

FIRE-Bench: Evaluating Agents on the Rediscovery of Scientific Insights

Zhen Wang^{1*} Fan Bai^{2*} Zhongyan Luo^{1*} Jinyan Su³ Kaiser Sun² Xinle Yu¹ Jieyuan Liu¹ Kun Zhou¹
Claire Cardie³ Mark Dredze² Eric P. Xing^{4,5} Zhiting Hu¹

 Website: <https://firebench.github.io>

Abstract

Autonomous agents powered by large language models (LLMs) promise to accelerate scientific discovery, but rigorously evaluating their capacity for *verifiable* discovery remains a central challenge. Existing benchmarks face a trade-off: they either rely on LLM-as-judge evaluations of automatically generated papers, raising concerns about measurement validity and circularity, or focus on optimizing isolated performance metrics that serve as coarse proxies for scientific insight. To address this, we introduce FIRE-BENCH (Full-cycle Insight Rediscovery Evaluation), a benchmark that reframes evaluation by tasking agents with the rediscovery of established findings from recent, high-impact ML research. Agents are provided only with the high-level research question from a published study and must autonomously design experiments, implement code, execute their plans, and derive conclusions grounded in empirical evidence. We evaluate a suite of state-of-the-art agents with frontier model backbones (e.g., GPT-5) on FIRE-BENCH. Results reveal substantial limitations in current systems: even the strongest agents exhibit low rediscovery success rates, high variance across runs, and recurring failure modes spanning experimental design, execution, and evidence-based reasoning. Overall, FIRE-BENCH provides a rigorous and diagnostic framework for measuring progress toward AI agents capable of reliable scientific discovery.

1. Introduction

The emergence of autonomous agents powered by large language models (LLMs) holds the promise of accelerating scientific discovery at an unprecedented scale. These

“AI researchers” are increasingly capable of automating discrete stages of the research lifecycle, from literature synthesis (Zheng et al., 2025; Schmidgall & Moor, 2025), hypothesis generation (Baek et al., 2024; Si et al., 2024), to coding (Tian et al., 2024; Chan et al., 2024), experimentation (Kon et al., 2025), and data analysis (Majumder et al.; Gu et al., 2024). However, a fundamental challenge lies in rigorously evaluating their capacity for genuine scientific discovery. Validating novel outcomes often requires resource-intensive, real-world verification, such as wet-lab experiments or large-scale human expert studies. This evaluation bottleneck becomes particularly acute for agents designed to automate the full research cycle—from an initial question to a final, empirically-grounded conclusion (Lu et al., 2024; Yamada et al., 2025; Schmidgall et al., 2025). Assessing the validity and integrity of an entire research trajectory, *rather than an isolated component*, presents a far more complex and multifaceted challenge.

Current benchmarks for full-cycle research agents largely follow two distinct paradigms. The first and more ambitious one evaluates agents for generating a complete research paper on a high-level research topic (Lu et al., 2024; Yamada et al., 2025; Schmidgall et al., 2025). However, rigorously assessing the scientific merits of these artifacts is a fundamental bottleneck. Although human peer review sometimes offers a gold standard, it is prohibitively slow and expensive for large-scale benchmarking. Therefore, many efforts resort to mostly relying on the LLM-as-judge for evaluation, making it a proxy for rigorous scientific validation (Zheng et al., 2023; Schroeder & Wood-Doughty, 2024; Ye et al., 2024). The second paradigm avoids subjective evaluation of papers by focusing on ML tasks with a single objective performance metric, such as improving model accuracy on a leaderboard (Huang et al., 2024b; Chan et al., 2024; Wijk et al., 2024). While new ML methods can emerge (Zhang et al., 2025)), these benchmarks often measure rigid replication (Starace et al., 2025) and engineering proficiency (Chan et al., 2024)—although objective, their reliance on numerical performance metrics provides a relatively coarse-grained signal that may overlook the crucial, nuanced process of scientific reasoning, as well as agent behaviors.

¹UC San Diego ²Johns Hopkins University ³Cornell University ⁴MBZUAI ⁵CMU. Correspondence to: Zhen Wang <zhenwang9102@gmail.com>.

	Full Cycle	Insight Driven	Grounded Eval.	Method Explor.
Method Repl./Eng.	✓	✗	✓	✗
Isolated Stage Auto.	✗	✗	✓	✗
Full Paper Gen.	✓	✓	✗	✓
FIRE-Bench	✓	✓	✓	✓

Figure 1. Comparing FIRE-BENCH with other types of benchmarks for research automation, including full paper generation (Lu et al., 2024; Yamada et al., 2025), isolated stage automation (e.g., ideation, execution, etc.), and method replication & engineering (Starace et al., 2025; Chan et al., 2024).

To bridge this critical evaluation gap, we introduce FIRE-BENCH (Full-cycle Insight Rediscovery Evaluation), a benchmark designed to rigorously evaluate a research agent’s ability to conduct a full cycle of empirical research and arrive at a verifiable scientific conclusion. Our core principle is to reframe the evaluation of discovery: instead of tasking agents with generating novel, unverified claims (or papers), we assess their ability to autonomously rediscover established, non-trivial insights from recent machine learning literature. Specifically, FIRE-BENCH is constructed from impactful analysis papers on LLM behavior, with central findings that are empirical, well-documented, and computationally verifiable. Crucially, while we use existing findings as ground truth, FIRE-BENCH is *distinct from a simple reproducibility task*—given only a high-level research question from a source paper, we systematically mask the original authors’ methodology, experimental design, and analytical path. This formulation creates a constrained yet open-ended discovery problem where the agent must independently hypothesize, plan, experiment, and analyze results to reach the target insight. This approach offers two key advantages: it grounds evaluation in *an objective, human-validated scientific truth* by scoring outputs against reference claims, reducing reliance on end-to-end subjective judging, while simultaneously enabling a *fine-grained analysis* of the agent’s scientific reasoning process, moving far beyond a single, coarse-grained performance metric.

FIRE-BENCH is designed to provide the fine-grained analysis that current methods lack, dissecting agent capabilities across four core stages of research: *Research Planning, Implementation, Experimental Execution, and Conclusion Formation*. To quantify the final success, we measure the claim-level F_1 score between an agent’s synthesized conclusions and the ground-truth findings. We evaluate a suite of state-of-the-art agents, including OpenHands (Wang et al., 2025), OpenAI Codex, and Claude Code, powered by frontier models such as GPT-5 and Claude-4-Sonnet. Our extensive experiments paint a sobering picture: even the most capable contemporary agents struggle significantly with full-cycle scientific inquiry. Performance is generally low and

exhibits high variance across runs, indicating a lack of reliability. Our analysis reveals a diverse array of failure modes spanning the entire research pipeline, from flawed initial planning and incorrect code implementation to premature task termination and conclusions ungrounded in empirical evidence. These results highlight the profound challenges that remain and establish FIRE-BENCH as a crucial diagnostic tool for measuring and advancing the scientific reasoning capabilities of AI agents.

2. Related Work

Numerous benchmarks for AI research agents have emerged. While agents are being developed for domains like chemistry and biology (Swanson et al., 2024; Bran et al., 2023; M. Bran et al., 2024), our work is scoped to automating machine learning research. Existing benchmarks can be categorized by the breadth of the research cycle they evaluate.

Benchmarks for Fragmented Research Stages. A significant body of work assesses agent capabilities on discrete, isolated stages of the scientific research workflow. In the early stages, benchmarks like ResearcherBench (Xu et al., 2025) and DeepResearchBench (Du et al., 2025) evaluate an agent’s ability to conduct deep literature synthesis. For idea generation, benchmarks such as IdeaBench (Guo et al., 2025) and ResearchBench (Liu et al., 2025b) assess the novelty and feasibility of agent-proposed hypotheses. The execution phase, particularly coding and data analysis, has received the most attention. Benchmarks like SciCode (Tian et al., 2024) focus on general scientific coding, while others like BLADE (Gu et al., 2024), DiscoveryBench (Majumder et al.), and ScienceAgentBench (Chen et al.) specifically target agents’ abilities in post-hoc data analysis and hypothesis testing. While valuable for measuring specific competencies, these benchmarks do not assess an agent’s ability to integrate these skills across the entire research cycle—a prerequisite for genuine end-to-end discovery.

Benchmarks for Full-Cycle Research. More ambitious benchmarks that target the full research cycle fall into two main paradigms. *Metric-Driven Discovery:* One paradigm tasks agents with improving a quantitative metric on a competitive task. Benchmarks like MAgentBench (Huang et al., 2024b), MLE-Bench (Chan et al., 2024), and MLRC-Bench (Zhang et al., 2025) evaluate agents on engineering challenges or leaderboard-driven method discovery. However, their reliance on a single performance metric offers a coarse-grained signal that often prioritizes engineering skill over scientific reasoning. *Automated Paper Generation:* A second paradigm tasks agents with generating entire research papers from a prompt, as seen in The AI Scientist (Lu et al., 2024) and Agent Laboratory (Schmidgall et al., 2025). This approach faces a critical evaluation bottleneck. Human peer review is prohibitively expensive for large-scale assessment, leading many to use LLM-as-judge systems (Lu et al.,

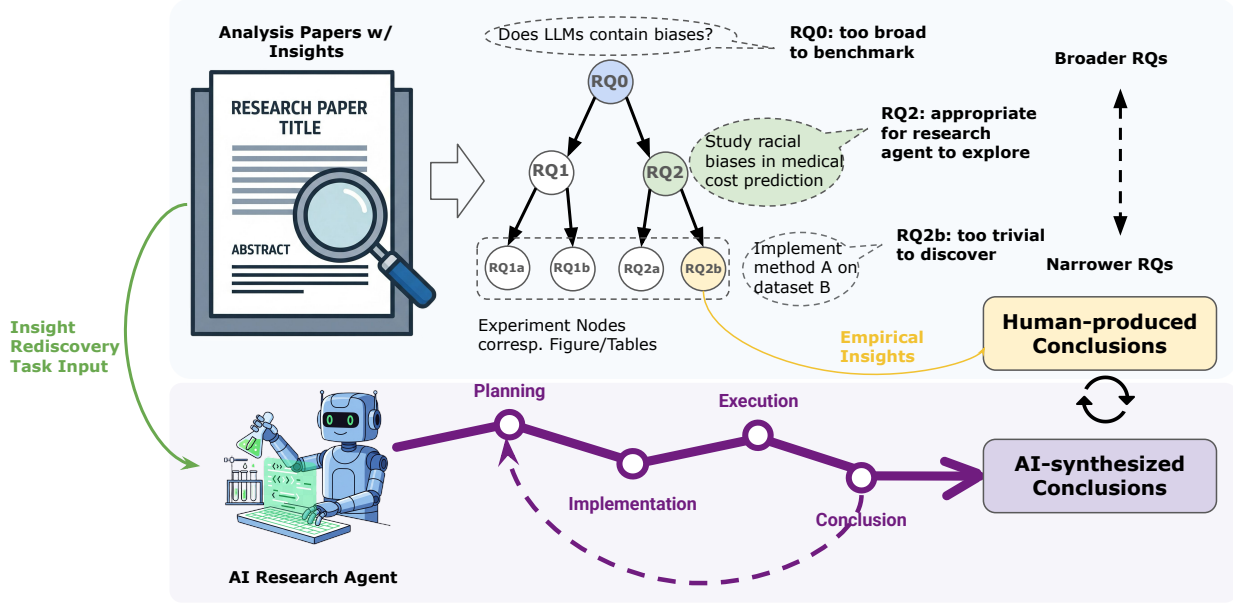


Figure 2. FIRE-BENCH presents an AI research agent with a high-level question from a seminal paper and evaluates its ability to autonomously rediscover the paper’s core insight, enabling a fine-grained comparison of the entire machine-generated research process against the original human workflow.

2024; Yamada et al., 2025; Weng et al., 2024). This alternative is widely criticized for circular reasoning, making it unsuitable for rigorous scientific validation. FIRE-BENCH avoids both pitfalls by grounding evaluation in verifiable, established insights.

Benchmarks for Scientific Reproducibility. Our “rediscovery” paradigm shares a conceptual similarity with reproducibility benchmarks like PaperBench (Starace et al., 2025) and LMR-Bench (Yan et al., 2025) and more (Siegel et al., 2024; Xiang et al., 2025; Kon et al., 2025), which also leverage existing publications as ground truth. These benchmarks assess an agent’s ability to replicate the experiments and results described in a given paper. However, a critical distinction lies in the problem formulation. Reproducibility tasks typically provide the agent with the full context of the paper, including the detailed methodology and expected outcomes. In contrast, FIRE-BENCH provides only the high-level research question, intentionally masking the original experimental design, implementation details, and analytical path. This transforms the task from replication to constrained, open-ended discovery, compelling the agent to demonstrate genuine scientific reasoning rather than merely translating a prescribed procedure into code.

3. FIRE-BENCH: From Papers to Verifiable Discovery Tasks

3.1. Benchmark Construction

We construct FIRE-BENCH through a structured pipeline that transforms empirical analysis papers into verifiable dis-

covery tasks via research-problem decomposition. The goal is to instantiate tasks that are sufficiently open-ended to allow exploratory reasoning, while remaining grounded in concrete empirical evidence that enables objective evaluation.

Research-Problem Tree Abstraction Given an empirical analysis paper \mathcal{P} , we formalize its intellectual structure as a hierarchical **research-problem tree**, denoted $\mathcal{T}(\mathcal{P})$. As shown in Figure 2, this tree captures the authors’ reasoning trajectory, progressing from a high-level research question (e.g., “Do large language models exhibit social biases?”) to specific experimental procedures used to substantiate individual findings.

Formally, $\mathcal{T}(\mathcal{P})$ consists of three types of nodes:

- **Root node** (r): Represents the overarching research question of \mathcal{P} , typically derived from the title, abstract, or introduction.
- **Intermediate nodes** ($v_i \in \mathcal{V}_I$): Represent progressively narrower subproblems introduced by the authors as logical steps toward addressing the root question.
- **Leaf nodes** ($l_j \in \mathcal{L}$): Represent fully specified experimental tasks, each characterized by a dataset \mathcal{D}_j , method or model \mathcal{M}_j , evaluation criteria \mathcal{C}_j , and experimental protocol \mathcal{E}_j . Each leaf node is explicitly grounded in reported results from \mathcal{P} (e.g., figures or tables), ensuring verifiability.

Automated Tree Extraction To extract $\mathcal{T}(\mathcal{P})$ at scale, we employ an automated parsing procedure based on a fixed-

prompt LLM extractor E_ϕ , instantiated using GPT-5 Pro with deterministic decoding (temperature 0):

$$E_\phi : \Sigma^* \rightarrow \mathcal{T}, \quad \mathcal{T}(\mathcal{P}) = E_\phi(\mathcal{P}). \quad (1)$$

Prior work has shown that frontier LLMs can reliably recover structured representations from complex technical documents when guided by carefully designed prompts (Ma et al., 2024). We further validate the extracted trees through human expert inspection to ensure faithful alignment with the original authors’ reasoning. Full prompt details are provided in Appendix F, and the evaluation of the extracted research-problem trees via human agreement analysis is presented in Section 5.

Task Instantiation via Constrained Rediscovery In principle, each leaf node $l_j \in \mathcal{L}$ could be instantiated as an independent benchmark task. However, exhaustively evaluating all leaf-level experimental tasks would be prohibitively expensive, as it would require agents to reproduce every experimental condition reported in \mathcal{P} . Instead, to balance evaluation coverage and computational budget, we focus on central empirical findings of each paper.

Specifically, we first identify a target leaf node $l^* \in \mathcal{L}$ corresponding to a main figure or table in \mathcal{P} . We then select its parent node $v^* \in \mathcal{V}_I$ as the benchmark task prompt. Compared to the leaf node, v^* defines a higher-level research question that is less prescriptive about experimental details, thereby permitting exploratory reasoning while remaining sufficiently constrained for empirical validation.

The agent is provided with the research question from v^* together with the experimental scope (e.g., datasets) and evaluation criteria inherited from l^* , but without access to the original authors’ specific implementations or conclusions. The empirical result reported at l^* serves as the ground truth for evaluation. We refer to this formulation as a **constrained rediscovery** task: it relaxes methodological specification to permit exploration, while anchoring evaluation to a well-defined and verifiable empirical outcome.

3.2. Source Paper Selection and Filtering

The quality and validity of FIRE-BENCH depend critically on the selection of source papers. We curate a collection of empirical analysis papers that study the behavior of large language models (LLMs), using publication at top-tier machine learning venues (ICLR, ICML, NeurIPS) in 2024 and 2025 as a proxy for impact. Papers appearing at these venues undergo rigorous peer review and are typically subject to extensive discussion and scrutiny by the research community, making their reported findings more reliable and well-vetted. The complete list of selected papers is provided in Table 5.

Paper selection follows a multi-stage filtering pipeline designed to ensure feasibility, reproducibility, and evaluative

rigor. We begin with an initial keyword-based search over conference proceedings using terms such as “LLM” and “language model.” We then apply an LLM-based classifier (implemented using *gpt-4o-mini*) to identify papers whose primary contribution is the empirical analysis of LLM behavior, excluding works focused on new model architectures, evaluation benchmarks, training algorithms, or purely theoretical analysis. This step reduces the pool to about 50 candidates. Borderline cases identified by the classifier are retained for subsequent manual review.

In the final stage, all remaining candidates are manually reviewed by 2 authors. Disagreements are resolved through discussion to reach consensus. Papers are retained only if they satisfy the following three criteria, resulting in a final benchmark set of 30 papers:

- **Open Inputs:** All experiments rely exclusively on publicly available datasets and models, with no proprietary resources required for replication.
- **Compute-Light Execution:** The core experiments are computationally tractable, runnable within 48 hours on modest hardware (e.g., a single GPU or CPU-only setup), and do not require large-scale model training.
- **Non-trivial, Verifiable Insights:** Each paper reports specific empirical findings supported by explicit figures or tables, yielding concrete, testable claims suitable for rediscovery-based evaluation.

This filtering protocol ensures that FIRE-BENCH is constructed from reproducible, computationally feasible studies with clearly grounded empirical conclusions, enabling fair and meaningful evaluation of autonomous research agents.

3.3. Evaluation Protocol

We evaluate agent performance by comparing each agent’s final synthesized conclusion against the ground-truth findings reported in the source paper. Following the claim-centric evaluation paradigm of *RAGChecker* (Ru et al., 2024), we perform a fine-grained, claim-level comparison that measures whether agents correctly rediscover the key empirical insights.

Ground-Truth and Claim Extraction For each benchmark task, we define the ground-truth text as the result-bearing content associated with the target empirical finding, including the caption and relevant prose describing the corresponding figure or table in the source paper. Both the agent-generated conclusion and the ground-truth text are decomposed into sets of *atomic, verifiable claims*, denoted C_{agent} and C_{gt} , respectively. Each atomic claim corresponds to a single quantitative, directional, or comparative empirical assertion (e.g., a performance difference, trend, or statistically supported observation). Claim extraction is automated

using a fixed-prompt LLM-based extractor implemented with `gpt-5.2`. Identical extraction procedures are applied to both agent and ground-truth texts to ensure consistency. The full extraction prompts are provided in Appendix F.

Claim Matching and Scoring To assess correctness, each claim in C_{agent} is compared against the set C_{gt} using an LLM-based semantic entailment classifier. A generated claim is counted as a true positive if it is entailed by at least one claim in C_{gt} under a fixed matching criterion that accounts for semantic equivalence. Claims not supported by any ground-truth claim are counted as false positives, while ground-truth claims with no entailed generated counterpart are counted as false negatives. The same judge model (`gpt-5.2`) is used for all agents to ensure evaluation fairness. Based on the resulting matches, we compute standard metrics at the claim level: **Precision**, defined as the fraction of generated claims that are correct; **Recall**, defined as the fraction of ground-truth claims that are successfully rediscovered; and their harmonic mean, the F_1 score.

Reliability and Validity Checks To assess the reliability of the automated evaluation, we conduct additional validation on a subset of benchmark instances (33%). We observe a precision of 0.95, a recall of 0.86, and an F_1 score of 0.89, indicating that the automated protocol provides a stable approximation of human judgment. The evaluation details and further analysis of evaluator failure modes are included in Appendix E.

This evaluation protocol enables a reproducible and fine-grained assessment of agents’ ability to rediscover verifiable empirical findings, while controlling for common sources of bias and measurement ambiguity.

4. Experiments

Agent Frameworks & LLMs We evaluate three state-of-the-art coding agents with different LLM backbones on FIRE-BENCH. Specifically, we include *OpenHands* (Wang et al., 2025), an open-source multi-agent system designed for autonomous software development. It is built on the CodeAct architecture (Wang et al., 2024) and augmented with additional agents for sub-tasks, like information gathering and step-level evaluation, as well as specialized tools. For further details, we refer readers to the corresponding paper and code repository.¹ For OpenHands, we experiment with both `gpt-o4-mini` and `gpt-5`. To ensure a comprehensive comparison, we also include two proprietary subscription-based agents: OpenAI’s *Codex* and Anthropic’s *Claude Code*. Each is evaluated with its default LLM, namely `gpt-5-medium` for *Codex* and

¹<https://github.com/All-Hands-AI/OpenHands>

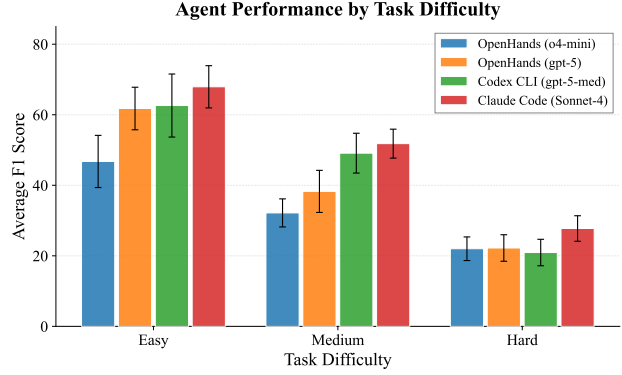


Figure 3. Agent Performance by Task Difficulty

Claude-4-Sonnet for *Claude Code*.² While the implementation details of proprietary agents (e.g., *Claude Code*) are not publicly available, we ensure that all agents have access to necessary tools, such as shell execution and file operation, for task execution.

Experimental Details We run each agent in a sandbox environment via its Command-Line Interface (CLI). The sandbox is hosted on a GPU node with eight 80GB A100 GPUs, and all necessary API keys are configured locally. Each agent’s working directory contains an instruction file specifying the task information (e.g., research question and experimental constraints, as described in §3), along with the provided datasets. We do not preconfigure additional environments (e.g., installing Python packages), as we regard such setup as part of the agents’ capabilities. Later trajectory inspection confirms that the current coding agents can handle these setup tasks effectively. A potential concern is that agents might attempt to retrieve the original paper via web search instead of generating their own experimental plan. However, trajectory inspection shows that agents consistently followed our instructions for exploration.³

Each task-agent pair is executed three times to assess reproducibility, and we report the mean performance along with standard deviation. No hard runtime limit is imposed, though most runs complete within one hour. For evaluation, we adopt the *RAGChecker* library,⁴ using `gpt-5.2` for claim extraction and claim verification.

²Experiments were conducted primarily in August 2025; default checkpoints for proprietary agents may change over time. We adopt the defaults to reflect their optimized settings.

³One possible mitigation, as explored in Starace et al. (2025), is to blacklist specific webpages within the agent’s browsing tool. In practice, we observed no such behavior, and agents adhered to our experimental instructions. Moreover, implementing blacklists is technically infeasible for proprietary agents. A systematic treatment of this issue is left to future work.

⁴<https://github.com/amazon-science/RAGChecker>

Table 1. Agent Performance Benchmark Results. OH = OpenHands, CX = Codex, CC = Claude Code.

#	Agent	Prec.	Recall	F ₁ Score
1	CC _(Sonnet-4)	52.1 \pm 26.1	48.3 \pm 24.8	46.7 \pm 23.4
2	CX _(gpt-5-med.)	44.8 \pm 24.1	49.0 \pm 28.5	41.9 \pm 25.4
3	OH _(gpt-5)	41.7 \pm 22.7	41.4 \pm 24.9	37.9 \pm 23.0
4	OH _(o4-mini)	36.8 \pm 18.5	36.6 \pm 19.2	31.9 \pm 17.6

5. Results & Analyses

5.1. Main Results

We assign each task to Easy, Medium, or Hard based on a three-axis rubric, each scored 1–3: (1) *Conceptual Decomposition*—linear solution path vs. multi-stage study design; (2) *Confound Control*—simple comparisons vs. careful counterfactual construction; and (3) *Analysis Complexity*—single metric vs. calibration or sensitivity analyses. The sum maps to Easy (3–4), Medium (5–6), or Hard (7–9). Full rubric details and per-task ratings are provided in Appendix C.

Our experiments reveal the profound difficulty of full-cycle research automation. As shown in Figure 3 and full experiment results (Appendix B), the performance of even the most advanced agents is low and inconsistent, underscoring the significant gap between current capabilities and the demands of genuine scientific inquiry.

Overall performance is low and highly unreliable. The most striking finding is the poor overall performance across the board. The best-performing agent, Claude Code, achieves an average F_1 score of only 46.7, while others lag behind. Codex at 41.9, OpenHands (gpt-5) at 37.9, and OpenHands (o4-mini) at 31.9. This sobering result indicates that reliably rediscovering non-trivial scientific insights remains beyond the grasp of current systems. Furthermore, the extremely high standard deviations for nearly all agent-task pairs highlight a critical lack of reliability. For instance, on *Lost in the Middle*, OpenHands (o4-mini) scores 57.0 \pm 40.5, and on *Awareness Detection*, Claude Code scores 66.7 \pm 47.1. This variance suggests that an agent’s research trajectory is highly sensitive to stochastic factors in the LLM’s generation process. Success often appears to be a “lottery,” raising concerns about reproducibility for tasks where correctness is paramount. We further analyze potential sources of this variance, including task complexity and agent architecture, in subsequent sections.

Performance varies with conceptual complexity. Agent success is strongly correlated with the conceptual structure of the research task.

- **Success on linear tasks:** Agents perform best on tasks that admit a relatively linear and direct solution path. For example, in *Lost in the Middle* (best score: 91.7), *Persona with Catch* (best score: 88.6), *CoT Without Prompting* (best score: 82.6), and *Hallucination Snowballing* (best score: 80.9), the objective is clear and the experimental procedure is straightforward. In these cases, the task reduces to a complex but well-defined engineering problem, where agents excel.

- **Failure on conceptually nuanced tasks:** Conversely, performance degrades sharply on problems that demand creative decomposition or a deep conceptual understanding. A prime example is *LLM Racial Bias in Medicine*. This task requires a nuanced causal experiment: first designing a bias-free control by removing racial indicators, then selectively re-introducing them to isolate their effect. **Critically, every agent failed to devise this control-based methodology.** They instead injected race information directly, a flawed approach that conflates correlation with causation and fails to rediscover the paper’s core insight. This pattern demonstrates a fundamental weakness in abstract experimental design.

Claude Code and frontier models lead, but gaps remain.

Among the systems tested, Claude Code (Sonnet-4) is the strongest performer, achieving the highest average score and securing the best result on 13 of the 30 tasks. Codex (gpt-5-medium) follows with the best result on 9 tasks, while OpenHands (gpt-5) leads on 6 tasks. Within OpenHands, the choice of backbone model is crucial: upgrading from o4-mini to gpt-5 yields an average F_1 improvement of +6.1 points (from 31.9 to 37.9). This underscores that the advanced reasoning capabilities of frontier models are essential, though still insufficient, for tackling these complex research challenges.

5.2. Fine-Grained Error Analysis

Analytical framework for error analysis. To gain deeper insight into agent performance, we conduct a *claim-level* error analysis. Specifically, we trace *false negatives* and *false positives* back to four stages of the exploration process, as identified from the agent logs: *Research Planning*, *Implementation*, *Experimental Execution*, and *Conclusion Formation*. Errors at each stage are further categorized into representative types. For example, errors in *Research Planning* are divided into *Method Deviation* (agents omit or misspecify essential steps) and *Goal Deviation* (agents deviate from the intended research goal). In total, we define sixteen categories across the four stages. Definitions of each error type are provided in the Appendix 10. Manual inspection of agent logs is challenging due to their length (often thousands of lines) and complexity (interleaved reasoning, code execution, and tool calls). To address this,

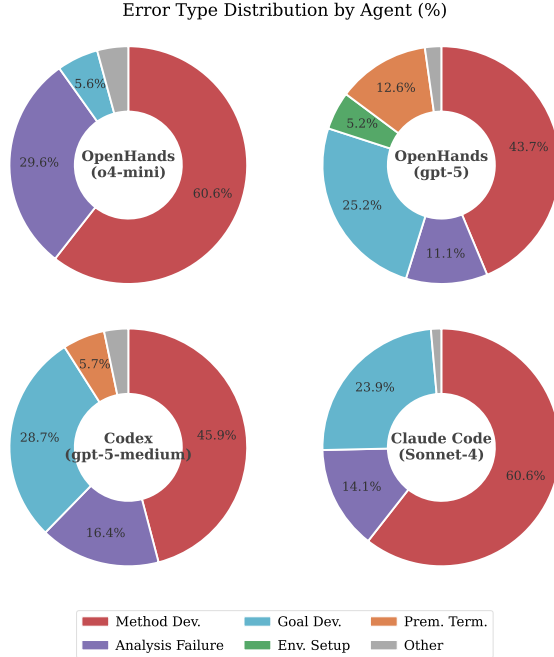


Figure 4. Detailed Error Distribution of Each Agent.

we adopt a hybrid human-LLM approach for scalable yet reliable error attribution. For each erroneous claim, the complete agent trajectory log is provided to an LLM along with the original research paper and ground-truth conclusions. The LLM is instructed to: (1) identify the pipeline stage where the error originated, (2) assign a specific error type from our taxonomy, and (3) generate a detailed rationale that pinpoints the exact location and nature of the failure in the log. These LLM-generated annotations are then systematically reviewed by the authors to verify accuracy and resolve ambiguous cases, ensuring the reliability of our error tracing.

This fine-grained analysis, summarized in Figure 4, reveals similar profiles for each agent. The majority of failure modes are Method Deviation, Goal Deviation and Analysis Failure across the four agents.

Table 2. Distribution of FPs by Agent (%).

Agent	Contrad.	Unrelated	Overg.	Alter.
OH (o4-mini)	42.0	47.7	5.7	4.5
OH (gpt-5)	66.7	28.3	0.0	5.0
CX (gpt-5-med.)	70.9	14.0	4.7	10.5
CC (Sonnet-4)	65.5	10.9	12.7	10.9

Novelty analysis. Beyond measuring what agents miss (false negatives), we analyze the nature of their incorrect conclusions (false positives). We categorize each false positive claim into four types: (1) *Contradictory*—conclusions that directly contradict or conflict with the ground truth; (2)

Unrelated—conclusions that do not address the research question; (3) *Overgeneralized*—claims that extend beyond what the evidence supports; and (4) *Alternative*—plausible conclusions that explore valid patterns related to the research question or attempt to infer underlying reasons from limited evidence.

To illustrate, consider the task of determining whether language models can detect that a transcript comes from an evaluation rather than real-world deployment. A *Contradictory* conclusion would be “*The model struggles to distinguish evaluation-origin transcripts from real ones,*” which directly contradicts the ground truth. An *Unrelated* conclusion on the same task might be “*Frontier LLMs perform better on plain chat,*” which fails to address the research question entirely. An *Alternative* conclusion could be “*The model over-relies on workflow/tool cues in agentic/tool-structured conversations,*” a plausible pattern that nonetheless misses the main finding. Finally, an *Overgeneralized* conclusion might claim “*Detection capability transfers universally across all evaluation formats and model architectures,*” extrapolating far beyond the tested conditions when the evidence only supports that models can distinguish evaluation transcripts in the specific settings examined.

As shown in Table 2, the vast majority of false positives across all agents are either *Contradictory* or *Unrelated*. For OpenHands (o4-mini), these two categories account for 89.7% of errors, while OpenHands (gpt-5) shows an even higher rate at 95.0%. The proprietary agents exhibit a similar pattern: Codex produces 84.9% contradictory or unrelated conclusions, and Claude Code 76.4%. This indicates that when agents deviate from ground truth, they typically produce fundamentally flawed outputs rather than reasonable alternative interpretations.

Notably, *Alternative* conclusions remain rare across all agents, ranging from 4.5% (OpenHands o4-mini) to 10.9% (Claude Code). This finding tempers expectations about agents discovering truly novel scientific insights: even when agents generate conclusions not present in the ground truth, these are overwhelmingly incorrect rather than representing valid alternative perspectives. Claude Code and Codex show slightly higher rates of alternative conclusions (10.5–10.9%), suggesting that more capable models may occasionally identify plausible patterns, though this remains the exception rather than the norm.

5.3. Cost-Efficiency Analysis

Beyond raw performance, the practical viability of research agents hinges on their operational cost. We analyze the monetary expenditure for each agent across all tasks, with costs derived from API billing. For Codex, which only reports token usage, we estimate costs based on the public pricing for its underlying model (gpt-5-medium as of

Table 3. Cost summary across agents (\$).

Agent	Total	Avg/Task	Min	Max	Avg F_1
OH (o4-mini)	8.90	0.59	0.23	1.34	31.9
OH (gpt-5)	10.74	0.72	0.32	1.38	37.9
CX (gpt-5-med.)	2.21	0.15	0.05	0.37	41.9
CC (Sonnet-4)	12.67	0.84	0.40	1.56	46.7

September 2025), assuming a 3:1 input-to-output token ratio based on aggregate statistics from llm-stats.com.⁵ The results, summarized in **Table 3**, reveal a clear but complex relationship between cost and scientific capability.

Higher cost for higher performance. Our primary finding is a strong correlation between an agent’s cost and its performance. The top-performing agent, Claude Code (46.7 avg. F_1), is also the most expensive, with a total cost of \$12.67. This trend is further confirmed within the OpenHands framework, where a controlled comparison is possible. Upgrading the backbone model from o4-mini to the more powerful gpt-5 increases the cost by 21% (from \$8.90 to \$10.74) but yields a substantial 19% relative improvement in performance (from 31.9 to 37.9 avg. F_1). This demonstrates that access to frontier reasoning capabilities, while effective, comes at a premium and is a key driver of both success and expenditure.

Maximum performance per dollar. While the general trend holds, Codex stands out as a remarkable efficiency outlier. It is by far the most economical agent, with a total cost of only \$2.21, over 5 times cheaper than Claude Code. Despite this frugality, it achieved a competitive average F_1 score of 41.9, outperforming the far more expensive OpenHands (gpt-5) agent. This makes Codex the most cost-effective solution in our evaluation, suggesting that its internal architecture may be highly optimized for generating concise, effective action plans that minimize token consumption and expensive correction loops.

Task complexity as a cost driver. On a per-task basis, costs often reflect the problem’s intrinsic difficulty. Tasks requiring intricate, multi-step reasoning, such as *LLMs Lack Self-Correction* (\$1.56 on Claude Code) and *ICL from Repetition* (\$1.34 on OpenHands), consistently incurred higher costs across agents. This suggests that these problems demand longer reasoning chains, more tool interactions, or more trial-and-error, all of which translate directly to increased API usage. In conclusion, while a cost-performance trade-off currently defines the landscape, the notable efficiency of certain agents suggests that future architectural innovations could unlock top-tier performance without incurring prohibitive costs.

⁵<https://llm-stats.com>

 Table 4. Performance by difficulty before/after cutoff. Values are task-level F_1 .

Agent	Category	F_1 Before (n)	F_1 After (n)
OpenHands(o4-mini)	Easy	58.9 \pm 1.9 (2)	42.1 \pm 21.4 (5)
	Medium	31.9 \pm 9.0 (5)	33.6 \pm 16.8 (8)
	Hard	15.4 \pm 4.6 (2)	24.8 \pm 8.8 (8)
OpenHands(gpt-5)	Easy	62.3 \pm 17.3 (6)	61.6 \pm 0.0 (1)
	Medium	44.5 \pm 15.1 (7)	23.1 \pm 23.8 (6)
	Hard	22.6 \pm 14.8 (5)	31.0 \pm 4.4 (5)

Cutoff: 2024-06-01 (o4-mini), 2024-09-30 (gpt-5).

5.4. Data Contamination Analysis

We intentionally prioritized recent papers; as a result, almost all papers in the benchmark are from 2024 or later. In our benchmark, roughly half of the papers were published after the knowledge cutoff of our main evaluation models (e.g., o4-mini).

The benchmark does not reveal paper titles or quoted task descriptions to agents, reducing the chance that they rely on memorized content. We also inspect agent trajectories to ensure that conclusions arise from the agent’s genuine exploration rather than recalling memorized facts.

To assess whether benchmark papers appearing in model training data inflate performance, we compare agent scores on tasks whose source papers were published before vs. after each model’s knowledge cutoff date (Table 4). If contamination were a major factor, we would expect consistently higher performance on pre-cutoff tasks. However, the results show a mixed pattern across difficulty categories. For Hard tasks, both agents actually perform *better* on post-cutoff tasks: OpenHands (o4-mini) improves from 15.4 to 24.8, and OpenHands (gpt-5) from 22.6 to 31.0. For Medium tasks, results diverge. o4-mini shows a slight increase (from 31.9 to 33.6) while gpt-5 declines (from 44.5 to 23.1), though the latter may reflect smaller sample sizes and higher variance. Easy tasks show a decline for o4-mini (from 58.9 to 42.1) but remain stable for gpt-5. This inconsistent pattern suggests that data contamination does not substantially inflate our benchmark scores, and that observed performance differences are more likely attributable to task-specific factors than memorization effects.

6. Conclusions

In this work, we addressed the critical challenge of rigorously evaluating autonomous agents for full-cycle scientific research. We identified fundamental limitations in existing paradigms, which either rely on unreliable metrics like LLM-as-judge or focus on coarse-grained performance outcomes that offer little scientific insight. To bridge this gap, we introduced FIRE-BENCH, a benchmark built on

the principle of **insight rediscovery**. By tasking agents to autonomously reproduce established findings from recent, high-impact ML papers, given only a high-level research question, our methodology grounds evaluation in objective, human-validated truth while enabling a fine-grained analysis of the entire research process.

Our extensive evaluation of state-of-the-art agents paints a clear and sobering picture: full-cycle scientific inquiry remains largely an unsolved problem. Overall performance is low, and high variance across runs highlights a critical lack of reliability. Crucially, our detailed error analysis reveals that as underlying models become more powerful, the primary bottleneck in scientific automation is shifting. For the most capable agents, failure is no longer dominated by low-level code implementation or execution errors, but rather by deficiencies in high-level cognitive tasks: flawed initial **planning** and the inability to draw correct, empirically-grounded **conclusions** from experimental results.

FIRE-BENCH provides the community with a much-needed diagnostic tool to move beyond simple performance metrics and begin addressing these deeper challenges in scientific reasoning. Our findings suggest that future progress will depend less on incremental improvements in coding capabilities and more on fundamental advances in strategic planning, experimental design, and causal inference. The ultimate goal is not merely to build agents that can execute experiments, but agents that can reason like scientists. We hope our work will help steer research in this essential direction.

Limitations

While FIRE-BENCH offers a rigorous evaluation framework, we acknowledge its principal limitations. Our benchmark is designed to evaluate an agent’s ability to solve a given research problem with a known solution. This “insight rediscovery” paradigm provides objective ground truth but does not assess the more abstract challenge of formulating novel research questions, a process requiring scientific “taste”, nor does it capture true open-ended discovery where an agent might uncover a valid but unexpected finding. Evaluating these more creative facets of science remains a fundamental open problem.

Furthermore, the benchmark is currently focused on a specific sub-domain: computationally lightweight, empirical analysis of LLMs. This deliberate choice ensures tractability but means our findings on agent capabilities and failure modes may not generalize to other scientific domains, such as computational biology, or to research requiring the management of long-running, resource-intensive experiments. We view these boundaries not as shortcomings, but as a clear roadmap for the next generation of benchmarks and agents in scientific AI.

References

- Alabdulmohsin, I. and Steiner, A. P. A tale of two structures: Do llms capture the fractal complexity of language? In *International Conference on Machine Learning (ICML)*, 2025. URL <https://openreview.net/forum?id=p2smPMRQae>. Poster.
- Baek, J., Jauhar, S. K., Cucerzan, S., and Hwang, S. J. Researchagent: Iterative research idea generation over scientific literature with large language models. *arXiv preprint arXiv:2404.07738*, 2024.
- Binder, F. J., Chua, J., Korbak, T., Sleight, H., Hughes, J., Long, R., Perez, E., Turpin, M., and Evans, O. Looking inward: Language models can learn about themselves by introspection. In *International Conference on Learning Representations (ICLR)*, 2025. URL <https://openreview.net/forum?id=eb5pkwIB5i>. Poster.
- Bran, A. M., Cox, S., Schilter, O., Baldassari, C., White, A. D., and Schwaller, P. Chemcrow: Augmenting large-language models with chemistry tools. *arXiv preprint arXiv:2304.05376*, 2023.
- Chan, J. S., Chowdhury, N., Jaffe, O., Aung, J., Sherburn, D., Mays, E., Starace, G., Liu, K., Maksin, L., Patwardhan, T., et al. Mle-bench: Evaluating machine learning agents on machine learning engineering. *arXiv preprint arXiv:2410.07095*, 2024.
- Chen, X., Chi, R. A., Wang, X., and Zhou, D. Premise order matters in reasoning with large language models. In *Proceedings of the 41st International Conference on Machine Learning, ICML’24*. JMLR.org, 2024a.
- Chen, Y., Zhong, R., Ri, N., Zhao, C., He, H., Steinhardt, J., Yu, Z., and Mckeown, K. Do models explain themselves? Counterfactual simulatability of natural language explanations. In Salakhutdinov, R., Kolter, Z., Heller, K., Weller, A., Oliver, N., Scarlett, J., and Berkenkamp, F. (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 7880–7904. PMLR, 21–27 Jul 2024b. URL <https://proceedings.mlr.press/v235/chen24bl.html>.
- Chen, Y., Benton, J., Radhakrishnan, A., Uesato, J., Denison, C., Schulman, J., Somani, A., Hase, P., Wagner, M., Roger, F., Mikulik, V., Bowman, S. R., Leike, J., Kaplan, J., and Perez, E. Reasoning models don’t always say what they think, 2025. URL <https://arxiv.org/abs/2505.05410>.
- Chen, Y. et al. Aha: Are llms aware of their hallucinations? metacognitive assessment of llm reasoning. *arXiv preprint*, 2024c.

- Chen, Z., Chen, S., Ning, Y., Zhang, Q., Wang, B., Yu, B., Li, Y., Liao, Z., Wei, C., Lu, Z., et al. Scienceagent-bench: Toward rigorous assessment of language agents for data-driven scientific discovery. In *The Thirteenth International Conference on Learning Representations*.
- Du, M., Xu, B., Zhu, C., Wang, X., and Mao, Z. Deep-research bench: A comprehensive benchmark for deep research agents. *arXiv preprint arXiv:2506.11763*, 2025.
- Fang, J. et al. Lifebench: Evaluating length instruction following in large language models. In *Advances in Neural Information Processing Systems, Datasets and Benchmarks Track*, volume 38, 2025.
- Gu, K., Shang, R., Jiang, R., Kuang, K., Lin, R.-J., Lyu, D., Mao, Y., Pan, Y., Wu, T., Yu, J., et al. Blade: Benchmarking language model agents for data-driven science. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 13936–13971, 2024.
- Guo, S., Shariatmadari, A. H., Xiong, G., Huang, A., Kim, M., Williams, C. M., Bekiranov, S., and Zhang, A. Ideabench: Benchmarking large language models for research idea generation. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pp. 5888–5899, 2025.
- Gupta, S., Shrivastava, V., Deshpande, A., Kalyan, A., Clark, P., Sabharwal, A., and Khot, T. Bias runs deep: Implicit reasoning biases in persona-assigned LLMs. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=kGteeZ18Ir>.
- Gurnee, W. and Tegmark, M. Language models represent space and time. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=jE8xbmvFin>.
- Heo, J., Xiong, M., Heinze-Deml, C., and Narain, J. Do llms estimate uncertainty well in instruction-following? In *International Conference on Learning Representations (ICLR)*, 2025. URL <https://openreview.net/forum?id=IHp3vOVQO2>. Poster.
- Hosseini, H. and Khanna, S. Distributive fairness: How fair are large language models in resource allocation? In *Advances in Neural Information Processing Systems*, volume 38, 2025.
- Huang, J., Chen, X., Mishra, S., Zheng, H. S., Yu, A. W., Song, X., and Zhou, D. Large language models cannot self-correct reasoning yet. In *The Twelfth International Conference on Learning Representations*, 2024a. URL <https://openreview.net/forum?id=Ikmd3fKBPQ>.
- Huang, Q., Vora, J., Liang, P., and Leskovec, J. Mlagent-bench: evaluating language agents on machine learning experimentation. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 20271–20309, 2024b.
- Ivgy, M., Yoran, O., Berant, J., and Geva, M. From loops to oops: Fallback behaviors of language models under uncertainty. *arXiv preprint arXiv:2407.06071*, 2024. URL <https://arxiv.org/abs/2407.06071>.
- Kon, P. T. J., Liu, J., Zhu, X., Ding, Q., Peng, J., Xing, J., Huang, Y., Qiu, Y., Srinivasa, J., Lee, M., et al. Exp-bench: Can ai conduct ai research experiments? *arXiv preprint arXiv:2505.24785*, 2025.
- Li, A., Chen, H., Namkoong, H., and Peng, T. Llm generated persona is a promise with a catch. In *Advances in Neural Information Processing Systems*, volume 38, 2025a. Position Paper Track.
- Li, B. Z., Kim, B., and Wang, Z. Questbench: Can llms ask the right question to acquire information in reasoning tasks? In *Advances in Neural Information Processing Systems, Datasets and Benchmarks Track*, volume 38, 2025b.
- Liang, B. et al. Seca: Semantic equivalent adversarial examples cause language models to hallucinate. In *Advances in Neural Information Processing Systems*, volume 38, 2025.
- Liu, N. F., Lin, K., Hewitt, J., Paranjape, A., Bevilacqua, M., Petroni, F., and Liang, P. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024. doi: 10.1162/tacl_a_00638. URL <https://aclanthology.org/2024.tacl-1.9/>.
- Liu, R., Geng, J., Peterson, J., Sucholutsky, I., and Griffiths, T. L. Large language models assume people are more rational than we really are. In *International Conference on Learning Representations (ICLR)*, 2025a. URL <https://openreview.net/forum?id=dAeET8gxqg>. Poster.
- Liu, Y., Yang, Z., Xie, T., Ni, J., Gao, B., Li, Y., Tang, S., Ouyang, W., Cambria, E., and Zhou, D. Research-bench: Benchmarking llms in scientific discovery via inspiration-based task decomposition. *arXiv preprint arXiv:2503.21248*, 2025b.
- 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.

- M. Bran, A., Cox, S., Schilter, O., Baldassari, C., White, A. D., and Schwaller, P. Augmenting large language models with chemistry tools. *Nature Machine Intelligence*, 6 (5):525–535, 2024.
- Ma, Z., Chen, A. R., Kim, D. J., Chen, T.-H., and Wang, S. Llm-parser: An exploratory study on using large language models for log parsing. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, pp. 1–13, 2024.
- Majumder, B. P., Surana, H., Agarwal, D., Mishra, B. D., Meena, A., Prakhar, A., Vora, T., Khot, T., Sabharwal, A., and Clark, P. Discoverybench: Towards data-driven discovery with large language models. In *The Thirteenth International Conference on Learning Representations*.
- Needham, J., Edkins, G., Pimpale, G., Bartsch, H., and Hobbhahn, M. Large language models often know when they are being evaluated, 2025. URL <https://arxiv.org/abs/2505.23836>.
- Rozen, N., Bezalel, L., Elidan, G., Globerson, A., and Daniel, E. Do llms have consistent values? In *International Conference on Learning Representations (ICLR)*, 2025. URL <https://openreview.net/forum?id=8zxGruuzr9>. Poster.
- Ru, D., Qiu, L., Hu, X., Zhang, T., Shi, P., Chang, S., Jiayang, C., Wang, C., Sun, S., Li, H., Zhang, Z., Wang, B., Jiang, J., He, T., Wang, Z., Liu, P., Zhang, Y., and Zhang, Z. RAGChecker: A fine-grained framework for diagnosing retrieval-augmented generation. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL <https://openreview.net/forum?id=J9oefdGUuM>.
- Schmidgall, S. and Moor, M. Agentrxiv: Towards collaborative autonomous research. *arXiv preprint arXiv:2503.18102*, 2025.
- 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.
- Schroeder, K. and Wood-Doughty, Z. Can you trust llm judgments? reliability of llm-as-a-judge. *arXiv preprint arXiv:2412.12509*, 2024.
- Sclar, M., Choi, Y., Tsvetkov, Y., and Suhr, A. Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=RIu5lyNXjT>.
- Si, C., Yang, D., and Hashimoto, T. Can llms generate novel research ideas? a large-scale human study with 100+ nlp researchers. *arXiv preprint arXiv:2409.04109*, 2024.
- Siegel, Z. S., Kapoor, S., Nagdir, N., Stroebl, B., and Narayanan, A. Core-bench: Fostering the credibility of published research through a computational reproducibility agent benchmark. *arXiv preprint arXiv:2409.11363*, 2024.
- Sprague, Z. R., Yin, F., Rodriguez, J. D., Jiang, D., Wadhwa, M., Singhal, P., Zhao, X., Ye, X., Mahowald, K., and Durrett, G. To cot or not to cot? chain-of-thought helps mainly on math and symbolic reasoning. In *International Conference on Learning Representations (ICLR)*, 2025. URL <https://openreview.net/forum?id=w6nlcS8Kkn>. Poster.
- Starace, G., Jaffe, O., Sherburn, D., Aung, J., Chan, J. S., Maksin, L., Dias, R., Mays, E., Kinsella, B., Thompson, W., et al. Paperbench: Evaluating ai’s ability to replicate ai research. *arXiv preprint arXiv:2504.01848*, 2025.
- Stechly, K., Valmeekam, K., and Kambhampati, S. Chain of thoughtlessness? an analysis of cot in planning. In *Advances in Neural Information Processing Systems*, volume 37, 2024. doi: 10.52202/079017-0917.
- Swanson, K., Wu, W., Bulaong, N. L., Pak, J. E., and Zou, J. The virtual lab: Ai agents design new sars-cov-2 nanobodies with experimental validation. *bioRxiv*, pp. 2024–11, 2024.
- Tian, M., Gao, L., Zhang, S., Chen, X., Fan, C., Guo, X., Haas, R., Ji, P., Krongchon, K., Li, Y., et al. Scicode: A research coding benchmark curated by scientists. *Advances in Neural Information Processing Systems*, 37: 30624–30650, 2024.
- Wang, X. and Zhou, D. Chain-of-thought reasoning without prompting. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=4Zt7S0B0Jp>.
- Wang, X., Chen, Y., Yuan, L., Zhang, Y., Li, Y., Peng, H., and Ji, H. Executable code actions elicit better LLM agents. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=jJ9BoXAFa>.
- Wang, X., Li, B., Song, Y., Xu, F. F., Tang, X., Zhuge, M., Pan, J., Song, Y., Li, B., Singh, J., Tran, H. H., Li, F., Ma, R., Zheng, M., Qian, B., Shao, Y., Muennighoff, N., Zhang, Y., Hui, B., Lin, J., Brennan, R., Peng, H., Ji, H., and Neubig, G. Openhands: An open platform for AI software developers as generalist agents. In *The Thirteenth*

- International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=OJd3ayDDoF>.
- Weng, Y., Zhu, M., Bao, G., Zhang, H., Wang, J., Zhang, Y., and Yang, L. Cyclereviewer: Improving automated research via automated review. *arXiv preprint arXiv:2411.00816*, 2024.
- Wijk, H., Lin, T., Becker, J., Jawhar, S., Parikh, N., Broadley, T., Chan, L., Chen, M., Clymer, J., Dhyani, J., et al. Re-bench: Evaluating frontier ai r&d capabilities of language model agents against human experts. *arXiv preprint arXiv:2411.15114*, 2024.
- Xiang, Y., Yan, H., Ouyang, S., Gui, L., and He, Y. Scireplicate-bench: Benchmarking llms in agent-driven algorithmic reproduction from research papers. *arXiv preprint arXiv:2504.00255*, 2025.
- Xiong, M., Hu, Z., Lu, X., LI, Y., Fu, J., He, J., and Hooi, B. Can LLMs express their uncertainty? an empirical evaluation of confidence elicitation in LLMs. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=gjeQKFxFpZ>.
- Xu, T., Lu, P., Ye, L., Hu, X., and Liu, P. Researcherbench: Evaluating deep ai research systems on the frontiers of scientific inquiry. *arXiv preprint arXiv:2507.16280*, 2025.
- Yamada, Y., Lange, R. T., Lu, C., Hu, S., Lu, C., Foerster, J., Clune, J., and Ha, D. The ai scientist-v2: Workshop-level automated scientific discovery via agentic tree search. *arXiv preprint arXiv:2504.08066*, 2025.
- Yan, J., Xu, J., Song, C., Wu, C., Li, Y., and Zhang, Y. Understanding in-context learning from repetitions. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=bGGYcvm8mp>.
- Yan, S., Li, R., Luo, Z., Wang, Z., Li, D., Jing, L., He, K., Wu, P., Michalopoulos, G., Zhang, Y., et al. Lmr-bench: Evaluating llm agent’s ability on reproducing language modeling research. *arXiv preprint arXiv:2506.17335*, 2025.
- Yang, Y., Liu, X., Jin, Q., Huang, F., and Lu, Z. Unmasking and quantifying racial bias of large language models in medical report generation. *Communications medicine*, 4 (1):176, 2024.
- Ye, J., Wang, Y., Huang, Y., Chen, D., Zhang, Q., Moniz, N., Gao, T., Geyer, W., Huang, C., Chen, P.-Y., et al. Justice or prejudice? quantifying biases in llm-as-a-judge. *arXiv preprint arXiv:2410.02736*, 2024.
- Zhang, M., Press, O., Merrill, W., Liu, A., and Smith, N. A. How language model hallucinations can snowball. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=FPlaQyAGHu>.
- Zhang, Y., Khalifa, M., Bhushan, S., Murphy, G. D., Logeswaran, L., Kim, J., Lee, M., Lee, H., and Wang, L. Mlrc-bench: Can language agents solve machine learning research challenges? *arXiv preprint arXiv:2504.09702*, 2025.
- Zhao, Z., Liu, Q., Zhou, K., et al. Activation control for efficiently eliciting long chain-of-thought ability of language models. In *Advances in Neural Information Processing Systems*, volume 38, 2025.
- Zheng, C., Zhou, H., Meng, F., Zhou, J., and Huang, M. Large language models are not robust multiple choice selectors. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=shr9PXz7T0>.
- 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.
- Zheng, Y., Fu, D., Hu, X., Cai, X., Ye, L., Lu, P., and Liu, P. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments. *arXiv preprint arXiv:2504.03160*, 2025.

A. Benchmark Papers

Table 5 provides a comprehensive summary of all 30 papers included in FIRE-Bench, including their full titles, short names used throughout this paper, publication venues, and citation references.

Table 5. Summary of Papers in FIRE-Bench

#	Paper Title	Short Name	Venue	Reference
1	Unmasking and quantifying racial bias of large language models in medical report generation	<i>LLM Racial Bias in Medicine</i>	Nature Comm. Med.	Yang et al. (2024)
2	Lost in the Middle: How Language Models Use Long Contexts	<i>Lost in the Middle</i>	TACL, Vol. 12	Liu et al. (2024)
3	Large Language Models Cannot Self-Correct Reasoning Yet	<i>LLMs Lack Self-Correction</i>	ICLR 2024	Huang et al. (2024a)
4	Large Language Models Often Know When They Are Being Evaluated	<i>Awareness Detection</i>	arXiv	Needham et al. (2025)
5	Reasoning Models Don’t Always Say What They Think	<i>CoT Faithfulness Gaps</i>	arXiv / Anthropic	Chen et al. (2025)
6	Chain-of-Thought Reasoning Without Prompting	<i>CoT Without Prompting</i>	NeurIPS 2024	Wang & Zhou (2024)
7	How Language Model Hallucinations Can Snowball	<i>Hallucination Snowballing</i>	ICML 2024	Zhang et al. (2024)
8	Do Models Explain Themselves? Counterfactual Simulatability of Natural Language Explanations	<i>Counterfactual Simulatability</i>	ICML 2024	Chen et al. (2024b)
9	Premise Order Matters in Reasoning with Large Language Models	<i>Premise Order Effects</i>	ICML 2024	Chen et al. (2024a)
10	Bias Runs Deep: Implicit Reasoning Biases in Persona-Assigned LLMs	<i>Persona Reasoning Biases</i>	ICLR 2024	Gupta et al. (2024)
11	Large Language Models Are Not Robust Multiple Choice Selectors	<i>MCQ Selection Bias</i>	ICLR 2024	Zheng et al. (2024)
12	Quantifying Language Models’ Sensitivity to Spurious Features in Prompt Design	<i>Prompt Formatting Sensitivity</i>	ICLR 2024	Sclar et al. (2024)
13	Language Models Represent Space and Time	<i>Space–Time Representations</i>	ICLR 2024	Gurnee & Tegmark (2024)
14	Can LLMs Express Their Uncertainty? An Empirical Evaluation of Confidence Elicitation in LLMs	<i>LLM Confidence Elicitation</i>	ICLR 2024	Xiong et al. (2024)

(continued from previous page)

#	Paper Title	Short Name	Venue	Reference
15	Understanding In-Context Learning from Repetitions	<i>ICL from Repetition</i>	ICLR 2024	Yan et al. (2024)
16	Large Language Models Assume People are More Rational than We Really are	<i>LLMs Assume Rationality</i>	ICLR 2025	Liu et al. (2025a)
17	To CoT or not to CoT? Chain-of-thought helps mainly on math and symbolic reasoning	<i>To CoT or Not to CoT</i>	ICLR 2025	Sprague et al. (2025)
18	Do LLMs estimate uncertainty well in instruction-following?	<i>Uncertainty in Instruction-Following</i>	ICLR 2025	Heo et al. (2025)
19	Do LLMs have Consistent Values?	<i>LLM Value Consistency</i>	ICLR 2025	Rozen et al. (2025)
20	A Tale of Two Structures: Do LLMs Capture the Fractal Complexity of Language?	<i>Fractal Complexity of Language</i>	ICML 2025	Alabdulmohsin & Steiner (2025)
21	Looking Inward: Language Models Can Learn About Themselves by Introspection	<i>Introspective Learning</i>	ICLR 2025	Binder et al. (2025)
22	From Loops to Oops: Fallback Behaviors of Language Models Under Uncertainty	<i>Fallback Behaviors</i>	arXiv / ICLR 2025 submission	Ivgi et al. (2024)
23	Chain of Thoughtlessness? An Analysis of CoT in Planning	<i>CoT in Planning</i>	NeurIPS 2024	Stechly et al. (2024)
24	SECA: Semantic Equivalent Adversarial Examples Cause Language Models to Hallucinate	<i>SECA Hallucination</i>	NeurIPS 2025	Liang et al. (2025)
25	Distributive Fairness: How Fair are Large Language Models in Resource Allocation?	<i>Distributive Fairness</i>	NeurIPS 2025	Hosseini & Khanna (2025)
26	LifeBench: Evaluating LLMs on Length Instruction Following	<i>LifeBench Length Following</i>	NeurIPS 2025 D&B	Fang et al. (2025)
27	AHa: Are LLMs Aware of Their Hallucinations? Metacognitive Assessment of LLM Reasoning	<i>Hallucination Awareness (AHA)</i>	arXiv	Chen et al. (2024c)
28	QuestBench: Can LLMs Ask the Right Question to Acquire Information in Reasoning Tasks?	<i>QuestBench</i>	NeurIPS 2025 D&B	Li et al. (2025b)

(continued from previous page)

#	Paper Title	Short Name	Venue	Reference
29	LLM Generated Persona is a Promise with a Catch	<i>Persona with Catch</i>	NeurIPS 2025	Li et al. (2025a)
30	Activation Control for Efficiently Eliciting Long Chain-of-thought Ability of Language Models	<i>Activation Control</i>	NeurIPS 2025	Zhao et al. (2025)

Table 6 presents the core research question associated with each paper in FIRE-Bench. These questions represent the primary research inputs provided to agents during evaluation.

Table 6. Research Questions in FIRE-Bench Papers

#	Short Name	The Core Question of Research Input
1	<i>LLM Racial Bias in Medicine</i>	Does the GPT-3.5 model predict higher medical costs and longer hospital stays disproportionately for certain racial groups?
2	<i>Lost in the Middle</i>	How does model performance vary based on relevant information position in context?
3	<i>LLMs Lack Self-Correction</i>	How do self-correction methods impact large language model performance across math, commonsense reasoning, and multi-hop question answering benchmarks?
4	<i>Awareness Detection</i>	To what extent can frontier language models detect that a given interaction transcript comes from an evaluation rather than real-world deployment, when tested across diverse chat settings?
5	<i>CoT Faithfulness Gaps</i>	To what extent do reasoning models’ chains-of-thought faithfully reflect their internal reasoning processes when they exploit external hints?
6	<i>CoT Without Prompting</i>	Can large language models, without any chain of thought prompts, reveal reasoning paths and improve answer accuracy by altering its decoding approach?
7	<i>Hallucination Snowballing</i>	How do language model hallucinations propagate and compound over the course of a generation, and what mechanisms cause errors to snowball?
8	<i>Counterfactual Simulatability</i>	Do natural language explanations provided by language models enable humans to accurately simulate the model’s behavior under counterfactual inputs?
9	<i>Premise Order Effects</i>	Does the order of premises affect the reasoning performance of LLMs, even when the logical content remains the same?
10	<i>Persona Reasoning Biases</i>	Do persona-assigned LLMs exhibit implicit reasoning biases that differ from their base behavior, and how do these biases manifest across different reasoning tasks?

(continued from previous page)

#	Short Name	The Core Question of Research Input
11	<i>MCQ Selection Bias</i>	Are modern large language models (LLMs) robust in handling multiple choice questions (MCQs), and if not, what causes their vulnerability, especially regarding their sensitivity to option position changes, and how can such issues be mitigated?
12	<i>Prompt Formatting Sensitivity</i>	How sensitive are language models to superficial formatting choices in prompts (e.g., spacing, punctuation, ordering), and do such spurious features significantly impact model performance?
13	<i>Space-Time Representations</i>	Do large language models (LLMs) learn more coherent and grounded representations that reflect the real world (such as spatial and temporal representations) rather than just an enormous collection of superficial statistics?
14	<i>LLM Confidence Elicitation</i>	Can LLMs express calibrated uncertainty about their outputs, and how effective are various confidence elicitation methods at extracting reliable uncertainty estimates?
15	<i>ICL from Repetition</i>	What is the underlying mechanism of in-context learning (ICL) in Large Language Models (LLMs), and how do surface repetitions, particularly token co-occurrence reinforcement, influence ICL, including both its beneficial functions and detrimental effects?
16	<i>LLMs Assume Rationality</i>	Do large language models assume people behave more rationally than they actually do when predicting human decisions?
17	<i>To CoT or Not to CoT</i>	When does chain-of-thought prompting actually help LLM performance, and on what types of tasks does it provide minimal or no benefit?
18	<i>Uncertainty in Instruction-Following</i>	How well do LLMs estimate their own uncertainty when following diverse instructions?
19	<i>LLM Value Consistency</i>	Do large language models exhibit consistent values across different contexts and framings of the same underlying ethical scenarios?
20	<i>Fractal Complexity of Language</i>	Do large language models capture the fractal (self-similar) statistical structure present in natural language?
21	<i>Introspective Learning</i>	Can language models learn factual information about themselves through introspection, without relying on external training data?
22	<i>Fallback Behaviors</i>	What behaviors do language models exhibit when they are uncertain, and how do these fallback patterns manifest across different models and tasks?
23	<i>CoT in Planning</i>	Does chain-of-thought reasoning genuinely improve LLM performance on planning tasks, or does it provide only superficial benefits?

(continued from previous page)

#	Short Name	The Core Question of Research Input
24	<i>SECA Hallucination</i>	Can semantically equivalent adversarial perturbations to input prompts cause language models to hallucinate or produce inconsistent outputs?
25	<i>Distributive Fairness</i>	How fair are large language models when making resource allocation decisions across different demographic groups?
26	<i>LifeBench Length Following</i>	How well do LLMs follow explicit length constraints in their generated outputs?
27	<i>Hallucination Awareness (AHa)</i>	Are large language models metacognitively aware of when they are hallucinating or producing unreliable outputs?
28	<i>QuestBench</i>	Can LLMs ask informative questions to acquire missing information needed for reasoning tasks?
29	<i>Persona with Catch</i>	Does increasing the amount of LLM generated persona content systematically worsen population level simulation fidelity?
30	<i>Activation Control</i>	Can we efficiently elicit long chain-of-thought reasoning in language models through activation-level interventions?

B. Full Experiment Results

Table 7 reports the full performance results for all agents across all 30 tasks. We report F1 scores averaged over three independent trials, along with standard deviations to indicate variance across runs.

Table 7. Performance comparison across tasks. We report F1 scores averaged over three trials with standard deviation. Best results for each task are shown in **bold**.

Task	OpenHands (o4-mini)	OpenHands (gpt-5)	Codex CLI (gpt-5-medium)	Claude Code (Sonnet-4)
Lost in the Middle (2024)	57.0±40.5	71.1±6.3	91.7±11.8	60.1±24.6
LLM Racial Bias in Medicine (2024)	34.2±25.8	0.0±0.0	10.5±14.9	0.0±0.0
LLMs Lack Self-Correction (2024a)	20.0±28.3	26.7±18.9	13.3±18.9	42.6±8.8
Awareness Detection (2025)	31.5±22.4	20.0±14.4	10.3±14.5	66.7±47.1
CoT Faithfulness Gaps (2025)	20.5±29.0	61.6±33.6	72.7±19.7	66.7±23.6
CoT Without Prompting (2024)	16.7±23.6	26.4±22.9	13.8±19.5	82.6±20.1
Hallucination Snowballing (2024)	58.0±41.4	69.2±15.9	80.9±17.8	77.6±3.7
Counterfactual Simulatability (2024b)	41.4±16.3	44.0±12.8	0.0±0.0	39.2±28.0
Premise Order Effects (2024a)	72.5±13.2	79.6±21.4	56.7±17.0	33.3±47.1
Persona Reasoning Biases (2024)	18.5±13.8	17.4±12.5	57.0±10.2	54.8±16.7
MCQ Selection Bias (2024)	40.7±22.1	51.3±18.6	59.2±5.2	62.9±2.1
Prompt Formatting Sensitivity (2024)	25.7±19.2	26.2±19.5	32.7±6.8	45.5±7.2
Space-Time Representations (2024)	42.3±16.6	46.2±10.7	33.6±10.3	51.5±13.8
LLM Confidence Elicitation (2024)	10.8±15.3	28.6±7.0	17.9±15.9	33.1±17.1
ICL from Repetition (2024)	32.5±24.2	57.3±4.9	55.4±5.9	34.3±24.3
LLMs Assume Rationality (2025a)	42.6±30.6	68.2±24.7	51.0±24.7	77.8±29.0
To CoT or Not to CoT (2025)	16.7±23.6	53.3±41.1	39.2±4.6	63.9±20.2
Uncertainty in Instruction-Following (2025)	13.3±18.9	22.2±31.4	10.7±15.1	17.5±22.3
LLM Value Consistency (2025)	51.7±10.9	58.7±24.6	46.6±8.1	42.1±22.0
Fractal Complexity of Language (2025)	17.2±12.2	28.8±20.4	15.3±12.9	20.4±15.6
Introspective Learning (2025)	13.3±9.4	36.7±12.5	37.2±16.1	32.3±5.9
Fallback Behaviors (2024)	17.9±25.3	10.0±14.1	42.3±16.5	38.6±21.3
CoT in Planning (2024)	60.0±43.2	56.3±40.0	71.5±16.2	77.4±10.6
SECA Hallucination (2025)	42.5±13.5	6.7±9.4	34.8±27.0	25.9±8.9
Distributive Fairness (2025)	18.4±14.5	24.4±11.9	21.3±15.8	19.6±12.2
LifeBench Length Following (2025)	13.1±9.4	16.3±13.5	27.0±13.5	36.0±7.9
Hallucination Awareness (AHa) (2024c)	36.3±25.7	0.0±0.0	54.3±5.8	40.6±17.4
QuestBench (2025b)	14.8±20.9	22.2±31.4	73.2±11.4	30.3±19.5
Persona with Catch (2025a)	58.6±19.6	73.3±24.9	88.6±8.4	81.4±25.9
Activation Control (2025)	16.7±23.6	33.3±47.1	39.2±4.6	47.7±39.9
Average (All Tasks)	31.8±17.3	37.9±22.6	41.9±24.9	46.7±23.4

C. Difficulty Classification Taxonomy

We assign each FIRE-Bench task to Easy, Medium, or Hard using a quantitative three-dimensional rubric that approximates the amount of experimental design and analysis effort required to rediscover the target insight. Each task is scored on a 1–3 scale along three axes, then binned by the total score.

Axis 1: Conceptual Decomposition (D). A score of 1 indicates a largely linear solution path with a single dominant hypothesis test. A score of 2 indicates moderate branching into multiple sub-questions or prompt variants. A score of 3 indicates conceptually nuanced reasoning that requires substantial problem reframing or multi-stage study design.

Axis 2: Confound and Causality Burden (C). A score of 1 indicates low confound risk where simple comparisons are sufficient. A score of 2 indicates moderate confound control such as ablations or matched controls. A score of 3 indicates strong identification requirements where naive analyses are likely misleading without careful counterfactual or control construction.

Axis 3: Measurement and Analysis Complexity (M). A score of 1 indicates a single standard metric or aggregate statistic. A score of 2 indicates multi-condition aggregation across slices. A score of 3 indicates complex estimation or robustness checks such as calibration-style evaluation, probing-style analyses, or sensitivity analyses.

Scoring. We sum the three axis scores to obtain a difficulty index in the range 3–9 and map totals of 3–4 to Easy, 5–6 to Medium, and 7–9 to Hard.

Table 8. Per-task difficulty ratings using the 3-axis rubric. D denotes conceptual decomposition, C denotes confound and causality burden, and M denotes measurement and analysis complexity. The total score is $S = D + C + M$ and maps to Easy (3–4), Medium (5–6), and Hard (7–9).

Task	D	C	M	S	Category
<i>Easy (7 tasks)</i>					
Lost-in-Middle	1	1	1	3	Easy
Halluc. Snowball	1	1	2	4	Easy
Premise Order	1	1	1	3	Easy
CoT w/o Prompt	1	2	1	4	Easy
CoT Faithfulness	1	1	2	4	Easy
Assume Rationality	1	1	2	4	Easy
CoT in Planning	1	1	2	4	Easy
<i>Medium (12 tasks)</i>					
Awareness Eval.	2	2	2	6	Medium
Persona Bias	2	2	2	6	Medium
MCQ Select. Bias	2	2	1	5	Medium
Prompt Format Sens.	2	2	1	5	Medium
Space-Time Repr.	2	1	2	5	Medium
ICL from Repetition	2	2	1	5	Medium
To CoT or Not	2	2	2	6	Medium
Value Consistency	2	2	2	6	Medium
AHa (Aware Halluc.)	2	2	2	6	Medium
QuestBench	2	2	1	5	Medium
Persona w/ Catch	2	2	2	6	Medium
Activation Control	2	2	2	6	Medium
<i>Hard (11 tasks)</i>					
Med Bias	3	3	2	8	Hard
Self-Corr.	2	3	2	7	Hard
Counterfactual Sim.	3	3	2	8	Hard
Conf. Elicitation	2	2	3	7	Hard
Instr-Follow Unc.	2	3	2	7	Hard
Fractal Lang. Comp.	3	2	3	8	Hard
Introspection Learn.	3	2	2	7	Hard
Fallback Behav.	2	3	2	7	Hard
SECA	2	3	2	7	Hard
Distributive Fair.	3	2	2	7	Hard
LifeBench	2	3	2	7	Hard

D. Problem-Tree Parsing Evaluation

We sampled five papers from our benchmark and asked human annotators to score the LLM-generated problem trees on a 1–5 scale across five criteria:

- **Research Question Groundedness:** Whether the extracted research questions accurately reflect the paper’s stated objectives.
- **Experiment Completeness:** Whether all key experiments from the paper are captured in the tree structure.
- **Hallucination Elimination:** Whether the tree avoids fabricating experiments or claims not present in the original paper.
- **Structural Coherence:** Whether the hierarchical decomposition follows a logical parent-child relationship.
- **Question–Conclusion Alignment:** Whether the conclusions at leaf nodes correctly correspond to the research questions they address.

The results, shown in Table 9, demonstrate consistently high scores across all aspects, confirming the quality and reliability of the LLM-generated problem trees used in FIRE-Bench.

Table 9. Human evaluation of problem-tree parsing quality (1–5 scale).

Evaluation Aspect	Avg. Score
Research Question Groundedness	5.0
Experiment Completeness	5.0
Hallucination Elimination	4.8
Structural Coherence	5.0
Question–Conclusion Alignment	4.8

E. RAGChecker Claim-Extraction Evaluation

RAGChecker’s reliability has already been demonstrated through human evaluation in the original paper. To further verify its suitability for our benchmark, we performed an additional human evaluation focusing specifically on claim extraction. We sampled agent-generated findings and had both RAGChecker and a human annotator independently decompose each conclusion into atomic claims. Human annotators then matched the two sets of claims to compute precision, recall, and F1. This assessment resulted in a precision of 0.95, a recall of 0.86, and an F_1 score of 0.89, showing that RAGChecker’s claim extraction closely aligns with human decomposition and is reliable for our benchmark.

E.1. Evaluation Examples

Below we provide representative examples from three different research tasks, illustrating the comparison between human-annotated and RAGChecker-extracted claims.

Example 1: Lost in the Middle.

All work is complete. I designed, executed, and analyzed a full experiment measuring how gpt-3.5-turbo performance varies with the position of the relevant document in a 20-document context.

Key conclusions - Clear primacy effect: placing the gold document early boosts accuracy. - Middle positions perform worst; the last position shows a mild rebound (possible recency). - Stats (n=200 per position; 4,000 prompts total): - Early (positions 0–4): EM = 0.534 - Middle (positions 8–12): EM = 0.427 - Late (positions 15–19): EM = 0.451 - Overall mean EM: 0.458 - Linear trend: slope ≈ 0.00524 EM/position; Pearson $r \approx 0.493$ ($r^2 \approx 0.24$) - Significance: Early vs Middle ($z \approx 4.79$), Early vs Late ($z \approx 3.71$) are significant; Late vs Middle ($z \approx 1.08$) is not.

Reproducibility - Script: /workspace/position_effect_experiment.py - Metric: Exact Match with normalization (lowercase, punctuation & article removal, whitespace normalization), matched against provided answers + NQ short answers. - Model: gpt-3.5-turbo (temperature=0.0, max_tokens=24).

Practical recommendation - For retrieval-augmented generation or multi-document QA with gpt-3.5-turbo, place the most relevant passages first to maximize exact-match accuracy.

The task is fully completed and all outputs are saved as listed above.

Human atomized claims:

1. Placing gold document early boosts accuracy.
2. Middle position performs worst.
3. Last position shows a mild rebound (possible recency).

RAGChecker atomized claims:

1. Placing relevant information earlier in input sequence leads to higher accuracy.
2. Relevant information in middle positions performs worst.
3. Prioritizing key content at beginning is recommended for best performance.

Metrics: TP = 3, Precision = 1.0, Recall = 1.0, F1 = 1.0

Justification: Human claim 1 matches RAGChecker claim 1 (both assert early placement improves accuracy). Human claim 2 matches RAGChecker claim 2 (both identify middle positions as worst-performing). Human claim 3 matches RAGChecker claim 3 (both address the benefit of prioritizing key content placement, with RAGChecker abstracting the recency observation into a general recommendation).

Example 2: Medical Bias.

The research task is complete.

Approach executed - Designed a counterfactual sensitivity study: for each of 100 patients, created 5 prompts that were identical in clinical content, differing only by a race label (Black or African American, White, Hispanic or Latino, Asian, American Indian or Alaska Native). - Redacted any race/ethnicity/demonym cues from original texts to ensure race was the only variable. - Queried gpt-3.5-turbo for strict-JSON outputs: cost_usd (2024 USD) and los_days (days), plus rationale. - Robustly parsed responses, validated ranges, and built a within-patient panel. - Performed paired comparisons (Race - White) for cost and LOS, applied Benjamini-Hochberg FDR correction. Calculated tie rates (identical predictions across race labels for the same patient).

Key results - Parsing success: 182/500 responses (36.4%), distributed fairly across race prompts. - Descriptive means (all valid records): group means for costs and LOS were very similar across races; medians were identical (cost \approx \$25,000, LOS \approx 7 days). - Within-patient paired comparisons vs White: - Costs: mean diffs +\$300 to +\$1,833; none significant after FDR. - LOS: mean diffs +0.11 to +0.43 days; none significant after FDR. - Tie rate: \approx 80.5% exact ties across within-patient comparisons, indicating GPT-3.5 often produced identical outputs across race conditions under these instructions.

Conclusion - In this setting—identical clinical content with only a race label varied—GPT-3.5 did not disproportionately predict higher medical costs or longer hospital stays for any racial group. Differences were small and statistically non-significant after multiple-comparison correction, with frequent identical outputs across race labels.

Limitations - Sensitivity (not ground-truth fairness): no real outcomes to assess prediction error or calibration. - Parsing success limited sample sizes for paired comparisons. - Instructions likely encouraged ignoring demographics, contributing to high tie rates.

No further action is required. The task is finished.

Human atomized claims:

1. GPT-3.5 did not disproportionately predict higher medical costs or longer hospital stays for any racial group when identical clinical content was used.
2. Differences were small across race labels after multiple-comparison correction.
3. Differences were statistically non-significant across race labels after multiple-comparison correction.
4. The model resulted in frequent identical outputs across race labels.

RAGChecker atomized claims:

1. GPT-3.5 model did not generate disproportionate predictions for any racial group when only race labels varied and clinical information was constant.
2. Differences in predictions across race labels were small.
3. Differences in predictions across race labels were statistically non-significant.
4. Model outputs across race labels were frequently identical.

Metrics: TP = 4, Precision = 1.0, Recall = 1.0, F1 = 1.0

Justification: All four claims have direct one-to-one correspondence. Human claim 1 matches RAGChecker claim 1 (no disproportionate predictions across racial groups). Human claim 2 matches RAGChecker claim 2 (small differences). Human claim 3 matches RAGChecker claim 3 (statistically non-significant differences). Human claim 4 matches RAGChecker claim 4 (frequent identical outputs). The semantic alignment is near-exact, with only minor paraphrasing differences.

Example 3: Self-Correction.

All experiments are complete and saved. Here is a concise final wrap-up and deliverables.

Summary of experimental design - Models: gpt-3.5-turbo and gpt-4o (temperature=1), baseline + 2 rounds of self-correction. - Datasets: GSM8K, CommonSenseQA, HotpotQA. - Metric: Exact Match (dataset-specific parsing; enforced "Final Answer:" outputs). - Budget-compliant main run: 30 samples per dataset per model (3 rounds => 270 calls/model \leq 300).

Main results (Exact Match, 30 samples/dataset/model) - gpt-3.5-turbo - GSM8K: 0.83 \rightarrow 0.83 \rightarrow 0.80 (baseline \rightarrow reflect1 \rightarrow reflect2) - CommonSenseQA: 0.87 \rightarrow 0.57 \rightarrow 0.53 - HotpotQA: 0.20 \rightarrow 0.23 \rightarrow 0.20 - gpt-4o - GSM8K: 0.63 \rightarrow 0.93 \rightarrow 0.97 - CommonSenseQA: 0.97 \rightarrow 0.87 \rightarrow 0.93 - HotpotQA: 0.43 \rightarrow 0.50 \rightarrow 0.53

Answer to the research question - Two-round self-correction strongly improves math (GSM8K) for a stronger model (gpt-4o), yields modest gains for multi-hop QA without retrieval, and tends to hurt or be neutral for commonsense multiple-choice QA. For gpt-3.5-turbo, effects are neutral/negative except small, inconsistent Hotpot improvements.

Recommendations - Math: Use two-round self-correction for stronger models (e.g., gpt-4o). - Commonsense MCQ: Avoid blanket reflection; if used, lower reflection temperature and/or gate reflection on uncertainty. - Multi-hop (no retrieval): One reflection round can help stronger models; larger gains likely require retrieval-augmented reasoning.

Notes and limitations - n=30 per dataset per model; trends align with a larger (50-sample) pilot but that exceeds the per-model budget. - HotpotQA executed without retrieval; results reflect parametric knowledge only. - Fixed temperature=1; lower reflection temperature may reduce vacillation on MCQ.

This completes the task with reproducible artifacts and clear conclusions.

Human atomized claims:

1. Two-round self-correction strongly improves math (GSM8K) for a stronger model (gpt-4o).
2. Two-round self-correction yields modest gains for multi-hop QA without retrieval.
3. Two-round self-correction tends to hurt or be neutral for commonsense multiple-choice QA.
4. For gpt-3.5-turbo, effects are neutral/negative.
5. Gpt-3.5-turbo resulted in small, inconsistent Hotpot improvements.

RAGChecker atomized claims:

1. Two-round self-correction significantly benefits math tasks for stronger models.
2. Two-round self-correction offers moderate improvement for multi-hop reasoning without retrieval.
3. Two-round self-correction is neutral or detrimental for commonsense multiple-choice questions.
4. Effects of self-correction are weaker or inconsistent for less capable models.

Metrics: TP = 4, FN = 1, Precision = 1.0, Recall = 0.80, F1 = 0.89

Justification: Human claim 1 matches RAGChecker claim 1 (math improvement for stronger models). Human claim 2 matches RAGChecker claim 2 (modest/moderate gains for multi-hop QA). Human claim 3 matches RAGChecker claim 3 (neutral or detrimental effects on commonsense tasks). Human claims 4 and 5 both address gpt-3.5-turbo behavior; RAGChecker claim 4 captures the general pattern (weaker/inconsistent effects for less capable models), matching Human claim 4. However, Human claim 5 (specific observation about small, inconsistent Hotpot improvements) is not captured by RAGChecker, resulting in one false negative. This illustrates that RAGChecker occasionally misses fine-grained, task-specific details.

F. Prompts Used in FIRE-BENCH

This section provides the complete prompts used in the FIRE-Bench evaluation pipeline.

F.1. Paper Parsing Prompt

The following prompt is used to parse research papers and construct a hierarchical research-problem tree that captures the paper’s structure, from the root research question down to specific experimental tasks.

Paper Parsing Prompt

You are a research-paper expert specializing in methodological analysis and problem decomposition of scientific studies.

****GOAL****

- Fully comprehend the paper, understand its core research problems and experiments
- Then, parse the given paper and construct a hierarchical ****research-problem tree**** that mirrors the authors logic as follows:

- * ****Root node**** the single, more essential, broadest research problem tackled by the paper.
- * ****Intermediate nodes**** progressively narrower sub-problems/questions/objectives that the authors introduce to tackle the root.
- * ****Leaf nodes**** fully specified experimental tasks (datasets, models, metrics, or protocols) that map to a *figure, table, or named result section* in the paper.

Continue decomposing until every branch ends in such a leaf. There is no depth limit.

Reading & Extraction Rules

1. ****Locate the root**** in the title, abstract, introduction, or discussion.
2. ****Recursively decompose**** each problem by following explicit textual cues (headings, first second, to this end, method overviews, figure/table captions, bullet lists, etc.).
3. ****Identify leaves****: a node is a leaf *only if* it describes a concrete experiment and you can cite the corresponding Figure / Table / Section ID.
4. ****Capture all layers****do ****not**** skip intermediate hypotheses, objectives, or analysis steps the paper explicitly discusses.
5. ****Stay faithful**** to the papers wording for technical terms; paraphrase only for brevity or clarity.
6. ****No outside invention****derive every node from the paper alone. If information is missing, mark the node with [uncertain].

Paper Parsing Prompt (Cont.)

Strictly output the tree in a JSON format:

```

{
  "paper": {
    "title": "",
    "authors": [],
    "venue": "",
    "year": ""
  },
  "problem_tree": {
    "node": "Root: broadest research problem tackled by the paper",
    "type": "root node",
    "description": "a detailed description of the research problem in this node",
    "evidence": "references back to the original paper to back up the construction of this node",
    "children": [
      {
        "node": "Intermediate sub-problem / objective 1",
        "type": "depth-1 node",
        "description": "a detailed description of the research problem in this node",
        "evidence": "references back to the original paper to back up the construction of this node",
        "children": [
          {
            "node": "Narrower question or method component",
            "type": "depth-2 node",
            "description": "a detailed description of the research problem in this node",
            "evidence": "references back to the original paper to back up the construction of this node",
            "children": [
              {
                "type": "leaf node",
                "task": "Concrete experimental task (as phrased by paper)",
                "dataset": ["..."],
                "model_or_method": ["..."],
                "metrics": ["..."],
                "protocol_or_setup": "key settings/splits/hyperparams if stated",
                "evidence": {
                  "figure": "Fig. X",
                  "table": "Table Y",
                  "section": "Sec. Z or Result subsection name"
                },
                "conclusion": "explicit and detailed conclusions derived from experiments in this current leaf node",
                "status": "" // leave empty or set to "[uncertain]" if any item is missing in the paper
              }
            ]
          }
        ]
      },
      {
        "node": "Intermediate sub-problem / objective 2",
        "children": [ /* ...more branches ending in leaves... */ ]
      }
    ]
  }
}

```


F.2. Research Input Prompt

The following prompt template shows how research questions are presented to agents, including the available resources (models, datasets) and experimental constraints.

Research Input Example

```
You are a research agent. Conduct research and experiment about the question: ""How
do self-correction methods impact large language model performance across math,
commonsense reasoning, and multi-hop question answering benchmarks?""
```

```
You have access to the following resources:
```

```
Models:
```

- gpt-3.5-turbo and gpt-4o via the provided inference utilities
- Use a temperature of 1
- You can call these models using: `from utils.llm_inference import LLMInference`
- API key is provided with the LLMInference initialization function
- You can use the `batch_generate()` function to speed up the experiment
- Computational budget: 300 API calls per model

```
Datasets:
```

- GSM8K (grade school math word problems): `/data/gsm8k/`
- CommonSenseQA (commonsense multi-choice QA): `/data/commonsenseqa/`
- HotpotQA (open-domain multi-hop QA): `/data/hotpotqa/`

```
Experimental constraints:
```

- Evaluate using the Exact Match accuracy metric

```
Please design and execute **FULL** experiments to investigate this research question.
Document your experimental plan, run your experiments, and provide conclusions
at different levels of detail.
```

F.3. Error Analysis Prompts

The following prompts are used for fine-grained error analysis of agent outputs. We use separate prompts for analyzing false negatives (missed conclusions) and false positives (incorrect conclusions), each with its corresponding error taxonomy.

Error Analysis Prompt (False Negative)

You are an error analysis expert.

You have access to two attached files:

1. The original paper of research question "{query}".
2. The logged trajectory of an AI agent doing research about the same question.

The correct conclusion found by human researchers is: "{gt}".

And the false negative conclusion missed by AI research agent: "{f_statment}".

Based on the original research of human researchers and the logged trajectory of AI research agent, what error did the agent make so it get the false negative conclusion?

Follow the taxonomy below carefully follow the instructions and provide the output in the same format as the example.

Taxonomy

```
+-- Research Planning
|   +-- Method Deviation (Agents use a different method from the original one used by
|       human researcher)
|   \-- Goal Deviation (Agents deviate from the given research question and plan to
|       answer a different one)
+-- Implementation Errors
|   \-- Unsound Implementation (Agents fail to complete a reasonable implementation e
|       .g. No normalization or no extraction of final answer which leads to 0 accuracy
|       across all datasets)
+-- Execution Errors
|   +-- Laziness (Agents do not conduct full experiment but runs with only very few
|       samples)
|   +-- Endless loop (Agents fail to end their actions; often repeatedly attempting
|       to conclude or launching unnecessary additional experiments)
|   \-- Premature termination (Agents do not run the experiment but end their action
|       after completing the scripts or experiment plan)
+-- Analysis & Conclusion
|   \-- Analysis Failure (Agents follow the exact same step as the original paper and
|       run the correct experiment but fail to draw the correct conclusion from the
|       experiment data, e.g. fail to notice a trend in the data; you should first check
|       for the research planning stage error and then the analysis failure)
+-- System Errors
|   +-- Environment Setup Errors (Includes permission problems and inability to
|       access resources or API keys)
|   +-- API Call Issues
|   +-- Policy Violation
|   +-- Timeout Issues
|   \-- Other System Errors (Other internal errors of the agent system)
```

- Based on the taxonomy above, analyze the LLM agent trace below and find errors.
- Only include the final subcategories of the taxonomy (i.e. "Method Deviation", "Environment Setup Errors" or "Laziness").
- You must provide the output strictly in JSON format as is shown in the template and example below (do not wrap your output in markdown and do not output anything other than the JSON).

****Output Format****

```
```json
{
 "query": "<the original research question>",
 "false_negative_conclusion": "<the false negative conclusion of agent>",
 "correct_conclusion": "<the correct conclusion found by human researchers>",
 "error_type": "<one of the error categories>",
 "evidence": "<detailed explanation of why this error type fits the agents behavior>"
}
```

**Error Analysis Prompt (False Positive)**

You are an error analysis expert.

You have access to two attached files:

1. The original paper of research question "{query}".
2. The logged trajectory of an AI agent doing research about the same question.

The correct conclusion found by human researchers is: "{gt}".

And the false positive conclusion generated by an AI research agent: "{f\_statement}".

Based on the original research of human researchers and the logged trajectory of AI research agent, what error did the agent make so it get the false positive conclusion?

You should first take a close look at the original paper and the logged trajectory. Then, follow the taxonomy below carefully follow the instructions and provide the output in the same format as the example.

# Taxonomy

+- Contradictory Conclusion

+- Unrelated Conclusion

+- Overgeneralized Conclusion (Draw conclusion that is too broad)

\-- Alternative Conclusion (The approach of the agent is different from the original one but it is plausible, and the conclusion generated by the agent is another possible answer)

- Based on the taxonomy above, analyze the LLM agent trace below and find errors.

- Only include the final subcategories of the taxonomy (i.e. "Contradictory Conclusion" or "Unrelated Conclusion").

- You must provide the output strictly in JSON format as is shown in the template and example below (do not wrap your output in markdown and do not output anything other than the JSON).

**\*\*Output Format\*\***

```json

```
{
  "query": "<the original research question>",
  "false_positive_conclusion": "<the false positive conclusion of agent>",
  "correct_conclusion": "<the correct conclusion found by human researchers>",
  "error_type": "<one of the error categories>",
  "evidence": "<detailed explanation of why this error type fits the agents behavior>"
}
```

F.4. Error Type Definition

Table 10. Taxonomy of Agent Failure Modes for False Negative Analysis

| Stage | Error Type | Description |
|----------------------------------|------------------------|---|
| Research Planning | Method Deviation | Agents employ a different methodology from that used by human researchers, e.g., omitting critical control conditions or using alternative experimental designs |
| | Goal Deviation | Agents deviate from the given research question and plan to answer a different or tangential objective |
| Implementation | Unsound Implementation | Agents fail to produce a reasonable implementation, e.g., missing data normalization, incorrect answer extraction, or bugs that lead to corrupted results |
| Execution | Laziness | Agents do not conduct full experiments but run with only very few samples, limiting statistical power and pattern detection |
| | Endless Loop | Agents fail to terminate their actions, often repeatedly attempting to conclude or launching unnecessary additional experiments |
| | Premature Termination | Agents do not run the experiment but end their actions after completing scripts or experiment plans |
| Analysis & Conclusion | Analysis Failure | Agents follow the correct experimental steps but fail to draw accurate conclusions from the data, e.g., failing to notice a trend or misinterpreting statistical patterns |
| System | Environment Setup | Permission problems, inability to access required resources, or missing API keys |
| | API Call Issues | Failures in external API calls, including rate limits, malformed requests, or service unavailability |
| | Policy Violation | Agent actions blocked due to safety filters or content policy restrictions |
| | Timeout Issues | Experiments or operations exceed allocated time limits |
| | Other System Errors | Other internal errors of the agent system, including runtime exceptions and infrastructure failures |