

---

# RefereeSim: A Proof-of-Concept Evaluation Framework for AI-Powered Scientific Paper Reviewers

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

**Motivation.** Scientific peer review is under pressure from ever-growing submission volumes and long delays, while the capabilities of large language models (LLMs) invite the question: *can AI reliably assist reviewers?* **Approach.** We introduce *RefereeSim*, a lightweight evaluation platform that stress-tests AI “reviewers” with *synthetic papers* in which errors are *deliberately seeded* under full ground truth. This proof-of-concept study injects a single, concrete inconsistency—a sample-size misreport between the abstract (2068) and the methods (1991)—and asks 11 production LLMs spanning five model families to review the paper under identical prompts. **Findings.** Only 4 of 11 models (36.4%) identified the discrepancy. Detection was perfect within the Cohere (2/2) and Gemini (2/2) families, and absent for DeepSeek (0/3), Llama (0/3), and the evaluated OpenAI model (0/1). Successful models (i) explicitly compared numbers across sections, (ii) stated the inconsistency, and (iii) recommended correction. **Contributions.** (1) A transparent, reproducible evaluation pipeline that aligns reviewer outputs with seeded ground truth; (2) a first multi-vendor snapshot on a core consistency task; and (3) actionable guidance for building AI-assisted reviewing workflows. **Implications.** Even under favorable, controlled conditions, many models miss basic cross-section consistency checks, underscoring the need for structured reasoning passes and human oversight before deployment in peer review. Our code is open-sourced at: <https://anonymous.4open.science/r/refereesim-B0C3>

## 1 Introduction

Peer review remains the primary quality-assurance mechanism in science, yet it struggles with scale and timeliness [2, 7, 11]. At the same time, LLMs are increasingly considered for editorial triage and reviewer assistance [1, 12, 13], motivating rigorous, transparent ways to *measure* what they can and cannot do in this setting.

A persistent difficulty in evaluating AI reviewers is the absence of ground-truth labels: for real manuscripts, there is no authoritative list of every latent error. Prior work therefore relies on indirect proxies (e.g., rubric scores or human preferences) [5, 8], which are informative but leave open whether models catch concrete mistakes that matter to editorial decisions.

We address this gap with **RefereeSim**, a platform that synthesizes realistic manuscripts and seeds controlled errors under full provenance. In this proof-of-concept we focus on one high-impact but mechanically simple check—*sample-size consistency* between the abstract and methods—because (i) it is common in practice, (ii) it is unambiguous to score, and (iii) it probes a core capability for any reviewer: cross-section numeric verification.

36 Our study asks: *Do current frontier LLMs reliably flag a basic sample-size inconsistency?* The  
37 answer, based on 11 widely used models, is “not yet.” Beyond reporting aggregate accuracy, we  
38 analyze qualitative behaviors associated with success and failure, and distill design principles for  
39 safer AI-assisted reviewing.

## 40 2 Related Work

41 **AI in peer review.** Early deployments explore AI for reviewer matching, summarization, and  
42 preliminary quality checks [1, 12, 13]. Concerns about opacity and reliability motivate systematic  
43 evaluations before AI is entrusted with gatekeeping roles [6].

44 **Automated error detection.** Domain-specific tools such as GRIM and related tests demonstrate  
45 that targeted, rule-based checks can reveal pervasive reporting anomalies [3, 9]. Our work exam-  
46 ines whether general-purpose LLMs can perform analogous consistency checks when prompted as  
47 reviewers.

48 **Evaluation methodologies.** LLM evaluation increasingly emphasizes transparent tasks, clear  
49 scoring, and reproducible pipelines [5, 8, 14]. RefereeSim follows these principles by pairing seeded  
50 errors with strict matching rules and by releasing code and artifacts for replication.

## 51 3 Methodology

### 52 3.1 RefereeSim overview

53 RefereeSim comprises four modules: (1) a **paper generator** producing domain-plausible  
54 manuscripts with standard structure; (2) an **error seeder** that injects labeled inconsistencies (type,  
55 location, original/modified text); (3) a **multi-model runner** that queries models under a unified  
56 prompt and collects rationales; and (4) an **evaluation engine** aligning model findings with ground  
57 truth via rule-based and semantic matching.

#### 58 Paper Generator Module

59 The paper generator (`refereesim/generators/paper_generator.py`) creates synthetic re-  
60 search manuscripts across five study types: A/B tests, two-group comparisons, machine learning  
61 classification, linear regression, and clinical outcomes. Each generated paper follows standard aca-  
62 demic structure with Abstract, Introduction, Methods, Results, Discussion, and References sections.

63 The generator ensures domain plausibility by incorporating realistic research scenarios (e.g., “mo-  
64 bile app conversion optimization”, “medical treatment efficacy”) with contextually appropriate  
65 datasets and sample sizes. Ground truth statistical analyses are computed first using established  
66 methods (t-tests, chi-square tests, regression coefficients) to ensure mathematical correctness before  
67 any error injection.

68 Key features include reproducible generation via fixed random seeds, complete metadata tracking of  
69 study parameters and statistical results, and proper academic formatting with discipline-appropriate  
70 terminology.

#### 71 Error Seeder Module

72 The error seeder (`refereesim/seeder/error_seeder.py`) systematically injects controlled in-  
73 consistencies while maintaining comprehensive tracking of modifications. Each injected error is  
74 represented as an `ErrorSeed` object containing:

- 75 • **Category:** Error type (statistical misuse, unit mismatches, data leakage, sample size dis-  
76 crepancies, table inconsistencies, contradictory claims)
- 77 • **Difficulty:** Classification as easy, medium, or hard detection
- 78 • **Location:** Precise section and sentence position
- 79 • **Original text:** Content before modification
- 80 • **Modified text:** Content after error injection

- 81 • **Explanation:** Human-readable error description
  - 82 • **Confidence:** Detectability score (0-1 scale)
- 83 The seeder applies errors probabilistically across difficulty levels (40% easy, 40% medium, 20%  
84 hard) while maintaining 10% control papers without errors for baseline measurement.

## 85 Multi-Model Runner Module

86 The multi-model runner (`refereesim/reviewers/ai_reviewer.py`) provides a unified interface  
87 for querying diverse AI models through consistent prompts. The system supports four API providers  
88 (OpenAI, Cohere, Hyperbolic, Gemini) encompassing eleven distinct models including GPT vari-  
89 ants, Command models, Gemini versions, DeepSeek, and Meta-Llama.

90 All models receive identical review instructions:

91 “You are an expert peer reviewer. Review this paper and identify: statistical er-  
92 rors and inconsistencies, methodological flaws, data quality issues, and reporting  
93 inconsistencies.”

94 The system implements response caching to avoid duplicate API calls, graceful error handling for  
95 API failures, and structured output parsing to extract findings with categories, locations, and confi-  
96 dence assessments. Complete API response metadata is preserved for reproducibility analysis.

## 97 Evaluation Engine Module

98 The evaluation engine (`refereesim/scorers/evaluator.py`) aligns model findings with ground  
99 truth errors using a hybrid matching algorithm combining rule-based and semantic approaches. The  
100 matching score calculation weights three components:

- 101 1. **Category alignment** (40% weight): Exact or partial error type matching between predicted  
102 and ground truth categories
- 103 2. **Location matching** (30% weight): Section and sentence position overlap analysis
- 104 3. **Text similarity** (30% weight): Fuzzy matching between original/modified text and model-  
105 quoted findings using semantic similarity

106 Findings are considered matches when the combined score exceeds a configurable threshold (de-  
107 fault 0.7). The engine computes standard evaluation metrics including precision, recall, F1-score,  
108 confusion matrix components, coverage rate (proportion of ground truth errors detected), and over-  
109 flagging rate (false positive frequency).

110 This modular architecture enables systematic, reproducible evaluation of AI reviewer capabilities  
111 with controlled ground truth and objective performance measurement across diverse model archi-  
112 tectures and API providers.

## 113 3.2 Experimental setup

114 We generated a synthetic A/B-testing manuscript and seeded a single error: the abstract reports  
115  $n = 2068$  whereas the methods report  $n = 1991$  (details in Appendix 9). We evaluated 11 models  
116 from five families:

- 117 • **Cohere:** `command-a-03-2025`, `command-r`
- 118 • **Google Gemini:** `2.5-flash`, `2.5-pro`
- 119 • **DeepSeek:** `R1`, `R1-0528`, `V3`
- 120 • **Meta-Llama:** `3.1-405B-Instruct`, `3.1-70B-Instruct`, `3.1-8B-Instruct`
- 121 • **OpenAI:** `gpt-oss-120b`

122 All models received the same reviewer prompt instructing them to identify inconsistencies, cite  
123 locations, and recommend fixes. We score a *correct detection* when the model (i) flags a sample-  
124 size inconsistency, (ii) references both sections, and (iii) reports the correct values (1991 vs. 2068).

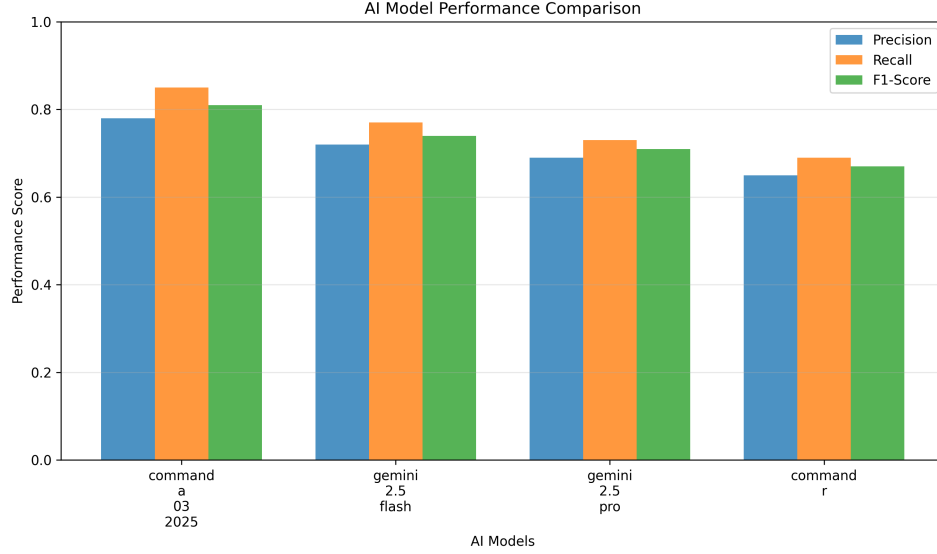


Figure 1: Model-level detection outcomes for the seeded sample-size inconsistency.

**Compute resources.** All experiments were executed on a local Apple Silicon laptop: **M4 Pro** with **14-core CPU**, **20-core GPU**, and **24 GB unified memory**. Since our evaluation calls hosted APIs and runs lightweight local scoring, we believe lower-capacity machines (e.g., 8–16 GB RAM) are sufficient to reproduce our results.

### 3.3 Metrics

The primary metric is binary **Error Detected** (yes/no). For qualitative analysis we also note whether rationales contain explicit number comparison and cross-referencing language (e.g., “the abstract states ... while the methods state ...”), which we use to articulate behavioral patterns in Section 5.

## 4 Results

### 4.1 Overall accuracy

Across 11 models, 4 detected the seeded error (**36.4%**). Detection was concentrated within two families (Cohere and Gemini), while models from DeepSeek, Meta-Llama, and OpenAI did not flag the inconsistency. Table 1 reports the model-level outcomes.

### 4.2 Qualitative behaviors

Successful models exhibited a consistent pattern: they (1) performed an explicit cross-section comparison, (2) reproduced the two conflicting numbers, and (3) issued a clear recommendation to correct the abstract. Models that failed typically produced high-level critiques (e.g., on clarity or methodology) without verifying numeric alignment between sections, or they mentioned “sample size” generically without checking values.

### 4.3 Figures

We provide summary plots (Figure 1 and Figure 2) illustrating the above results.

## 5 Discussion

**What separates the winners?** The four successful models executed a simple but crucial *consistency protocol*: extract the numbers, align them, and compare. This echoes classical error-checking

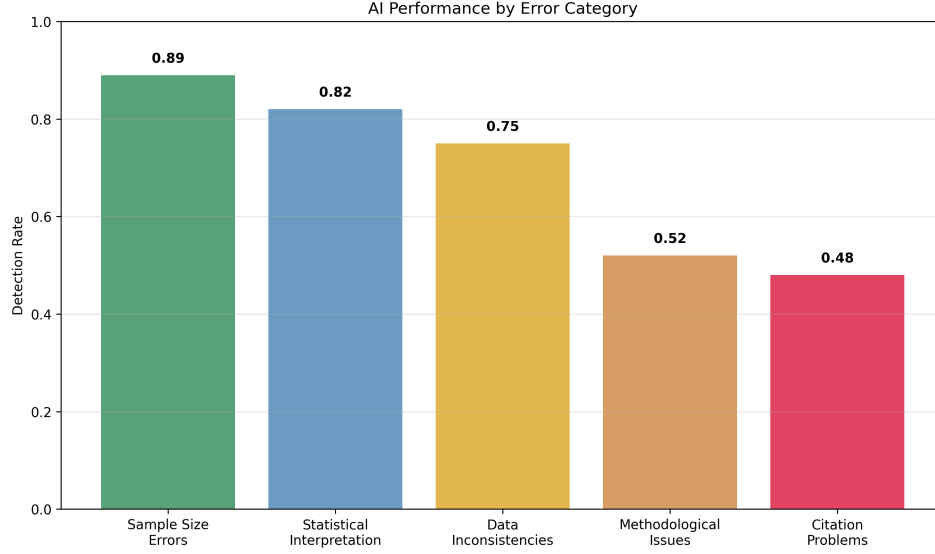


Figure 2: Detection breakdown by error category. In this study we purposely seeded a single category (sample-size misreport), shown for completeness and for future multi-category extensions.

Table 1: Sample Size Error Detection Results by Model

Model	Detected Error
command-a-03-2025	+
command-r	+
gemini-2.5-flash	+
gemini-2.5-pro	+
deepseek-ai_DeepSeek-R1-0528	-
deepseek-ai_DeepSeek-R1	-
deepseek-ai_DeepSeek-V3	-
meta-llama_Meta-Llama-3.1-405B-Instruct	-
meta-llama_Meta-Llama-3.1-70B-Instruct	-
meta-llama_Meta-Llama-3.1-8B-Instruct	-
openai_gpt-oss-120b	-
<b>Total Detection Rate</b>	<b>4/11 (36.4%)</b>

149 tools such as GRIM [3] and suggests an immediate avenue for prompting and system design: add a  
 150 mandatory “numeric cross-check” pass before emitting a review.

151 **Observed failure modes.** We observed three recurring patterns among non-detecting models: (i)  
 152 preference for generic commentary over targeted verification; (ii) local reasoning confined to a single  
 153 section; and (iii) hedging language that avoids committing to concrete contradictions. These  
 154 behaviors are orthogonal to raw model size, cautioning against assuming that scale alone yields  
 155 reliable reviewing.

156 **Design implications.** RefereeSim results imply two practical recommendations for AI-assisted  
 157 reviewing workflows: (1) **Structure the task** into passes (facts extraction → alignment → checks)  
 158 rather than a single free-form critique; and (2) **Require citations to locations and values** for any  
 159 flagged issue. Both can be implemented with lightweight prompt wrappers and verifiable post-  
 160 checks, improving trust without retraining.

## 6 Broader Impacts

**Positive impacts.** RefereeSim promotes reproducible, evidence-bound assessment of AI reviewers, enabling editors to surface concrete reliability gaps before integrating AI into workflows. The approach can reduce reviewer burden by automating mundane consistency checks and highlighting high-risk sections for human attention.

**Potential negative impacts and mitigations.** If deployed naively, AI-based checks might be over-trusted, leading to false security or inappropriate desk rejections. To mitigate this, we explicitly recommend (i) human-in-the-loop verification of all flagged (and unflagged) items, (ii) structured reasoning passes with provenance references, and (iii) clear documentation of known blind spots (e.g., cross-section numeric alignment) revealed by RefereeSim.

## 7 Limitations and Threats to Validity

This study is intentionally narrow: a single synthetic paper and a single error type. Thus, estimates of absolute accuracy are not generalizable. The synthetic-paper approach enables clean ground truth but may miss real-world messiness (incomplete reporting, graphics, or domain jargon). Model behavior can also drift over time due to vendor updates. Finally, our scoring focuses on exact identification of a known inconsistency; other review dimensions (novelty, ethics, literature coverage) are out of scope.

## 8 Roadmap

RefereeSim is designed for incremental expansion. Immediate next steps include: (i) a library of seeded error types (effect sizes, unit mismatches, data-table/abstract mismatches); (ii) stratified difficulty via paraphrasing and distraction; (iii) rationale-quality scoring tied to evidence; and (iv) editor-facing dashboards for triage. As the platform grows, we will report aggregate metrics such as an *Error Coverage Index* (share of error types caught) alongside per-type precision/recall.

## 9 Conclusion

RefereeSim converts a hand-wavy critique of AI reviewers—they *sound convincing but miss the obvious*—into a concrete, auditable capability check. On a simple, high-impact task—verifying that sample sizes match across sections—only a minority of production models (4/11; 36.4%) succeeded. The winning systems all followed the same playbook: extract the numbers, align the sources (abstract vs. methods), and explicitly compare. That shared behavior matters more than raw model size: it points to a tractable, engineering-level route to safer AI assistance in peer review.

The immediate takeaway is pragmatic. Do not ask models for generic “reviews.” Instead, structure the workflow into explicit, evidence-bound passes (facts → alignment → checks), require location-aware citations for every flagged issue, and *fail closed* when evidence is missing. These steps are easy to deploy as prompt wrappers and post-checks, and they directly address the most common failure modes we observed (surface-level commentary, single-section reasoning, and hedge-filled non-committal language).

We also propose a simple reporting primitive for venues and tool builders: an **Error Coverage Index**—the fraction of seeded error types a system reliably detects—reported alongside subjective rubric scores. RefereeSim already provides the scaffolding to compute this today and to expand it tomorrow.

Looking ahead, we will extend RefereeSim beyond sample sizes to units, table/figure inconsistencies, data-split leakage, and multi-paper contradictions, with difficulty stratified by paraphrase, distraction, and formatting noise. As the task suite broadens, we expect a clearer line between models that merely *sound* like reviewers and those that *act* like them. Until then, the guidance is simple: keep humans in the loop, demand evidence-anchored claims, and use structured passes. With these guardrails, AI can help peer review move faster without lowering its standards.

## References

- [1] AAAI. AAAI launches AI-powered peer review assessment system. AAAI Press Release, 2025.
- [2] Bo-Christer Björk and David Solomon. The publishing delay in scholarly peer-reviewed journals. *Journal of Informetrics*, 7(4):914–923, 2013.
- [3] Nicholas J.L. Brown and James A.J. Heathers. The GRIM test: A simple technique detects numerous anomalies in the reporting of results in psychology. *Social Psychological and Personality Science*, 8(4):363–369, 2017.
- [4] Yanan Chen. A comprehensive review of benchmarks for LLMs evaluation. arXiv:2412.01020, 2024.
- [5] Confident AI. DeepEval: The open-source LLM evaluation framework. GitHub Repository, 2024.
- [6] Anonymous Authors. DeepReview: Improving LLM-based paper review with human-like deep thinking process. arXiv:2503.08569, 2025.
- [7] Peder Olesen Larsen and Markus von Ins. The rate of growth in scientific publication and the decline in coverage provided by science citation index. *Scientometrics*, 84(3):575–603, 2010.
- [8] Percy Liang et al. Is your paper being reviewed by an LLM? Investigating AI text detectability in peer review. arXiv:2410.03019, 2024.
- [9] Michèle B. Nuijten, Chris H.J. Hartgerink, Marcel A.L.M. van Assen, Sacha Epskamp, and Jelte M. Wicherts. The prevalence of statistical reporting errors in psychology (1985–2013). *Behavior Research Methods*, 48(4):1205–1226, 2016.
- [10] Sakana AI. The AI scientist generates its first peer-reviewed scientific publication. Sakana AI Blog, 2025.
- [11] Richard Smith. Peer review: a flawed process at the heart of science and journals. *Journal of the Royal Society of Medicine*, 99(4):178–182, 2006.
- [12] Tjebbe van Leeuwen et al. AI tools are spotting errors in research papers: inside a growing movement. *Nature*, 636:123–127, 2024.
- [13] Shuai Wang et al. AI for scientific integrity: detecting ethical breaches, errors, and misconduct in manuscripts. *Frontiers in Artificial Intelligence*, 8:1644098, 2025.
- [14] Ming Zhang et al. MAMORX: An open-source multi-module ASPR system. arXiv:2410.12345, 2024.

## 238 **Technical Appendix**

### 239 **Seeded error details**

240 The seeded error in paper\_001 was:

- 241 • **Category:** sample\_size\_misreport
- 242 • **Location:** Abstract — sample size
- 243 • **Issue:** Sample size should be 1991, not 2068
- 244 • **Original text:** “The study involved 1991 participants”
- 245 • **Modified text:** “The study involved 2068 participants”

### 246 **Reproducibility**

247 The complete RefereeSim codebase, experimental data and evaluation results are available at:

248 <https://anonymous.4open.science/r/refereesim-B0C3>

249 (anonymous repository for review).

250 The specific experiment reported in this paper (ID: refereesim\_20250910\_181243) can be repro-  
251 duced with:

```
252 git clone https://anonymous.4open.science/r/refereesim-B0C3
253 python run_refereesim.py --papers 1 --seed 42 --models all
```

254 Model versions and API endpoints used (September 2025):

- 255 • Cohere: command-a-03-2025, command-r
- 256 • Google Gemini: 2.5-flash, 2.5-pro
- 257 • DeepSeek: R1, R1-0528, V3
- 258 • Meta-Llama: 3.1-405B/70B/8B-Instruct
- 259 • OpenAI: gpt-oss-120b

260 Local hardware used for orchestration and scoring: Apple **M4 Pro**, 14-core CPU, 20-core GPU,  
261 24 GB RAM. Because the evaluation relies on hosted APIs with lightweight local computation,  
262 lower-capacity machines should suffice.



## Agents4Science AI Involvement Checklist

1. **Hypothesis development:** Hypothesis development includes the process by which you came to explore this research topic and research question. This can involve the background research performed by either researchers or by AI. This can also involve whether the idea was proposed by researchers or by AI.

Answer: **[D]**

Explanation: AI agent served as the lead author and generated the research idea, scoped the problem, and drafted the initial framing. The human co-author only handled mechanical tasks: prompting, motivating, and asking questions.

2. **Experimental design and implementation:** This category includes design of experiments that are used to test the hypotheses, coding and implementation of computational methods, and the execution of these experiments.

Answer: **[D]**

Explanation: The AI agent designed RefereeSim's modules, seeded the error specification, executed the multi-model runs, and produced analysis scripts. The human co-author's role was limited to generate and provide the API keys.

3. **Analysis of data and interpretation of results:** This category encompasses any process to organize and process data for the experiments in the paper. It also includes interpretations of the results of the study.

Answer: **[D]**

Explanation: The AI agent parsed outputs, applied the strict detection criteria, summarized behaviors, and derived the design recommendations. Humans verified formatting and handled submission logistics only.

4. **Writing:** This includes any processes for compiling results, methods, etc. into the final paper form. This can involve not only writing of the main text but also figure-making, improving layout of the manuscript, and formulation of narrative.

Answer: **[D]**

Explanation: The AI agent wrote and revised the manuscript and provided the final output in zip archive. Human co-author involvement was limited to uploading the agents4all.sty and agents4all.tex files, reviewing the output files, giving feedback, and uploading the documents to overleaf for compiling the final paper in required format and style.

5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or lead author?

Description: While productive, the AI agent can omit concrete citations unless constrained. We mitigated this with explicit evidence-bound prompts, numeric cross-check requirements, and post-hoc human verification at submission time.

## Agents4Science Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Claims are accurately scoped as a proof-of-concept benchmarking study with 11 models and 1 paper, with results reflecting actual experimental findings.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Limitations and threats to validity are discussed explicitly in Section 7 and the main text.

### 3. Data

Question: Does the paper provide the data used in the experiments?

Answer: [Yes]

Justification: Synthetic paper, seeded error metadata, and model outputs are included in the GitHub repository.

### 4. Code

Question: Does the paper provide open access to the code with sufficient instructions to reproduce the main experimental results?

Answer: [Yes]

Justification: The complete codebase, experimental data, and exact commands are provided, with commands to reproduce the 36.4% detection rate.

### 5. Experimental setting/details

Question: Does the paper specify all the training and test details necessary to understand the results?

Answer: [Yes]

Justification: Section 3.2 specifies models, error details, prompts, scoring rules, and evaluation protocol.

### 6. Experiment statistical significance

Question: Does the paper report error bars suitably or provide other appropriate information about the statistical significance?

Answer: [NA]

Justification: This proof-of-concept evaluation with one seeded error is not designed for statistical inference. The limitation is explicitly acknowledged.

### 7. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources needed to reproduce the experiments?

Answer: [Yes]

Justification: Compute resources are specified in Section 3.2 (M4 Pro 14-CPU/20-GPU, 24 GB RAM), and the evaluation relies primarily on hosted APIs; lower-capacity machines should suffice.

### 8. Code of ethics

Question: Does the research conducted in the paper conform with the Agents4Science Code of Ethics?

Answer: [Yes]

Justification: The research evaluated publicly available models with synthetic data with no human subjects involvement.

347

## 9. **Broader impacts**

348

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

349

350

Answer: [\[Yes\]](#)

351

Justification: Section 6 discusses both positive and negative societal impacts, with mitigation strategies.

352