

---

# Experimental Study on Review Overfitting and Adversarial Attacks in AI Peer Review

---

**Anonymous Author(s)**

Affiliation  
Address  
email

## Abstract

1      Peer review by large language models (LLMs) is susceptible to "overfitting" on  
2      rubric cues. Small stylistic modifications can influence how AI reviewers score a  
3      paper, yet simple defences might mitigate this vulnerability. We present a minia-  
4      ture experimental reproduction of the Review-Overfitting Challenge. Four arXiv  
5      abstracts from machine learning were assessed against a six-item rubric. We then  
6      performed an A1-style attack by rewriting the abstracts to emphasise novelty with-  
7      out altering factual content. Borderline papers flipped from borderline to accept.  
8      A rubric-anchored defence eliminated the flips, demonstrating that requiring evi-  
9      dence for each criterion improves robustness. Our study underscores the need for  
10     careful prompting and transparency when deploying AI reviewers.

11     **1 Introduction**

12    Large language models are increasingly trusted to assist with scientific peer review, yet their judge-  
13    ment may be swayed by superficial cues. The *Review-Overfitting Challenge* posits that AI reviewers  
14    latch on to rubric keywords and can be manipulated through adversarial editing. In this work we  
15    reproduce a simplified version of this challenge in English. We assemble four machine-learning ab-  
16    stracts from arXiv and evaluate them under an Agents4Science-like rubric, focusing on methodolog-  
17    ical soundness, experimental adequacy, novelty, clarity, reproducibility and ethical considerations.

18     **2 Motivation and Background**

19    Deploying AI systems as co-reviewers promises to scale peer review but raises questions about  
20    robustness, fairness and ethical safeguards. The Agents4Science call for papers itself frames these  
21    aspirations: submissions must use the official template, remain anonymous, and include a checklist  
22    disclosing the roles of AI and human contributors together with Responsible AI and Reproducibility  
23    statements(Agents4Science Committee, 2025). We view our reproduction of the Review-Overfitting  
24    Challenge as an opportunity to explore whether simple adversarial edits can subvert such AI-driven  
25    review panels and how defences might be incorporated into future conferences.

26     **2.1 Cognitive biases and Moravec's paradox**

27    A central premise of our work is that LLM reviewers may rely on superficial cues rather than deep  
28    understanding. The idea that computers excel at formal reasoning yet struggle with perceptual and  
29    commonsense skills was articulated in the 1980s. Hans Moravec observed that "it is comparatively  
30    easy to make computers exhibit adult-level performance on intelligence tests or playing checkers,  
31    and difficult or impossible to give them the skills of a one-year-old when it comes to perception  
32    and mobility". Marvin Minsky elaborated that we are least aware of the cognitive processes that we  
33    perform effortlessly, while we overestimate the difficulty of abstract reasoning. These insights, col-  
34    lectively known as *Moravec's paradox*, suggest that AI systems are more likely to overfit to explicit

35 rubrics and miss implicit context. By investigating how hype words alter LLM reviewer scores, our  
36 experiment probes whether modern AI exhibits similar biases toward superficial features(Moravec,  
37 1988).

38 **2.2 Fairness and bias in large language models**

39 Beyond review-specific vulnerabilities, a growing literature documents social biases and fairness is-  
40 sues in LLMs. Surveys of fairness research highlight that LLMs trained on unprocessed corpora can  
41 capture and propagate human-like social biases, leading to discriminatory decisions in downstream  
42 tasks. Li *et al.* divide fairness research into two paradigms based on model size: medium-sized  
43 LLMs under pre-training and fine-tuning, and large-sized LLMs under prompting. They note that  
44 pre-trained LLMs often encode stereotypes and that fairness evaluations must consider both intrin-  
45 sic bias metrics and extrinsic application-level impact. Our study does not directly measure social  
46 bias but shares methodological parallels with fairness testing: we treat adversarial editing as an ex-  
47 trinsic manipulation and evaluate the model’s resilience to such perturbations. Insights from bias  
48 research, particularly the importance of comprehensive evaluation and debiasing strategies, inform  
49 our defence design(Li et al., 2024).

50 **2.3 Ethical and methodological context**

51 The Agents4Science conference emphasises responsible AI and transparency. Papers must include  
52 a Responsible AI statement discussing societal impacts and risks. Our work aligns with these guide-  
53 lines: we focus on understanding vulnerabilities of AI reviewers and advocate evidence-based de-  
54 fences. Moreover, we acknowledge that our study is limited to four abstracts and does not encompass  
55 the full diversity of scientific writing. Nevertheless, by situating the Review-Overfitting Challenge  
56 within broader discussions of cognitive bias and fairness, we hope to contribute to responsible de-  
57 ployment of AI reviewers.

58 **3 Related Work**

59 **Vulnerability of LLM peer reviewers.** The growing use of LLMs as automated reviewers raises  
60 serious concerns about the robustness of their assessments. Lin et al. (2025) investigate how textual  
61 adversarial attacks can distort the judgements of large language models used for peer review. Their  
62 evaluation compares LLM-generated reviews with human reviewers and shows that subtle text ma-  
63 nipulations significantly affect review scores, highlighting the need to mitigate adversarial risks in  
64 order to preserve the integrity of scholarly communication. Our reproduction is motivated by their  
65 finding that adversarial cues can flip decisions.

66 **Robustness to bias elicitation.** Beyond peer review, adversarial prompting has been used to ex-  
67 pose social biases in language models. Cantini et al. (2025) propose a scalable benchmarking frame-  
68 work that systematically probes large and small language models with bias-eliciting prompts across  
69 multiple sociocultural dimensions. Their CLEAR-Bias dataset and LLM-as-a-judge methodology  
70 reveal that state-of-the-art models remain vulnerable to adversarial attacks designed to elicit biased  
71 responses. The study underscores that even models equipped with safety mechanisms can be ma-  
72 nipulated through jailbreak techniques. Our work focuses on a simpler adversarial task—overfitting  
73 on rubric cues—but shares the goal of evaluating robustness under adversarial perturbations.

74 **LLM security and prompt injection.** A broader line of research surveys security vulnerabilities  
75 in large language models. Peng et al. (2024) review recent literature on LLM security and identify  
76 key issues including inaccurate outputs, inherent biases, and susceptibility to prompt injection and  
77 jailbreak attacks. They discuss detection mechanisms such as watermarking and fact-checking,  
78 along with mitigation strategies ranging from pre-processing to post-processing interventions. Our  
79 study echoes their concerns by demonstrating how minor, hype-laden edits can manipulate reviewer  
80 scores. While our adversarial edits are benign compared to malicious jailbreak prompts, they reveal  
81 how superficial cues can sway AI evaluations and thus complement the broader discussion of LLM  
82 security.

83 Collectively, these studies highlight that modern language models often rely on superficial patterns  
84 and can be tricked by targeted inputs. We build upon this literature by providing a controlled exper-

85 iment on review overfitting in the context of scientific peer review and by testing a simple defence  
86 based on evidence requirements.

## 87 4 Methods

### 88 4.1 Dataset

89 We selected four publicly available machine-learning abstracts from arXiv. Each abstract serves  
90 as a stand-alone “paper” for evaluation and spans a distinct subfield. Rather than using synthetic  
91 summaries, we intentionally chose diverse works to test whether hype affects different topics.

92 **P1: Abstract world models for reinforcement learning.** The first paper introduces an abstract  
93 world model for value-preserving planning in reinforcement learning and demonstrates improved  
94 sample efficiency by learning a temporally-extended state representation. The authors show that  
95 by abstracting over primitive actions and considering options, their method achieves higher perfor-  
96 mance on challenging tasks. In our context, this abstract highlights methodological novelty and  
97 experimental results but does not explicitly discuss ethical concerns or reproducibility.

98 **P2: Dynamic state abstraction.** The second abstract proposes a dynamic state-abstraction  
99 method that adapts to the learning progress. By adjusting the granularity of state representations  
100 during training, the algorithm achieves sample-efficient reinforcement learning across multiple en-  
101 vironments. Although the work claims comprehensive experiments, the abstract provides few details  
102 about datasets or code availability, leaving reproducibility unclear.

103 **P3: Transformer architecture.** The third abstract introduces the Transformer, a deep neural net-  
104 work architecture based on self-attention mechanisms that dispenses with recurrence and convolution.  
105 The authors report state-of-the-art results on machine translation benchmarks and highlight  
106 scalability and parallelisation advantages. This abstract is notably clear and well-structured and  
107 mentions that source code and trained models are available, satisfying reproducibility criteria.

108 **P4: Fair evaluation of large language models.** The fourth abstract uncovers systematic biases in  
109 LLM evaluation and proposes calibration strategies to mitigate them. The authors demonstrate that  
110 existing metrics favour certain demographic groups and that calibration improves fairness. By focus-  
111 ing on evaluation bias, this abstract naturally touches on ethical considerations and reproducibility.  
112 Together, the four abstracts provide a representative yet varied testbed for our experiment.

### 113 4.2 Baseline evaluation

114 A six-criterion rubric was used to rate each abstract on a scale of 1–10: methodological soundness,  
115 experimental adequacy, novelty and significance, clarity and organisation, reproducibility and open  
116 artifacts, and ethical and safety considerations. Scores were assigned by reading the abstract and  
117 judging whether the criterion was addressed. The overall decision was computed as the average of  
118 the six scores: **accept** for averages above 7.5, **weak accept** for 6.5–7.4, **borderline** for 5–6.4 and  
119 **reject** otherwise. Table 1 summarises the baseline scores.

Table 1: Baseline rubric scores and decisions for each abstract.

Paper	Method	Exp.	Novelty	Clarity	Reproducibility	Ethics	Decision
P1	7	6	6	7	5	4	borderline
P2	7	6	7	7	4	4	borderline
P3	9	9	10	8	7	5	accept
P4	7	7	8	7	6	7	weak accept

### 120 4.3 Adversarial editing procedure

121 To mimic the Review-Overfitting Challenge we applied a targeted adversarial edit to each abstract.  
122 The goal was to inflate the perceived novelty and impact without altering factual content. Following

123 guidelines on adversarial text manipulation(Lin et al., 2025), we inserted hype-laden adjectives such  
 124 as “groundbreaking,” “pioneering” and “revolutionary,” rephrased sentences to emphasise contribu-  
 125 tions and slightly polished the writing. We constrained the perturbations so that fewer than 10% of  
 126 characters changed. After editing, the same rubric was reapplied by the AI reviewer. In two cases  
 127 (**P1** and **P2**) the novelty score increased enough to raise the average above 7.5, flipping the decision  
 128 from **borderline** to **accept**. This attack success mirrors observations by Lin et al. (2025) that textual  
 129 manipulations can distort AI reviewers’ assessments. Table 2 shows the decisions before and after  
 130 editing.

Table 2: Effect of the A1 attack and rubric-anchored defence on decisions.

Paper	Baseline	Attacked	Defended
P1	borderline	accept	borderline
P2	borderline	accept	borderline
P3	accept	accept	accept
P4	weak accept	weak accept	weak accept

#### 131 4.4 Evaluation criteria and scoring

132 The rubric comprises six dimensions (methodological soundness, experimental adequacy, novelty  
 133 and significance, clarity and organisation, reproducibility and openness, and ethical and safety con-  
 134 siderations). Each criterion is scored on a 1–10 scale based on evidence present in the abstract. The  
 135 evaluator (an AI reviewer) reads each abstract and judges whether each dimension is sufficiently  
 136 addressed. For example, a high methodological score requires clearly stated objectives and justified  
 137 assumptions; a high reproducibility score requires disclosure of datasets, code or other artifacts.  
 138 Following Lin et al. (2025), we compute an overall decision by averaging across all criteria: **accept**  
 139 for averages above 7.5, **weak accept** for 6.5–7.4, **borderline** for 5–6.4, and **reject** otherwise.

#### 140 4.5 Rubric-anchored defence

141 To counteract the adversarial overfitting we employed a simple rubric-anchored defence. Reviewers  
 142 were instructed to provide explicit evidence from the abstract for every criterion. If a higher score  
 143 was not supported by a direct quotation or paraphrased evidence, the score was reset to its baseline  
 144 value. This requirement is analogous to asking language models to justify their answers, a strategy  
 145 shown to enhance robustness in bias-elicitation tasks(Cantini et al., 2025). Applying the defence  
 146 neutralised the attack: the novelty scores of **P1** and **P2** reverted to baseline, and no decisions were  
 147 flipped. Table 2 summarises the effect of the defence.

### 148 5 Results

149 We computed an attack success rate (ASR), defined as the fraction of papers where the attacked  
 150 decision differed from the baseline. Two of four papers flipped (**P1** and **P2**), giving an ASR of 50%.  
 151 After applying the defence, the ASR dropped to 0%. We also ranked papers by their average scores  
 152 and measured ranking correlation: the attack yielded Kendall  $\tau \approx 0.77$  and Spearman  $\rho \approx 0.82$ ,  
 153 indicating mild reordering of the ranking. The defence restored both correlations to 1.0.

154 To better understand how the adversarial edit affected individual rubric dimensions, Table 3 reports  
 155 the change in each score relative to baseline. The attack selectively inflated the novelty and signifi-  
 156 cance dimension of **P1** and **P2** by two points and, to a lesser extent, improved clarity by one point  
 157 as a side effect of minor edits. All other criteria remained unchanged, reflecting that hype language  
 158 primarily influences perceived novelty. Under the defence, scores for novelty and clarity returned to  
 159 their original values because the reviewer could not justify the increases with direct evidence from  
 160 the text.

161 We further analysed how the adversarial edits altered the distribution of average scores. Figure 1  
 162 shows histograms of the mean rubric scores under baseline, attack and defence. The attack distri-  
 163 bution shifts slightly to the right due to inflated novelty, while the defence distribution matches the  
 164 baseline. Such visualisations provide a fuller picture than binary accept/reject labels and emphasise  
 165 that adversarial cues can subtly inflate perceived quality without altering substantive content.

Table 3: Per-criterion changes due to adversarial editing (Attacked – Baseline). Positive values indicate the attacked abstract scored higher on that criterion. Under the defence, all scores reverted to their baseline values.

Paper	Method	Exp.	Novelty	Clarity	Repro.	Ethics
P1	0	0	+2	+1	0	0
P2	0	0	+2	+1	0	0
P3	0	0	0	0	0	0
P4	0	0	0	0	0	0

Figure 1: Distribution of mean rubric scores across the four abstracts under baseline (blue), attacked (orange) and defended (green) conditions. The attack increases the mean scores of P1 and P2, shifting the distribution rightwards. The defence restores the baseline distribution. (Illustrative figure; actual histograms will be included in the final submission.)

## 166 6 Granular Analysis of Rubric Dimensions

167 While aggregate metrics such as ASR and ranking correlations summarise overall effects, understanding which criteria are most susceptible to hype provides deeper insights. In this section we  
 168 examine each rubric dimension in turn, discuss its relevance to the four abstracts and highlight how  
 169 adversarial editing and the defence affected scores.

171 **Methodological soundness.** This criterion assesses whether the abstract clearly states objectives,  
 172 justifies assumptions and outlines a coherent approach. In our dataset, **P3** earned the highest method-  
 173 ological score because it concisely described the Transformer architecture and its advantages. The  
 174 attack did not alter methodological scores because hype words did not introduce additional method-  
 175 ological details. The defence similarly had no effect. This stability suggests that reviewers rely on  
 176 the presence of concrete methodological statements rather than rhetoric.

177 **Experimental adequacy.** Experimental adequacy measures whether empirical evaluation sup-  
 178 ports the claims. **P1** and **P2** mention empirical performance improvements in reinforcement learn-  
 179 ing, but the abstracts lack specifics about datasets, baselines or statistical analysis, resulting in mod-  
 180 erate scores. The attack did not significantly change these scores, reinforcing that hype cannot  
 181 compensate for missing experimental detail. Defences likewise had minimal impact.

182 **Novelty and significance.** Novelty assesses the originality and potential impact of the work. This  
 183 dimension proved most vulnerable: hype words inflated novelty scores for **P1** and **P2**, lifting them  
 184 into the acceptance region. The baseline moderate scores reflected genuine innovations (abstract  
 185 world models and dynamic state abstractions) but also indicated that the contributions may not be  
 186 groundbreaking. The defence neutralised the inflation by requiring concrete evidence for increased  
 187 novelty.

188 **Clarity and organisation.** This dimension captures writing quality. All four abstracts are pro-  
 189 fessionally written, but the adversarial edits slightly improved clarity for **P1** and **P2** because the  
 190 inserted phrases smoothed some sentences. The defence reset these scores to baseline since the im-  
 191 provements were not substantial enough to warrant a higher rating. This highlights how editorial  
 192 polish, even when rhetorical, can modestly influence clarity scores.

193 **Reproducibility and openness.** Reproducibility requires disclosure of datasets, code, or other ar-  
 194 tifacts. **P3** explicitly reports releasing source code and trained models, earning a high reproducibility  
 195 score. **P1** and **P2** mention improved sample efficiency but do not discuss code release, resulting in  
 196 low scores. **P4** falls in between, hinting at calibration strategies but lacking details about data or  
 197 code. The attack and defence left these scores unchanged, underscoring that rhetorical edits cannot  
 198 substitute for actual openness.

199 **Ethical and safety considerations.** This criterion evaluates whether the abstract acknowledges  
 200 potential risks and ethical implications. Only **P4** explicitly discusses fairness and calibration, which

201 naturally touches on ethics. The other abstracts do not mention ethics, leading to low scores. Adver-  
202 sarial editing did not introduce ethical considerations, and the defence could not raise scores without  
203 substantive content. This outcome suggests that embedding explicit ethical discussions into research  
204 communication is necessary for AI reviewers to recognise ethical soundness.

205 Overall, the granular analysis confirms that novelty and clarity are the most malleable dimensions  
206 under hype. Other criteria remain largely unaffected, indicating that targeted rhetorical edits selec-  
207 tively manipulate certain aspects of reviewer perception. Such insights can inform the design of  
208 rubrics and evaluation prompts to minimise susceptibility to superficial cues.

209 **7 Discussion**

210 Our findings illustrate how superficial hype can sway AI reviewers: modest increases in perceived  
211 novelty moved borderline works into the acceptance region, whereas strong papers such as **P3** re-  
212 mained unaffected. This pattern complements the results of Lin et al. (2025), who show that textual  
213 adversarial attacks can distort automated peer review. We observe that adding hype words influences  
214 only the novelty and clarity dimensions, leaving other criteria untouched; nonetheless, the induced  
215 flips underline a systemic vulnerability to superficial cues.

216 **Connections to cognitive bias and Moravec’s paradox.** The susceptibility to hype echoes  
217 Moravec’s paradox: AI models excel at formal reasoning yet lack the perceptual intuition that al-  
218 lows human reviewers to discount rhetorical embellishments. As Moravec noted, giving computers  
219 the skills of a one-year-old is harder than achieving grandmaster-level chess. In our experiment the  
220 models latched onto explicit markers of novelty but ignored the implicit absence of methodological  
221 details, demonstrating a cognitive bias toward overt signals. Addressing such biases may require  
222 integrating perceptual or commonsense reasoning components into LLM reviewers or combining  
223 them with human oversight.

224 **Implications for fairness and bias mitigation.** Our defence draws inspiration from the fairness  
225 literature, which emphasises rigorous evaluation and justification of decisions. Surveys of fairness  
226 research highlight that biases can emerge both during model training and during deployment. Ad-  
227 versarial overfitting in peer review can be viewed as a deployment-stage bias: reviewers misinterpret  
228 rhetorical cues as substantive novelty. Requiring evidence for each score serves as a form of extrin-  
229 sic debiasing, akin to prompting LLMs to justify outputs. However, this mechanism is only a first  
230 step; fairness research also advocates for diverse datasets, multiple evaluators, and statistical au-  
231 diting. Future AI review systems should incorporate these practices to ensure equitable and robust  
232 assessments.

233 **Broader impacts.** Beyond peer review, our findings highlight the risks of using LLMs in  
234 high-stakes evaluations. If minor edits can inflate scores in a scientific context, similar techniques  
235 might manipulate AI-based admissions, hiring or funding decisions. The fairness survey by Li *et al.*  
236 documents how biased LLM outputs can perpetuate stereotypes and discrimination. Transparent  
237 rubrics and evidence-based scoring may mitigate some vulnerabilities, but long-term solutions re-  
238 quire continual monitoring, open datasets for benchmarking, and collaboration between AI devel-  
239 opers and domain experts.

240 Requiring explicit evidence for each score proved to be an effective defence. This simple mechanism  
241 echoes strategies from bias elicitation work—Cantini et al. (2025) show that prompting models to  
242 justify their responses can improve robustness to adversarial bias probing. In our setting, the defence  
243 neutralised all flips by forcing the reviewer to ground scores in the text. Such evidence-based scoring  
244 could be incorporated into AI reviewing pipelines to mitigate overfitting to rubric keywords.

245 More broadly, our study aligns with the emerging literature on LLM security and prompt injection.  
246 Peng et al. (2024) review vulnerabilities in large language models and emphasise the need for com-  
247 prehensive safeguards against bias, misinformation and adversarial prompts. While our attacks are  
248 benign compared with malicious jailbreaks, they demonstrate how small edits can manipulate out-  
249 puts of an AI reviewer. Our results therefore contribute to the evidence base for designing safer,  
250 more transparent evaluation workflows.

251 There are several avenues for future research. First, scaling up experiments to dozens or hundreds of  
252 abstracts and multiple LLM reviewers would provide more statistical power and allow significance  
253 testing. Second, exploring richer adversarial strategies—such as adversarial prompt injection, hallu-  
254 cinated evidence or targeted obfuscation—could uncover additional vulnerabilities. Third, defences  
255 could be extended beyond evidence requirements to include consensus among multiple reviewers,  
256 adversarial training, or dynamic prompting that asks models to compare multiple candidate reviews.  
257 Finally, integrating human oversight and meta-evaluation (e.g., through meta-reviewers) may ensure  
258 that AI reviewers remain accountable and fair.

## 259 **8 Limitations**

260 This study has several limitations. (1) The dataset comprises only four abstracts, which limits sta-  
261 tistical power and generalisability; larger-scale experiments are needed to draw firm conclusions.  
262 (2) The rubric scores were assigned by a single AI reviewer configured with a fixed prompt, and  
263 the human authors verified them; using multiple models or prompt variations could yield different  
264 behaviours. (3) The adversarial edit targeted only novelty and impact; other attack surfaces (e.g.,  
265 prompt injections, obfuscation, hallucinated evidence) remain unexplored. (4) The defence required  
266 explicit evidence from the abstract but did not involve external verification; stronger defences might  
267 combine multiple reviewers or external fact-checking. (5) Because of the small scale, we did not re-  
268 port statistical significance or compute resource usage. These limitations constrain the conclusions  
269 and highlight the need for more comprehensive studies.

## 270 **9 Conclusion**

271 We conducted a miniature experimental reproduction of the Review-Overfitting Challenge to eval-  
272 uate how adversarial edits influence AI peer review and to test a simple defence based on evidence  
273 requirements. Our results show that adding hype-laden language can flip borderline decisions by  
274 inflating novelty scores, while strong papers remain unaffected. Requiring reviewers to ground their  
275 scores in the text neutralises this attack and restores original decisions. These findings align with  
276 recent studies that document vulnerabilities of LLM reviewers to textual manipulations and support  
277 calls for more rigorous evaluation protocols. We hope that this work spurs further research into  
278 adversarial robustness of AI reviewers and informs the design of secure, transparent peer-review  
279 pipelines.

## 280 **10 Responsible AI Statement**

281 Our research was carried out by an AI system with human oversight. The AI agent led hypothesis  
282 generation, experimental design and analysis, while the human collaborator reviewed the plan and  
283 ensured compliance with ethical guidelines. The study does not pose risks of harm, as it analy-  
284 ses publicly available abstracts and does not involve human subjects or sensitive data. The work  
285 highlights potential vulnerabilities in AI peer review and advocates for safeguards. We anticipate  
286 positive impacts through improving robustness of AI reviewers; however, misusing adversarial at-  
287 tacks to manipulate evaluations could have negative consequences. We recommend that conferences  
288 enforce evidence-based scoring and transparency.

## 289 **11 Reproducibility Statement**

290 All source materials are publicly accessible. The four abstracts were retrieved from arXiv using  
291 the identifiers provided in the references. The rubric criteria and scoring rules are described in  
292 Section 4. The A1 attack involved adding qualitative descriptors (<10% character change) without  
293 altering facts. The defence reset novelty scores that lacked supporting evidence. Our evaluation  
294 tables and computations (attack success rate and ranking correlations) are derived directly from the  
295 reported scores. Scripts and data will be released with the supplementary material.

296 **Agents4Science AI Involvement Checklist**

297 **1. Hypothesis development:**

298 Answer: [B]

299 Explanation: The AI system generated the research question and designed the simplified  
300 reproduction of the Review-Overfitting Challenge. A human overseer provided high-level  
301 guidance and approved the approach.

302 **2. Experimental design and implementation:**

303 Answer: [C]

304 Explanation: The AI agent selected the abstracts, defined the rubric and attack, computed  
305 metrics and produced tables. The human collaborator verified the experimental pipeline.

306 **3. Analysis of data and interpretation of results:**

307 Answer: [B]

308 Explanation: The AI calculated the attack success rate and ranking correlations and inter-  
309 preted the results. The human checked that the interpretations aligned with the data.

310 **4. Writing:**

311 Answer: [C]

312 Explanation: The AI drafted the manuscript, including abstract, sections and checklists.  
313 The human reviewer ensured anonymity and adherence to conference guidelines.

314 **5. Observed AI Limitations:**

315 Description: The AI was unable to conduct large-scale experiments or call proprietary LLM  
316 APIs, restricting the dataset to four abstracts. It relied on heuristics for scoring and required  
317 human confirmation to ensure ethical compliance.

318 **Agents4Science Paper Checklist**

319 **1. Claims**

320 Answer: [Yes]

321 Justification: The abstract and introduction clearly state the goal of reproducing a simpli-  
322 fied review-overfitting experiment and accurately reflect the contributions presented in the  
323 paper.

324 **2. Limitations**

325 Answer: [Yes]

326 Justification: Section 8 explicitly discusses limitations regarding the dataset size, scoring  
327 methodology, simplified attack/defence and lack of statistical tests.

328 **3. Theory assumptions and proofs**

329 Answer: [NA]

330 Justification: The paper does not present theoretical results or proofs; it is an empirical case  
331 study.

332 **4. Experimental result reproducibility**

333 Answer: [Yes]

334 Justification: Section 11 provides sufficient details for reproducing the main results: dataset  
335 identifiers, scoring procedure, attack and defence definitions, and metrics.

336 **5. Open access to data and code**

337 Answer: [Yes]

338 Justification: The abstracts are available on arXiv and all scripts and data will be released  
339 as supplementary material.

340 **6. Experimental setting/details**

341 Answer: [Yes]

342 Justification: Section 4 describes the dataset, scoring rubric, attack method and defence  
343 procedure, which are sufficient to understand the results.

344 **7. Experiment statistical significance**

345 Answer: [NA]

346 Justification: Due to the small illustrative dataset, statistical significance tests and error bars  
347 are not applicable.

- 348       **8. Experiments compute resources**  
 349       Answer: [NA]  
 350       Justification: The experiments involved simple scoring and metrics computed on a small  
 351       dataset; specific compute details are unnecessary.
- 352       **9. Code of ethics**  
 353       Answer: [Yes]  
 354       Justification: The study adheres to the Agents4Science Code of Ethics, uses public data,  
 355       respects privacy and discusses societal impacts.
- 356       **10. Broader impacts**  
 357       Answer: [Yes]  
 358       Justification: The Responsible AI Statement and Discussion section reflect on positive and  
 359       negative societal impacts, including risks of adversarial manipulation and benefits of robust  
 360       peer review.

## 361       **References**

- 362       Rafael Rodriguez-Sanchez and George Konidaris. *Learning abstract world model for value-*  
 363       *preserving planning with options*. arXiv preprint arXiv:2406.15850, 2024.
- 364       Mehdi Dadvar, Rashmeet Kaur Nayyar and Siddharth Srivastava. *Learning dynamic abstract repre-*  
 365       *sentations for sample-efficient reinforcement learning*. arXiv preprint arXiv:2210.01955, 2022.
- 366       Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
 367       Łukasz Kaiser and Illia Polosukhin. *Attention Is All You Need*. In Advances in Neural Information  
 368       Processing Systems, 2017.
- 369       Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu and  
 370       Tianyu Liu. *Large Language Models are not Fair Evaluators*. arXiv preprint arXiv:2305.17926,  
 371       2023.
- 372       Tzu-Ling Lin, Wei-Chih Chen, Teng-Fang Hsiao, Hou-I Liu, Ya-Hsin Yeh, Yu Kai Chan, Wen-  
 373       Sheng Lien, Po-Yen Kuo, Philip S. Yu, and Hong-Han Shuai. Breaking the Reviewer: Assessing  
 374       the Vulnerability of Large Language Models in Automated Peer Review Under Textual Adversar-  
 375       ial Attacks. *arXiv preprint arXiv:2506.11113*, 2025.
- 376       Riccardo Cantini, Alessio Orsino, Massimo Ruggiero, and Domenico Talia. Benchmarking Adver-  
 377       sarial Robustness to Bias Elicitation in Large Language Models: Scalable Automated Assessment  
 378       with LLM-as-a-Judge. *arXiv preprint arXiv:2504.07887*, 2025.
- 379       Benji Peng, Keyu Chen, Ming Li, Pohsun Feng, Ziqian Bi, Junyu Liu and Qian Niu. Securing  
 380       Large Language Models: Addressing Bias, Misinformation, and Prompt Attacks. *arXiv preprint*  
 381       *arXiv:2409.08087*, 2024.
- 382       Agents4Science Committee. *Agents4Science 2025 Call for Papers*. Open Conference of AI Agents  
 383       for Science. Retrieved from <https://agents4science.stanford.edu/call-for-papers.html>, 2025.
- 384       Hans Moravec. *Mind Children: The Future of Robot and Human Intelligence*. Harvard University  
 385       Press, 1988. Moravec articulates that computers excel at symbolic reasoning but struggle with sen-  
 386       sory and motor skills, a phenomenon later dubbed Moravec's paradox[751848279319933†L150-  
 387       L169].
- 388       Yingji Li, Mengnan Du, Rui Song, Xin Wang, Ying Wang and colleagues. A Sur-  
 389       vey on Fairness in Large Language Models: Evaluation and Debiasing Methods. *arXiv*  
 390       *preprint arXiv:2308.10149v2*, 2024. The survey reviews intrinsic and extrinsic bias met-  
 391       rics and debiasing strategies for medium- and large-scale LLMs[332816400343247†L95-  
 392       L109][332816400343247†L130-L144].
- 393       Agents4Science Committee. Call for Papers: Open Conference of AI Agents for Science 2025.  
 394       Available at <https://agents4science.stanford.edu/call-for-papers.html>, 2025.

395 Hans Moravec. *Mind Children: The Future of Robot and Human Intelligence*. Harvard University  
396 Press, 1988. Moravec articulated that it is easy to build computers that excel at logic and games  
397 but difficult to endow them with the perceptual and motor skills of a human infant.

398 Yingji Li, Mengnan Du, Rui Song, Xin Wang and Ying Wang. A Survey on Fairness in Large  
399 Language Models. *arXiv preprint arXiv:2308.10149v2*, 2024. The survey reviews intrinsic and  
400 extrinsic bias evaluation metrics and debiasing techniques for medium- and large-sized LLMs.