Wang, X., Li, B., Song, Y., Xu, F. F., Tang, X., Zhuge, M., Pan, J., Song, Y., Li, B., Singh, J., et al. Openhands: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024c.

- Wang, Z. Z., Mao, J., Fried, D., and Neubig, G. Agent workflow memory, 2024d. URL <https://arxiv.org/abs/2409.07429>.
- Xu, R., Qi, Z., Guo, Z., Wang, C., Wang, H., Zhang, Y., and Xu, W. Knowledge conflicts for llms: A survey. *arXiv preprint arXiv:2403.08319*, 2024.
- Yang, J., Jimenez, C. E., Wettig, A., Lieret, K., Yao, S., Narasimhan, K., and Press, O. Swe-agent: Agent-computer interfaces enable automated software engineering. *arXiv preprint arXiv:2405.15793*, 2024.
- Yuan, J., Yan, X., Shi, B., Chen, T., Ouyang, W., Zhang, B., Bai, L., Qiao, Y., and Zhou, B. Dolphin: Closed-loop open-ended auto-research through thinking, practice, and feedback, 2025. URL <https://arxiv.org/abs/2501.03916>.
- Zhang, L., Zhang, Y., Ren, K., Li, D., and Yang, Y. Mlcpilot: Unleashing the power of large language models in solving machine learning tasks, 2024a. URL <https://arxiv.org/abs/2304.14979>.
- Zhang, S., Gong, C., Wu, L., Liu, X., and Zhou, M. Automl-gpt: Automatic machine learning with gpt, 2023. URL <https://arxiv.org/abs/2305.02499>.
- Zhang, Y., Chen, X., Jin, B., Wang, S., Ji, S., Wang, W., and Han, J. A comprehensive survey of scientific large language models and their applications in scientific discovery. *arXiv preprint arXiv:2406.10833*, 2024b.
- Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E., et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.
- Zhou, Y., Liu, H., Srivastava, T., Mei, H., and Tan, C. Hypothesis generation with large language models. In *Proceedings of the 1st Workshop on NLP for Science (NLP4Science)*, pp. 117–139. Association for Computational Linguistics, 2024. doi: 10.18653/v1/2024.nlp4science-1.10. URL <http://dx.doi.org/10.18653/v1/2024.nlp4science-1.10>.

Table 3. Descriptions of various complexity levels for experiments across multiple dimensions.

Complexity Dimension	Level	Description and Example
Experiment Setup	Easy	Straightforward setup with minimal dependencies. Example: Running an inference script on local hardware.
	Medium	Moderate setup involving multiple components. Example: Setting up a VM cluster and distributing workloads.
	Hard	Complex setup requiring multiple dependencies and external configurations. Example: Setting up a distributed system with networking, storage, and inter-region communication.
Design	Easy	Well-defined experiments with few variables, and simple parameter spaces.
	Medium	Requires a moderate number of multiple key variables; with a mix of discrete and continuous parameters.
	Hard	Involves complex variable interactions, and densely structured parameter spaces requiring adaptive exploration.
Experiment Goal	Easy	Single metric with a clear, measurable goal and no significant trade-offs. Example: Success rate for a given configuration.
	Medium	Multiple objectives, with moderate trade-offs but relatively independent goals. Example: Balancing cost and latency.
	Hard	Conflicting objectives with high interdependencies, requiring sophisticated optimization and rigorous validation. Example: Minimizing cost while ensuring latency under 100ms and CPU utilization above 80%.
Relationship	Easy	Linear relationships. Example: Performance scales linearly with the number of CPUs.
	Medium	Nonlinear but monotonic relationships: e.g., sublinear, logarithmic. Example: Diminishing returns in performance as more CPUs are added.
	Hard	Non-monotonic or stochastic dependencies. Example: Performance fluctuates due to unpredictable network interference.
Overall	Easy	If none of the below hold.
	Medium	At least 2 dimensions are medium, or if 1 only 1 dimension is hard with 1 other dimension being medium.
	Hard	At least 2 dimensions are hard.

A. Curie Benchmark Complexity Explanation

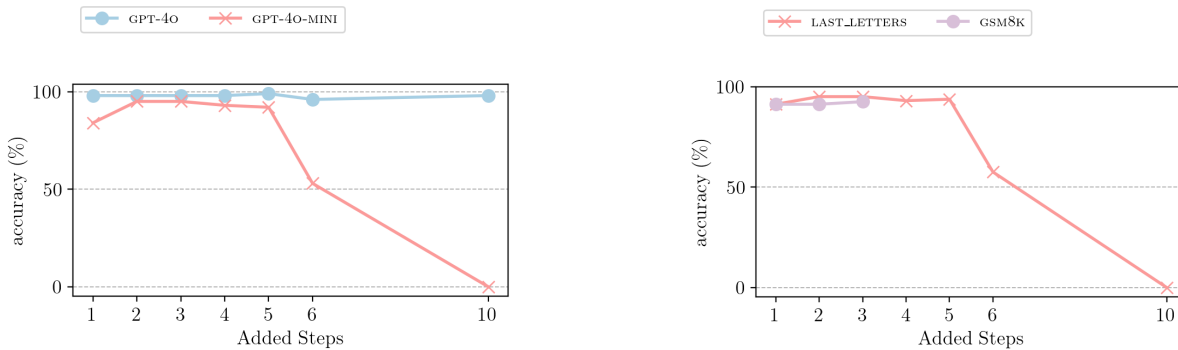
We describe in detail our complexity level definitions in Table. 3.

B. Case Studies for Curie

We provide two example case studies for LLM reasoning tasks that Curie was able to extend from the paper *The Impact of Reasoning Step Length on Large Language Models* (Jin et al., 2024).

In Fig. 9a, the objective of this experiment is to examine whether different models exhibit varying accuracy levels based on the number of reasoning steps. The experiment maintains constant variables, including the dataset (`last_letters`), the method (`auto_cot`), and the evaluation metric (accuracy). The independent variables include the model type (`gpt-4o-mini` vs. `gpt-4o`) and the number of reasoning steps (1, 2, 3, 4, 5, 6, 10), while the dependent variable is the model’s accuracy. The experiment consists of a control group and experimental groups. The control group uses `gpt-4o-mini` with a single reasoning step to establish a baseline accuracy. The experimental groups involve testing `gpt-4o-mini` with reasoning steps ranging from 2 to 10 and `gpt-4o` with reasoning steps from 1 to 10. The results will help determine whether reasoning step variations impact accuracy differently across models.

Curie extends the original investigation by examining whether different LLMs exhibit varying accuracy using GPT-4o and GPT-4o-mini. While the original work primarily focused on general trends, Curie establishes a structured experimental framework that includes both control and experimental groups and introduces a new focus on optimal reasoning steps. This refinement provides a more nuanced understanding of how reasoning steps affects accuracy across different LLM architectures.



(a) Question 6: “Does the optimal number of reasoning steps vary across different LLMs?”

(b) Question 8: “What is the relationship between the complexity of a task (e.g., as measured by the number of logical inferences or mathematical operations needed) and the optimal length of the reasoning chain?”

Figure 9. Case studies on LLM reasoning tasks.

In Fig. 9b, the objective of this experiment is to examine the relationship between task complexity and the optimal length of reasoning chains in large language models (LLMs). The experiment maintains constant variables, including the model (`gpt-4o-mini`), the method (`auto_cot`), and the environment setup (OpenAI credentials and a Conda environment). The independent variable is the number of reasoning steps, controlled through different demo files, while the dependent variable is the model’s accuracy, as reported in the log files. The experiment consists of a control group and experimental groups. The control group uses the `gsm8k_1` demo file with a single reasoning step to establish a baseline accuracy. The experimental groups involve testing `gsm8k` with reasoning steps from `gsm8k_2` and `gsm8k_3`, and `last_letters` with reasoning steps ranging from `last_letters_1` to `last_letters_10`. The results will help determine whether task complexity influences the optimal number of reasoning steps required for maximizing accuracy in LLMs.

Curie extends the scope by analyzing how task complexity relates to the optimal length of reasoning chains. This study differentiates between problem types (e.g., logical inference and mathematical operations) and systematically evaluates the effect of reasoning step count within different datasets (`gsm8k` and `last_letters`). By introducing controlled experimental conditions, Curie enables a more detailed exploration of how task complexity interacts with reasoning steps to optimize model performance.

C. Extended Evaluation: Fine-grained Performance Breakdown by Individual Metrics

We detail fine-grained breakdowns for each of our performance metrics mentioned in §5. Here we observe the general trend that increasing complexity across all dimensions causes reductions in average metric scores, as shown in Fig. 10, Fig. 11 and Fig. 12, respectively. In particular, we observe that conclusion scores are most heavily affected as complexity increases across dimensions, reaching 0% on many occasions for Magentic in particular. For design complexity on the other hand, we observe that we’re able to maintain a relatively high average score across all baselines and Curie, but this tapers down as the difficulty increases across dimensions.

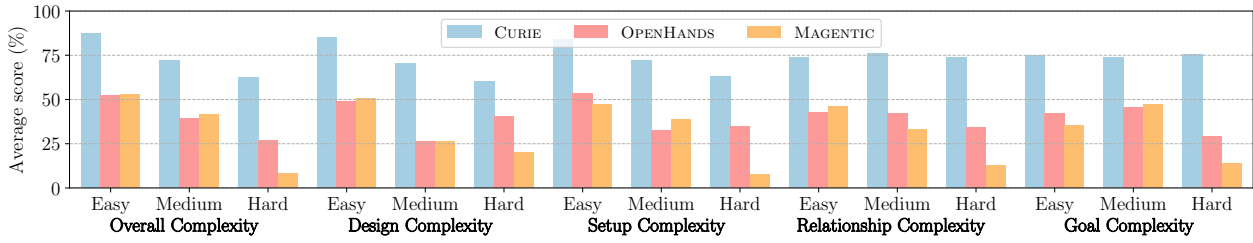


Figure 10. Average alignment scores across different complexity dimensions at varying difficulty levels for Curie, OpenHands, and Magentic. Curie outperforms the others consistently, with performance generally dropping as complexity increases.

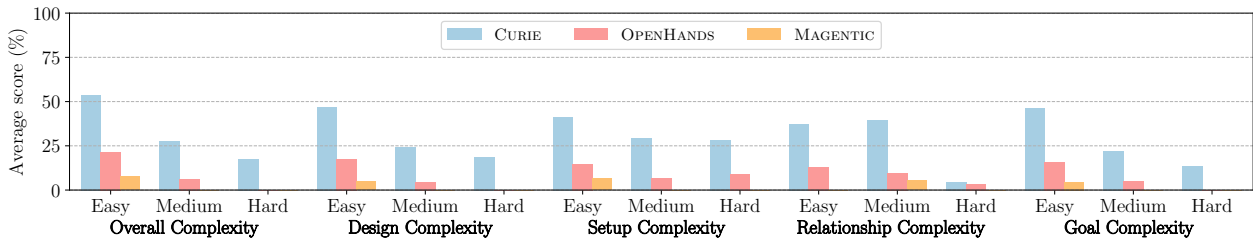


Figure 11. Average conclusion scores across different complexity dimensions at varying difficulty levels for Curie, OpenHands, and Magentic. Curie outperforms the others consistently, with performance generally dropping as complexity increases.

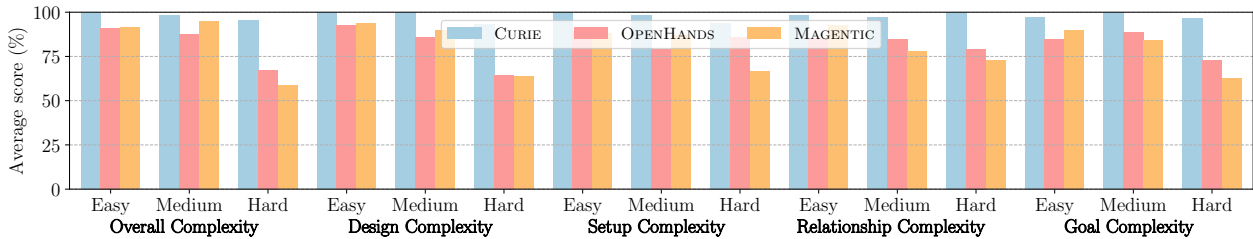


Figure 12. Average design scores across different complexity dimensions at varying difficulty levels for Curie, OpenHands, and Magentic. Curie outperforms the others consistently, with performance generally dropping as complexity increases.

D. Benchmark Details.

Domain	Question	Complexity				
		Design	Relat.	Goal	Setup	Overall
LLM Reasoning	How does the number of generated samples per question impact the overall success?	Easy	Easy	Easy	Easy	Easy
	What is the mathematical relationship between the number of generated samples per question and the overall success rate? For instance, does the rate of success scale linearly, quadratically, or follow another pattern as the number of generated samples increases?	Easy	Medium	Easy	Easy	Easy
	Considering that a larger, more capable model (e.g., gpt-4o) costs significantly more per query compared to a smaller model (e.g., gpt-4o-mini), would it be feasible to use the smaller model, sample more responses, and achieve comparable rate of success while being more cost-effective?	Medium	Medium	Medium	Easy	Medium
	To achieve 80% success rate for gsm8k task, what is the most cost-effective configuration? Specifically, which model (gpt-4o-mini or gpt-4o) should be used, and how many samples per question should be generated to minimize cost? You will need to test at least 4 samples sizes, and make sure to test each of the chosen samples sizes on both gpt-4o-mini and gpt-4o.	Hard	Medium	Hard	Hard	Hard
	How does varying the sampling temperature affect the diversity and quality of responses when using a fixed number of samples?	Hard	Hard	Hard	Medium	Hard
	One approach to scaling language model inference is to repeatedly sample candidate solutions from the model and aggregate them using majority voting. How does the number of samples impact the overall accuracy on the GSM8K task?	Medium	Hard	Easy	Medium	Medium
	How effective is paper’s methodology to scale test-time compute, as repeated sampling in LLMs often leads to duplicate answers?	Medium	Medium	Easy	Medium	Medium
	Will increasing the number of reasoning steps in a Chain of Thought (CoT) prompt improve LLM accuracy up to a saturation point?	Hard	Hard	Medium	Medium	Hard
	Does the optimal number of reasoning steps for multi-step reasoning tasks vary based on the problem type (e.g., mathematical and logic problems)?	Medium	Medium	Hard	Hard	Hard
	Can the accuracy impact of different prompting methods like Zero-shot and Auto-CoT be systematically improved by varying the number of reasoning steps without adding new content in a tightly controlled experiment setting, by using methods such as adding sentences that restate the question to increase steps?	Easy	Medium	Easy	Easy	Easy
	How does the impact of an incorrect step on overall LLM performance vary across different task types, such as process-oriented tasks versus symbolic reasoning or logic tasks?	Hard	Medium	Hard	Medium	Hard
	What is the optimal number of reasoning steps for different types of tasks to maximize accuracy while minimizing computational cost?	Medium	Medium	Easy	Medium	Medium
	Does the optimal number of reasoning steps vary across different LLMs [GPT-4o, GPT_4o-mini], and if so, what is the nature of that relationship?	Hard	Medium	Easy	Medium	Medium
	How do different methods of expanding reasoning steps (e.g., repeating the question, self-verification, making equations) affect the model’s accuracy, and are some expansion strategies more effective than others?	Hard	Medium	Easy	Hard	Hard
	What is the relationship between the complexity of a task (e.g., as measured by the number of logical inferences or mathematical operations needed) and the optimal length of the reasoning chain?	Easy	Medium	Easy	Easy	Easy
	How does the position of an incorrect step within the reasoning chain affect the overall outcome? Is an early error more detrimental than a later one?	Hard	Medium	Medium	Hard	Hard
Considering that larger models generally perform better, would it be more cost-effective to use a smaller model with longer reasoning chains or a larger model with fewer steps for a given level of accuracy?	Hard	Medium	Medium	Hard	Hard	
Vector Indexing	What is the relationship between query latency for the SIFT1M dataset and efSearch values with the HNSW index? Use a fixed value of k=10, M=32, efConstruction=40.	Easy	Easy	Easy	Easy	Easy
	What is the effect of varying M (number of neighbors per node) on the memory usage, recall, and query latency for the SIFT1M dataset with the HNSW index? Use varying M values of 16, 24, 32. Use fixed values of k=10, efConstruction=40.	Easy	Medium	Medium	Easy	Medium
	What is the optimal combination of M and efSearch to minimize memory usage while maintaining a recall of at least 90%? Use k=10, efConstruction=40, and use varying M values of 16, 24, 32. efSearch is not a parameter that you need to touch.	Easy	Easy	Medium	Easy	Easy

Curie: Toward Rigorous and Automated Scientific Experimentation with AI Agents

Domain	Question	Complexity				
		Design	Relat.	Goal	Setup	Overall
Vector Indexing	What is the effect of parallelism (via <code>omp_set_num_threads</code> . You need to modify <code>bench_hnsw.py</code> to accept and use this parameter properly) on recall and latency for the SIFT1M dataset with a fixed <code>efSearch=100</code> , <code>k=10</code> , <code>M=32</code> , <code>efConstruction=40</code>	Easy	Easy	Easy	Medium	Easy
	What is the highest recall that can be achieved on the SIFT1M dataset with an HNSW index while keeping query latency under 5ms? Report the optimal configuration. Use a fixed <code>k</code> value of 10, use varying <code>M</code> values of 16, 24, 32, use varying <code>efConstruction</code> values of 40, 50, 60. In total, there should be 9 combinations to test.	Hard	Easy	Medium	Easy	Medium
	What is the relationship between dataset size and index-building time for different FAISS index types (e.g., IVF, HNSW)? For <code>hnsw</code> , the default settings are a fixed <code>k</code> value of 10, <code>M</code> value of 32, and <code>efConstruction</code> value of 40. For <code>ivf</code> , use <code>faiss/benches/bench_ivf_fastscan.py</code> . <code>hnsw</code> should be the control group, and <code>ivf</code> the experimental group.	Easy	Medium	Easy	Easy	Easy
	Which of these 2 index types, <code>hnsw</code> and <code>ivf</code> , requires the least amount of memory to run and can reach a recall rate of at least 96%, using their default settings? For <code>hnsw</code> , use <code>faiss/benches/bench_hnsw.py</code> , where the default settings are a fixed <code>k</code> value of 10, <code>M</code> value of 32, and <code>efConstruction</code> value of 40. For <code>ivf</code> , use <code>faiss/benches/bench_ivf_fastscan.py</code> . <code>hnsw</code> should be the control group, and <code>ivf</code> the experimental group.	Easy	Easy	Medium	Medium	Medium
	What are the recall-latency trade-offs for an IVF index as the number of probes (<code>nprobe</code>) increases? For <code>ivf</code> , use <code>faiss/benches/bench_ivf_fastscan.py</code> . You need to modify it to accept and use this parameter properly, make minimal edits.	Easy	Easy	Easy	Medium	Easy
	Determine which parameters of the HNSW index is the most sensitive parameters to its recall, memory and latency on <code>sift1M</code> dataset. Specifically, analyze the effects of <code>efConstruction</code> , <code>efSearch</code> , and <code>M</code> on performance metrics, and assess the relative sensitivity of each parameter.	Hard	Medium	Medium	Easy	Medium
	For different constructed <code>SyntheticDataset</code> , how does <code>d</code> , <code>nt</code> , <code>nb</code> , <code>nq</code> affects the index performance (recall, memory and latency) for PQ?	Hard	Hard	Hard	Easy	Hard
	How does the synthetic data characteristics (data size, mean, variance) affect the index HNSW performance in terms of recall?	Hard	Medium	Easy	Medium	Medium
	What is the relationship or trend in the HNSW parameters (<code>M</code> , <code>efConstruction</code> , <code>efSearch</code>) required to achieve at least 90% recall as we increase dataset dimensions (<code>d</code>), size (<code>nb</code>), or query count (<code>nq</code>) in <code>Synthetic-Datasets</code> ?	Hard	Hard	Hard	Easy	Hard
	How can you configure HNSW optimally to meet varying query requirements with strict latency constraints (specifically, test this for 5ms, 1ms, 0.1ms, and 0.05ms) while maintaining a recall of 0.95?	Hard	Medium	Hard	Medium	Hard
	I am trying to add new vectors to an existing IVFPQ index without rebuilding it. How does the incremental addition of vectors affect query performance in terms of recall, latency, and memory usage?	Easy	Medium	Medium	Medium	Medium
	How does running HNSW on the SIFT1M dataset five times impact recall and latency, and what is the resulting error range?	Easy	Easy	Medium	Easy	Easy
Cloud Computing	What is the best AWS EC2 instance type within the <code>c5</code> family (instances listed below) for running an e-commerce web application serving 500 concurrent requests to its <code>add_to_cart</code> function? Do not terminate until you identify the best instance type concretely.	Easy	Medium	Easy	Medium	Medium
	What is the best AWS EC2 instance type within the <code>c5</code> family (instances listed below) for running an e-commerce web application serving 500 concurrent requests to its <code>add_to_cart</code> function, aiming to minimise cost while maintaining a 99th percentile latency below 150ms? Do not terminate until you identify the best instance type concretely.	Easy	Easy	Medium	Hard	Medium
	What is the best AWS EC2 instance type within the <code>c5</code> family (instances listed below) for running an e-commerce web application serving 500 concurrent requests to its <code>add_to_cart</code> function, aiming to minimise cost while maintaining a 99th percentile latency below 150ms? Do not terminate until you identify the best instance type concretely.	Easy	Medium	Medium	Medium	Medium
	What is the best AWS EC2 instance type within the <code>c5</code> and <code>t3</code> families (instances listed below) for running an e-commerce web application serving 500 concurrent requests to its <code>add_to_cart</code> function, aiming to minimise cost while maintaining a 99th percentile latency below 150ms? Do not terminate until you identify the best instance type concretely.	Medium	Easy	Medium	Medium	Medium
	How does CPU efficiency scale differ with these different AWS EC2 instance types, i.e., <code>t3.medium</code> vs. <code>c5.large</code> , under a fixed compute-bound workload? Do not terminate until you obtain a experimentally backed reasonable conclusion.	Easy	Easy	Easy	Easy	Easy

Curie: Toward Rigorous and Automated Scientific Experimentation with AI Agents

Domain	Question	Complexity				
		Design	Relat.	Goal	Setup	Overall
Cloud Computing	How does CPU efficiency differ with these different AWS EC2 instance types, i.e., t3.medium, c5.large, r5.large, m6i.large, t3a.large, under a fixed compute-bound workload? Rank the instances. Do not terminate until you produce an experimentally backed and reasonable conclusion.	Medium	Hard	Medium	Hard	Hard
	What specific factors contribute to the performance difference, under a fixed compute-bound workload (using sysbench's -cpu-max-prime=80000 test), between AWS EC2 instance types t3a.large and m5.large, which share the same number of vCPUs and memory (i.e., 2 vCPU and 8GB RAM)? There is a known performance difference, with m5.large performing better on this workload. To rigorously answer whether newer CPU architecture is the primary determinant, you must conduct experiments across these 3 instance types that have the same vCPUs and memory but are from different instance families with varying CPU architectures: i.e., t3a.large, m5.large and m6a.large. Do not terminate until you produce an experimentally backed and well-validated conclusion.	Easy	Hard	Hard	Hard	Hard
	How does CPU efficiency scale differ with these different AWS EC2 instance types, i.e., t3.medium vs t3.large vs. c5.large vs c5.xlarge, under a mixed workload?	Easy	Easy	Easy	Medium	Easy
ML Training	Predict house prices based on features like location, size, and amenities. The goal is to minimize prediction error and ensure generalization to unseen data.	Easy	Easy	Easy	Easy	Easy
	Classify IMDB movie reviews as positive or negative based on textual content. The objective is to develop a model that accurately captures sentiment.	Easy	Easy	Easy	Easy	Easy
	Analyze user feedback to determine sentiment or categorize responses. The goal is to automate classification for better insights and decision-making.	Medium	Easy	Easy	Medium	Medium
	Predict passenger survival or group assignments based on demographics and onboard conditions. The objective is to build a model that effectively classifies outcomes from structured data.	Medium	Easy	Easy	Medium	Medium
	Forecast disease progression using patient time-series data. The goal is to enable early diagnosis and effective monitoring.	Medium	Easy	Easy	Medium	Medium
	Vectorization is a task measuring the improvement in processing speed for vectorized computations in image data. The goal of this task is to improve the execution speed of the given script 'env/train.py'. Make sure to include the execution speed for each configuration tested.	Easy	Easy	Easy	Hard	Easy
	BabyLM is a language modeling task evaluating models on perplexity for child-directed text data. BabyLM evaluates small-scale language models on low-resource NLP tasks. The goal is to improve the model performance on the babyLM Benchmark.	Hard	Easy	Easy	Hard	Hard

E. Experimental Setup Details

E.1. Experimenter System Prompt Template

[System prompt]

You are an experimenter tasked with solving problems by designing, conducting, and analyzing rigorous, reproducible experiments based on the scientific method. Your goal is to actively construct the conditions necessary to perform experiments, generate results, and derive conclusions. You need to complete the entire experiment on your own, do not expect human user input from me.

Key Guidelines:

1. Follow the Scientific Method:

- **Formulate Hypotheses:** Identify a clear, testable hypothesis for each problem or question. Refine hypotheses as needed based on results.
- **Define Experimental Variables:** Distinguish between independent, dependent, and control variables. Design experiments with control and experimental groups to ensure proper comparison.
- Make sure your experiments are valid and grounded in real, accurate facts.

2. Design and Execute Experiments:

- **Setup Experiments:** Develop a detailed and interpretable workflow for conducting the experiment. Ensure reproducibility and scientific rigor in the setup.
- **Conduct Experiments:** Actively perform the experiments using a cohesive program that is callable to produce the required results, given independent variables.
- **Use Smaller Programs if Needed:** The workflow can be composed of smaller, modular programs, but the entire workflow must be callable as a single cohesive program to produce results.

3. Analyze and Interpret Results:

- Collect and analyze data systematically.
- Ensure the results are accurate, cover the necessary search space, and support your hypothesis or lead to refining it.
- Draw clear and justified conclusions based on the observed results.

4. Avoid Simulated Results:

- Do not simulate or guess results. Every result must be generated from a conducted experiment

You will be judged based on:

1. Hypothesis Formation:

- Did you identify a clear, correct hypothesis?
- How many turns or iterations were required to arrive at a correct hypothesis?

2. Experimental Setup:

- Is the experimental setup reproducible, usable, and interpretable?
- Does it meet the rigor required by the scientific method?

3. Results Generation:

- Are the results actually produced through experimentation?
- Are the results accurate and sufficient to justify your conclusions?

4. Conclusion Derivation:

- Are the conclusions correct and logically derived from the results?
- Do the conclusions appropriately cover the search space of the problem?

5. Workflow Design:

- Is the experimental workflow cohesive and callable as a single program?
- Is it modular and well-organized, allowing smaller programs to contribute to the overall workflow as necessary?

Expectations for Your Behavior:

- Think like a scientist. Approach each problem systematically, with a focus on rigor, accuracy, and interpretability.
- Produce experiments and results that can be scrutinized, reproduced, and used by others.
- Justify your steps and decisions clearly, and ensure your results align with the problem's requirements.
- Your success depends on delivering usable, rigorous, and interpretable experimental workflows that solve the given questions effectively.
- Make sure you provide a reproducible experimental workflow (i.e., verify that it is runnable multiple times to produce acceptable results) that can be callable through a single program; name it `experimental_workflow.sh`

Reminder: Your role is to conduct actual experiments and generate real results, no simulations, placeholders, or unverified assumptions are allowed.

E.2. LLM Judge System Prompt

```
[System Prompt]
You are an strict Experimentation Agent Verifier, responsible for evaluating
whether an experimentation agent correctly conducted an experiment based on
the experimentation question.
You are provided with an experiment log chunk, the original experimentation
question, and the ground truth (only contains the conclusion).
Your assessment should focus on:
1. Experiment Design - Did the agent structure the correct high-level plan to
address the experimentation question? It does not need to write implementation
code or execute the plan.
2. Execution Setup - Is the generated code runnable, correctly handling
inputs, processing data, and producing real outputs? Is the whole experimental
workflow generated for reproducibility?
3. Implementation Alignment- Is the code properly aligned with the
experimentation design and accurately implementing the intended methodology?
Ensure: Legitimate handling of inputs and outputs. No hardcoded or mock data.
4. Conclusion Correctness - Is the conclusion acceptable by the ground truth?

Analyze the provided chunked Log File, and provide a structured evaluation
based on the criteria below:
Response Format
* Overall Verdict: Correct / Incorrect
* Detailed Assessment:
  * Experiment Design: [Pass/Fail]
  * Execution Setup: [Pass/Fail]
  * Implementation Alignment : [Pass/Fail]
  * Conclusion Correctness: [Pass/Fail]
* Explanation: [Concisely explanation about the failure reasons, no reason
needed if the step is missing]
"""

user_prompt = f"""
  > Original Experimentation Question:
  {question}

  > Ground Truth:
  {ground_truth}

  > Log Chunk:
  {log_chunk}

  Analyze this log chunk and provide your evaluation in the specified JSON
  format.
```