



PaperDoctor: Evidence-Grounded and Actionable Feedback for Scientific Papers in Progress

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Autoresearch agents are reshaping the research pipeline, but they also let flawed
2 claims enter the literature at scale. Human advisors catch such issues on in-progress
3 drafts through careful, traceable feedback, yet advisor-style assess requires exten-
4 sive manual effort and does not scale. To shift automated paper assessment from a
5 judge to a diagnostician, we introduce PaperDoctor, an agent framework for pre-
6 submission feedback with three key innovations: **(i) Holistic hierarchical pipeline.**
7 Every paper is assessed across writing, layout, references, code, theory, prior work,
8 and experiments through three layers: L1 paper-only screening runs cheaply on
9 every submission; L2 typed verifiers route each claim to the skill that owns its evi-
10 dence; and L3 reproducers rerun experiments by priority. **(ii) Evidence-grounded**
11 **actionable feedback.** Each finding is a triple of an observation (Why), a pointer to
12 a specific location such as a sentence, equation, or code line (Where), and a revision
13 suggestion (How), making every critique auditable and actionable. **(iii) Effective**
14 **experimental reproduction.** Beyond reading the paper, PaperDoctor selectively
15 rebuilds and reruns experiments based on claim importance and compute budget,
16 surfacing reproducibility gaps and quantitative limitations that are invisible from
17 the manuscript alone. We evaluate PaperDoctor on human studies which are junior
18 students on their pre-submission papers, yield 85% accuracy, notably high in claim
19 assumption validation. We further evaluate PaperDoctor on 60 manuscripts across
20 machine learning, natural and social-sciences, covering human- and agent-author
21 papers with code. Overall, PaperDoctor produces more grounded feedback than
22 human and LLM reviewers, its findings track paper quality, and its reproduction
23 stage surfaces limitations missed in the paper’s main body. To empower the com-
24 munity, we develop a demo interface that lets authors browse findings grounded
25 on their paper. PaperDoctor reframes automated paper assessment as a diagnostic
26 process rather than a verdict, taking a concrete step toward AI advisors that help
27 authors raise the quality of scientific writing in the autoresearch era.

28 1 Introduction

29 “Don’t find fault, find a remedy.” — Henry Ford

30 Autoresearch agents that turn a proposed idea into a full manuscript that [22, 45, 28, 34, 2] increases
31 the risk of producing work whose quality cannot be guaranteed, such as these claims are often
32 difficulty to verify. Moreover, most human-written drafts are now at least partly AI-assisted, whether
33 in writing, figure drawing, or literature survey. This trend raises paper volume while leaving draft
34 quality uncontrolled. Existing automated reviewers [8, 18, 12, 13, 49, 32] do not close this gap.
35 By returning a decision (“accept” or “reject”) with a justification, they primarily serve to filter
36 submissions, which is useful for the reviewer’s side but uninformative for the author.

37 Human advisors catch exactly these issues during discussion with students [4, 31]. They go through
38 the draft carefully: circling a claim with no experiments validation, underlining a theorem whose

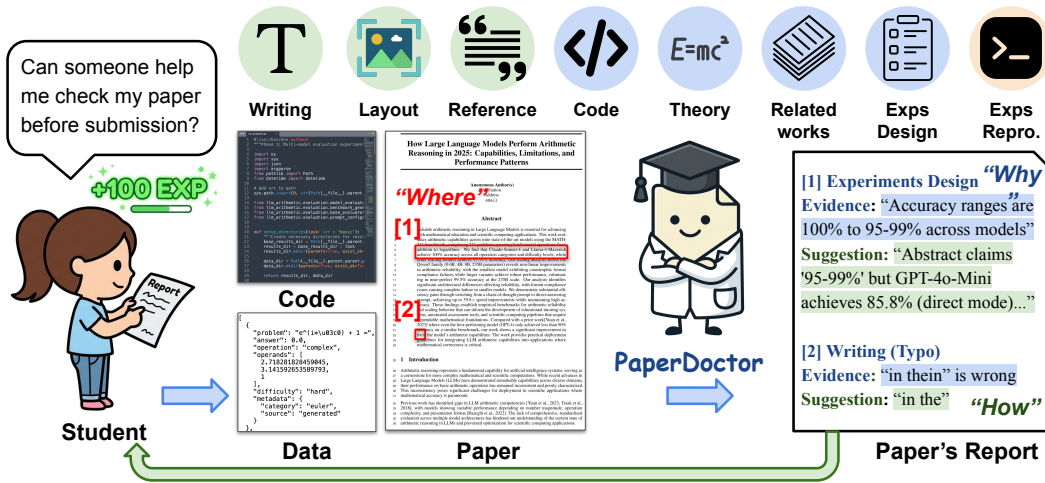


Figure 1: **PaperDoctor Agent provides evidence-grounded, actionable feedback.** Given a paper along with its code and datasets, PaperDoctor returns a holistic report of findings across multiple dimensions (writing, citations, figures, equations, code, reproduction, and related work). Each finding pairs an error type (*Why*) with a concrete suggestion (*How*) and is grounded (*Where*) in the paper itself, allowing authors to audit and learn from every critique.

39 assumption breaks, or flagging a figure that is unclear. Each mark on the page is a small diagnosis. As
 40 illustrated in Figure 1, it points to *where* the symptom is in the paper, explains *why* it is a problem, and
 41 prescribes *how* to fix it; in short, it is **evidence-grounded** and **actionable**. This is how students learn
 42 from an advisor’s diagnosis on the page, which is fundamentally different from a reviewer’s judgment.
 43 Yet such depth comes at a cost: it takes hours per paper, making it impossible to keep pace with
 44 agent-assisted writing [22, 45, 28, 34]. Current agents either write the paper for the author [34, 2]
 45 or grade it at the door [21, 27, 18, 33, 20]. Neither is what a student receives from an advisor: a
 46 careful diagnosis, like that of a doctor, that points to a specific section and says what to change.
 47 However, delivering this level of feedback is non-trivial. A paper is a multi-faceted artifact spanning
 48 presentation, implementation, experimental analysis, and more. This is why authors typically rely on
 49 feedback from a range of people, such as advisors, peers, and industry mentors, each catching issues
 50 the others miss.

51 Motivated by this, we ask: *can an agent provide diagnosis-level feedback at scale?* We introduce
 52 PaperDoctor, a research agent framework that delivers constructive feedback to human authors.
 53 Notably, PaperDoctor highlights the following: (i) **Holistic, hierarchical pipeline.** PaperDoctor
 54 decomposes a paper, together with its code and datasets, along dimensions ranging from writing to
 55 experimental reproduction, and organizes the assessment into three hierarchical levels of increasing
 56 cost and subjectivity. **L1 (surface screening)** runs cheaply on every submission and targets objective
 57 issues in writing, layout, references, and figures. **L2 (claim verification)** extracts atomic claims
 58 and routes each one to the skill that owns its evidence—web search for prior-work checks, a vision
 59 language model (VLM) for figure assessment, code analysis for implementation claims, and theory
 60 verifiers for derivations. **L3 (experimental reproduction)** selectively execute experiments, validating
 61 reproducibility and correctness at execution time. This design spends effort proportional to the cost
 62 of verification and keeps the full pipeline tractable. (ii) **Evidence-grounded, actionable feedback.**
 63 Each finding is a triple of an observation with a reason (*Why*), a pointer to a specific location in the
 64 paper (*Where*, *i.e.*, a sentence, equation, code line, or external URL), and a concrete suggestion for
 65 revision (*How*). This makes every critique both auditable and directly actionable for human authors.
 66 (iii) **Effective experimental reproduction.** Beyond reading the manuscript, PaperDoctor selectively
 67 prioritize and reruns experiments based on claim importance and compute budget. By executing code
 68 rather than only reading it, PaperDoctor surfaces reproducibility gaps and quantitative limitations that
 69 are invisible from the manuscript alone.

70 To validate the effectiveness of PaperDoctor, we first conduct a human study with junior student
 71 participants, collecting their in-progress paper drafts and asking them to rate the feedback produced
 72 by PaperDoctor. PaperDoctor reaches 85% agreement with their ratings. Moreover, we evaluate
 73 PaperDoctor on 60 manuscripts covering machine learning, the natural sciences, and the social
 74 sciences, including both human- and agent-written papers with accompanying code and data. This
 75 diverse benchmark allows us to systematically assess PaperDoctor’s robustness across disciplines,

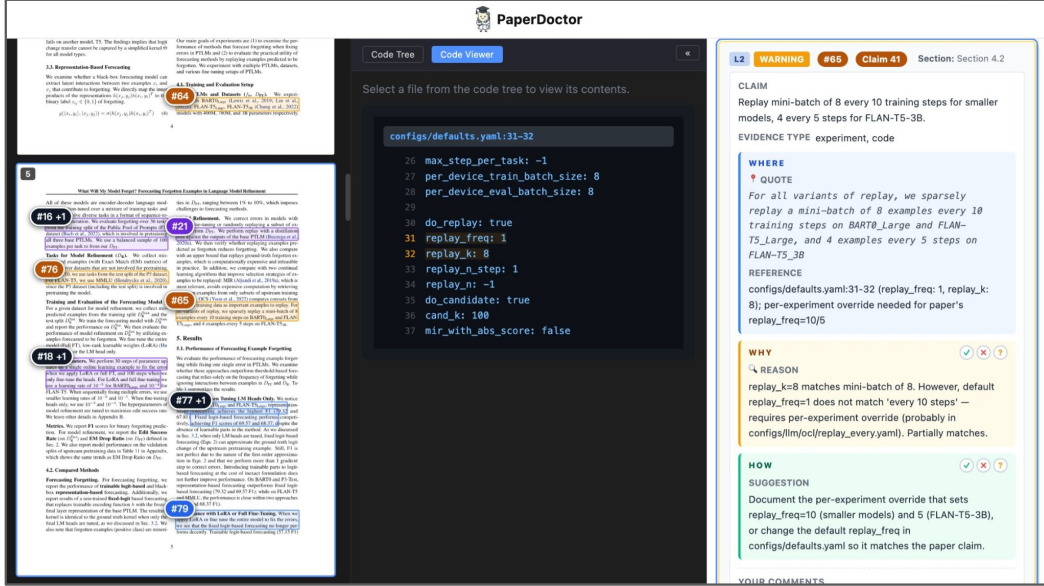


Figure 2: **The PaperDoctor Diagnosis Interface.** Authors upload their paper and code, and PaperDoctor returns the holistic report. **Left:** the paper, with each finding overlaid as a numbered, color-coded bubble anchored to the exact span it critiques (Where). **Middle:** the code viewer, so findings about implementation can be audited against the source. **Right:** the finding card for the selected highlight, showing the observation (Why) and a concrete suggestion for revision (How). In the example shown, PaperDoctor identifies an implementation-level mismatch: the paper claims a replay mini-batch of 8 every 10 training steps, but the default config (`configs/defaults.yaml:31-32`) sets `replay_freq=1`.

76 writing styles, and varying levels of methodological rigor. We found that (i) PaperDoctor produces
 77 more grounded feedback than human and LLM reviewers, its findings track paper quality, and its
 78 full coverage surfaces limitations missed from the paper main body. In particular, by executing the
 79 accompanying code and re-running key experiments, PaperDoctor uncovers discrepancies between
 80 reported and actual results that are otherwise difficult to detect through reading-only review.

81 Lastly, to empower the research community, we release a demo interface that lets authors browse
 82 findings anchored directly on their own papers (Fig. 2), making it easy to trace each critique back to
 83 its exact location.

84 2 PaperDoctor

85 2.1 Overview

86 Given a paper and its code, PaperDoctor aims to produce a set of findings $\{\mathcal{F}_i\}_{i=1}^N$, where each
 87 finding is defined as

$$88 \mathcal{F}_i = (f_i, e_i, s_i). \quad (1)$$

88 Here, f_i is the finding itself, e_i is the *evidence* grounding it to a specific location (a sentence,
 89 code line, or external URL), and s_i is a concrete *suggestion* for revision. The pair (e_i, s_i) lets the
 90 author audit and act on each finding without searching through the full paper. Notably, we distinguish
 91 findings into two severity levels. An *error* indicates that PaperDoctor is confident the issue is incorrect,
 92 while a *warning* indicates uncertainty and flags the finding for further clarification, such as a human
 93 check.

94 **Paper-Code Parsing.** Producing $\{\mathcal{F}_i\}$ by feeding the entire paper and codebase into a single model
 95 is infeasible: a paper is a long, multimodal document and a codebase is itself a large, structured
 96 artifact. Moreover, most downstream skills only need a targeted slice of the inputs (*i.e.*, reference
 97 verification needs only the bibliography; figure assessment needs only the rendered pages). We
 98 therefore run a single parsing step that produces three decomposed reusable artifacts:

$$(\mathcal{P}_t, \mathcal{P}_v, \mathcal{C}_t, \mathcal{B}) \leftarrow (\text{Paper}, \text{Code}), \quad (2)$$

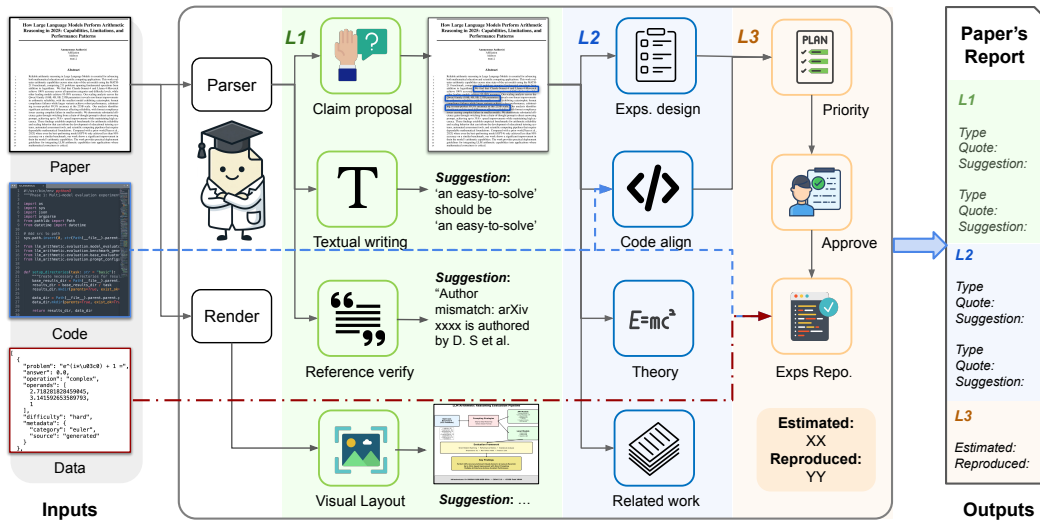


Figure 3: **Illustration of PaperDoctor pipeline.** A preparation step parses and renders the paper and code into a structured format and page images. (i) **The L1 surface-screening** stage runs four paper-only skills in parallel, including claim extraction. (ii) **The L2 claim-verification** stage dispatches each claim to the verifiers matching its evidence type and runs them in parallel. (iii) **The L3 experiment-reproduction** stage executes a prioritised, human-approved subset of the reproduction plan.

99 where \mathcal{P}_t is the paper as section-organized markdown (via Mathpix¹), \mathcal{P}_v is the same paper rendered
 100 page-by-page as images for downstream Vision-Language Model (VLM) use, and \mathcal{C}_t is the code
 101 indexed with tree-sitter² into per-file units. \mathcal{B} denotes the parsed bibliography of the paper (*i.e.*, .bib
 102 file). Every downstream skill reads from this shared representation and requests only the section,
 103 page image, or code snippet it needs.

104 **Hierarchical Pipeline.** A paper spans many dimensions, and processing all of them in a single pass
 105 is infeasible. Different aspects demand different forms of evaluation, and these evaluations vary
 106 widely in cost: a writing check is a single LLM call, whereas experiment reproduction can consume
 107 hours of GPU execution. We therefore organize PaperDoctor as a hierarchical pipeline of three levels
 108 (L1–L3) that spends effort proportional to the cost of verification. As illustrated in Figure 3, L1
 109 handles the most concrete, surface-level checks (such as presentation); L2 verifies individual claims
 110 along specific dimensions (such as theory or comparison with prior work); and L3 runs the most
 111 expensive stage, full experimental reproduction.

112 2.2 L1 – Surface Screening

113 This stage focuses on straightforward issues that can be easily addressed by browsing the paper.

114 **Writing Review.** We ask an LLM to read the paper section by section \mathcal{P}_t and flag writing issues
 115 as it goes. Clear mistakes such as typos or grammatical errors are marked as *errors*, since they are
 116 unambiguously wrong. Stylistic issues, where the text is understandable but could be phrased more
 117 clearly, are marked as *suggestions* instead, leaving the final decision to the author. For every issue,
 118 we require the LLM to quote the original sentence verbatim as evidence, ensuring each finding can
 119 be traced back to a specific location in the paper.

L1: Writing

Evidence: Page 3 “We trian the model on a large corpus of academic papers.”
Suggestion: There is a typo: “trian” should be “train”.

120
 121 **Figure Review.** A paper is as much a visual artifact as a textual one: its figures, tables, and overall
 122 layout are carefully curated by the authors and judged by readers at a glance. Yet most of this visual
 123 information is lost in markdown extraction. We therefore treat visual inspection as a separate check
 124 in PaperDoctor: we render the paper into page images \mathcal{P}_v and ask a VLM to review them directly.
 125 Clear visual defects, such as figures overflowing the text margin or overlapping captions, are flagged
 126 as *errors*. More subjective issues, such as undersized fonts or insufficient color contrast, are flagged

¹<https://mathpix.com/>

²<https://github.com/tree-sitter/tree-sitter>

127 as *warnings* for the author to judge. As with writing review, every finding must be grounded to a
128 specific page or figure index as evidence.

L1: Figure

Evidence: Page 5, Figure 3 extends beyond the right text margin.

Suggestion: Rescale the figure width to `\linewidth`.

129

130 **Citation Check.** AI-assisted manuscripts routinely contain references that do not resolve to any real
131 paper, and this is tedious to catch by reading the bibliography alone. Conditioned on the \mathcal{B} , we pair
132 the agent with a web search backend: for each reference, the LLM issues up to three search queries
133 and records a resolver URL only when the backend returns a genuine match, never synthesizing itself.

L1: Citation Check

Evidence: Reference [12] “Smith et al., Neural Reasoning in Transformers, NeurIPS 2023” returns no
match.

Suggestion: The reference appears to be hallucinated. Please verify and replace with a valid source.

134

135 **Claim Extraction.** A paper makes dozens of arguments across its sections—novelty claims in the
136 introduction, methodological choices in the method section, performance numbers in the experiments,
137 and so on—and the value of PaperDoctor comes from checking each of them against its own evidence.
138 A claim that is never extracted can never be verified. We therefore ask an LLM to densely extract
139 every verifiable assertion the authors make, covering [theory, code, experiments(designs),
140 literature]; a single claim may be tagged with multiple evidence types. Each claim is then
141 dispatched to the corresponding L2 branches for verification.

L1: Claim Extraction

Claim: “Our method achieves 92.3% accuracy on ImageNet, outperforming prior state-of-the-art by 3.1
points.”

Evidence: Introduction

Claim Type: [Experiment, Related Work]

142

143 Owing to the modular design, all four skills in L1 can run *in parallel*.

144 2.3 L2 – Claim Verification

145 In this stage, each claim by L1 claim extraction is sent to every verifier for further verification.

146 **Code Verification.** A common failure mode of AI-assisted drafts is that the described optimizer,
147 architecture, or training setup does not match the released code. These mismatches are almost
148 invisible to human reviewers, who rarely open the repo while reading. We therefore ask PaperDoctor
149 to check each code-tagged claim directly against the source ($\mathcal{P}_t, \{\mathcal{C}_j\}$). For example, hyperparameter
150 claims are often best verified by first inspecting configuration files before descending into the Python
151 implementation: if the paper claims AdamW but the config specifies Adam, we flag a *warning*—the
152 mismatch is real, but may be a stale config or a last-minute switch the author should confirm. If a
153 component described in the paper is missing from the code altogether, we flag it as an *error*.

L2: Code Verification

Claim: “We train all models using the AdamW optimizer.”

Evidence: `configs/train.yaml` line 14 specifies `optimizer: Adam`, which is conflict with the paper
experiment settings

Suggestion: Mismatch between paper and code. Verify which optimizer was actually used.

154

155 **Theory Verification.** Errors in theoretical derivations are among the hardest to catch: a proof that
156 reads smoothly often hides missing assumptions, skipped steps, or notation drift across equations, and
157 even careful readers can miss these on a first pass. We therefore ask PaperDoctor to re-derive each
158 argument in \mathcal{P}_t step by step rather than summarize it. PaperDoctor jointly examines all theoretical
159 content, including equations, variables, and the notation that links them across the paper, and records
160 the full trace so that a superficial check is itself visible as a superficial trace. Specifically, PaperDoctor
161 examines four aspects in turn: correctness of each step, hidden assumptions that the paper does not
162 state, boundary or edge-case behavior, and notation consistency across derivations. A loss whose
163 expectation silently drops between its definition and its final form is caught here, before it propagates
164 into code or experiments.

L2: Theory Verification

Claim: “The expected loss reduces to $\mathbb{E}[\|x - \hat{x}\|^2]$ (Eq. 7).”

Evidence: Step from Eq. 6 to Eq. 7 drops the cross-term $\mathbb{E}[x^\top \hat{x}]$ without justification.

Suggestion: Missing assumption that x and \hat{x} are uncorrelated. Please state explicitly or correct the derivation.

165

166

167

168

169

170

171

172

173

Literature Check. A common issue in scientific drafts is overstated novelty or weak engagement with the literature (“*no prior work addresses X*” when an earlier paper already does), and this bias is easier to catch by searching than by reading. Based on $(\mathcal{P}_t, \mathcal{B})$, PaperDoctor pairs an LLM with a web search backend to separate baseline comparisons, cited facts, and novelty assertions, so that each novelty assertion can be further labelled as `novel`, `incremental`, or `prior_art_exists`.

Note that this differs from the reference verifier in L1: the reference verifier only checks whether a cited paper exists and is correctly attributed, while this stage examines whether the surrounding literature *content* actually supports the paper’s novelty and positioning claims.

L2: Literature Check

Claim: “No prior work addresses multi-modal reasoning over long video sequences.”

Evidence: Web search returns Chen et al. (2023), “LongVid-Reasoner”, which targets the same setting.

Suggestion: Novelty overstated. Relabel as `prior_art_exists` and cite Chen et al.

174

175

176

177

178

179

180

181

182

183

184

185

Experiment Design. While some claims can be closed-loop verified against external sources (literature) or formal content (theory, code), most claims in a paper rest on *experimental evidence* and crucially, whether the experiments themselves are well-designed determines whether the contribution is genuinely grounded. Experimental issues thus fall into two regimes: design mistakes (missing ablations, missing experiments for a stated contribution) that can be found from the paper alone, and reproduction mistakes that can only be found by running. This module handles the first, before any execution. Agent reviews whether each argument is supported by a corresponding experiment, along with fairness, ablation sufficiency, statistical rigour, baseline recency, and cherry-picking risk. It then emits a `reproduction-plan` entry per experiment listing the command, priority, feasibility, run mode (evaluation or training), and the numeric target lifted from the paper. No experiments run here; the plan is a declarative contract for L3.

L2: Experiment Design

Claim: “Our cross-modal attention module is the key component driving the gains over prior work.”

Evidence: Table 3 reports only the full method versus the baseline; no ablation removes the cross-modal attention module to isolate its contribution.

Suggestion: Add an ablation that disables the cross-modal attention module and reports performance on the same benchmark.

Reproduction plan: `bash eval/ablate_attn.sh, target $\Delta \geq 1.0$ point drop, priority high.`

186

187

188

Notably, the L2 stage remains highly modular: all four verifiers, as well as the per-claim dispatch within each verifier, run *in parallel*.

189

190

191

192

193

194

2.4 L3 – Experiments Reproduction

After L1 and L2, PaperDoctor has already covered most aspects that can be assessed from the paper’s main body. The remaining equally important is reproduction, which is challenging yet essential. Different papers come with very different experimental setups: some require only inference, others involve full training, and many depend on substantial resources such as GPU compute, datasets, and storage. We therefore design a separate L3 stage dedicated to reproduction.

195

196

197

198

199

200

201

Priority Ordering. It is worth noting that not every claim deserves an equal degree of attention. For example, an experiment that backs a main claim in the abstract carries far greater weight than a hyper-parameter sensitivity study, even though the latter may be cheaper to run. We rank experiments by their importance to the paper’s central contributions and by expected *feasibility check*. PaperDoctor assigns each entry in the reproduction plan a priority label of `{high, medium, low}`. Under a compute budget, high-priority items run before medium and low, evaluations before training, and ready experiments before blocked ones.

202

203

204

Manual Approval. Reproduction consumes real compute and storage, and executes code that may affect the environment. A mis-typed claim could silently trigger a multi-hour training run. PaperDoctor therefore presents the L2 plan, annotated with estimated GPU hours, dataset size,

205 and storage footprint, for the author to approve. Only then does an LLM-driven executor handle
206 environment setup, dispatch, and log parsing.

L3: Reproduction

Claim: “Our method achieves 78.4% accuracy on MMLU.”

Evidence: Reproduced run yields 71.2% (logs/mmlu_eval.log), a 7.2-point gap.

Suggestion: Discrepancy exceeds the 1–2% tolerance. Flagged as *error*; verify evaluation protocol or reported number.

207
208 **Report Results.** Each executed experiment will receive a decision by comparing the reproduced
209 value against the paper’s reported one. Rather than imposing a fixed numeric threshold, PaperDoctor
210 judges the verdict in context: it considers the metric type, the typical variance reported in the paper,
211 and the magnitude of the original gap, and decides whether the result counts as a *pass*, a *warning*
212 (partial match), or an *error* (OOM, divergence, or substantial mismatch). This gives the author a
213 direct view of which claims hold up under execution and which diverge from what the paper reports.

214 3 Experiments

215 We design experiments to answer four research questions. **Q1.** How do human (such as students)
216 perceive PaperDoctor’s feedback? **Q2.** How does PaperDoctor perform across papers of varying
217 sources and quality? **Q3.** How do PaperDoctor’s findings compare to those of human and LLM
218 reviewers? **Q4.** What unique findings does PaperDoctor uncover during experiment reproduction?

219 3.1 Settings

220 We conduct human studies with seven PhD students, collecting their pre-submission paper drafts.
221 For the PaperDoctor generated report, we let these authors write and label every finding as *accepted*,
222 *rejected*, or *uncertain*.

223 Furthermore, to study the performance across various paper types, we collect 60 papers from four
224 diverse sources (as summarized in Tab. 1): twenty come from Agents4Science [2], half of them
225 accepted as orals or spotlights and half rejected, so we could ablate AI-written outcome. Another
226 twenty are from PaperBench [30]: human-written papers at top AI venues (ICML), all with released
227 code. The remaining twenty are human-written papers from Nature Communications and Nature
228 Human Behaviour (ten each), which pulls the evaluation outside machine learning and into the natural
229 and social sciences.

Table 1: Papers used in PaperDoctor evaluation.

Author	Source	Domains	Size	Status	Code
AI	Agents4Science	AI & Science	20	Oral/Spot. & Reject	✓
Humans	PaperBench	AI (ICML)	20	Oral/Spotlight	✓
Humans	Nature Communications	Natural Science	10	Public	✓
Humans	Nature Human Behaviour	Social Science	10	Public	✓

230 We use Claude-opus-4.7 as the base model and implement PaperDoctor as skills.

231 3.2 How do author assess PaperDoctor’s feedback?

232 As shown in Fig. 5, we use PaperDoctor to generate 349 findings across 7 pre-submission papers,
233 and ask each author to label every finding on their own paper. Authors accept 85% of the findings on
234 average.

235 **Where do authors agree most?** As shown in Fig. 5, authors accept the majority of PaperDoctor’s
236 findings across all six dimensions, with acceptance rates ranging from 55% (figure) to 88% (theory)
237 without rejection. L2 dimensions, which target specific claim types such as theory and related work,
238 achieve consistently higher acceptance (77–88%) than L1 surface-level checks (55–74%), owing to
239 their grounding in external evidence (*i.e.*, web search) and full-paper consistency over equations and
240 notation. Among them, Experiment review yields the largest number of findings per paper, providing
241 authors with a systematic reminder to revisit their experimental design.

242 **Where do authors reject most?** Figure dimension is the most rejected dimension, with nearly a
243 third of findings marked as Rejected. Two reasons stand out: (i) current vision models still struggle
244 to parse complex scientific figures and tend to hallucinate, and (ii) figure design is highly subjective,
245 so automated judgments often diverge from what authors actually intend.

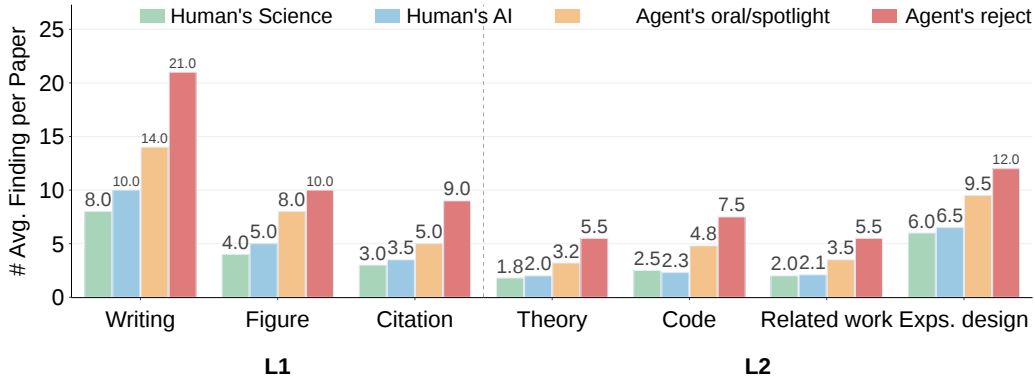


Figure 4: **Average findings per paper across paper categories.** We split papers into four groups by author and quality: human-authored science papers (Nature Communication, Nature Human Behavior), human-authored ML papers (PaperBench), and agent-written papers from Agents4Science, separated into oral/spotlight and rejected.

246 3.3 Are PaperDoctor’s findings correlated with the paper category and quality?

247 **Is PaperDoctor share the same role as Humans?** Figure 4 breaks PaperDoctor’s findings down
 248 by skill across four paper groups: human-written Nature papers, human-written AI papers from
 249 PaperBench, and agent-written papers that were either accepted (oral/spotlight) or rejected at
 250 Agents4Science. Several observations stand out.

251 **Human-written papers have fewer sugges-**
 252 **tions than Agent’s papers.** Across every di-
 253 mension, agent papers get more findings than
 254 human ones, suggesting PaperDoctor’s finding
 255 counts track paper quality. Among human pa-
 256 pers, Nature submissions sit at the low end on
 257 most dimensions, likely because their stricter re-
 258 view process catches surface issues earlier than
 259 open conference review. **Agent rejected pa-**
 260 **pers have more suggestions than Agent top**
 261 **papers.** The gap is consistent across both L1
 262 and L2 dimensions. Within L2, Experiment
 263 Design draws the most, ahead of theory, code,
 264 and related work. Many of these come from
 265 PaperDoctor checking the paper’s internal self-
 266 consistency, for instance whether a number in
 267 the table matches the text, or whether the re-
 268 ported metric matches the one defined in the
 269 method.

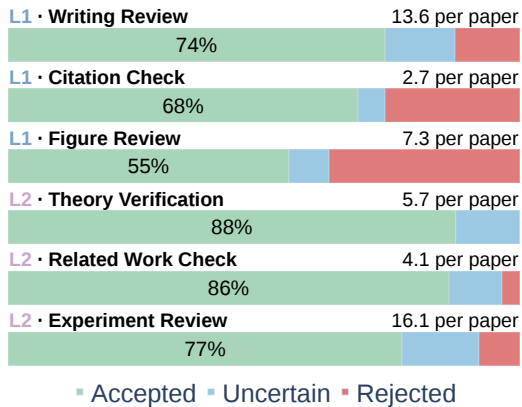


Figure 5: Author votings (Accepted / Uncertain / Rejected) on PaperDoctor’s findings, broken down by individual modules.

270 Overall, PaperDoctor’s finding counts align well with paper category and quality. Publicly accepted
 271 human papers carry the fewest issues, reflecting rounds of polishing, while agent-written papers,
 272 especially rejected ones, leave many issues a human collaborator would normally catch.

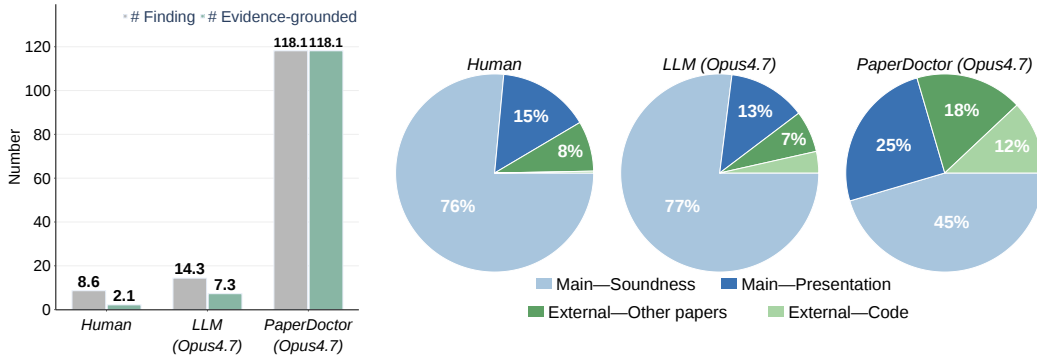
273 3.4 Comparison with Human Reviewers and LLM Reviewers

274 We study how human reviewers and LLM reviewers behave on the PaperBench papers [30] and obtain
 275 reviews directly from their authors. These reviews are not publicly available, ensuring that no current
 276 LLM has seen them during training.

277 **How well do humans and LLMs ground their findings?** We first ask whether human and LLM
 278 reviewers actually ground their findings in the paper, using an LLM judge to label each finding as
 279 evidence-grounded or not. We report results in Fig. 6a. We observe two gaps between baseline
 280 reviewers and PaperDoctor.

281 **(i) Volume:** Human reviewers produce only 8.6 findings per paper on average, and an LLM reviewer
 282 (Opus 4.7) with full access to the paper writes 14.3, far below PaperDoctor’s 118.1. The gap is by
 283 design: rather than taking each argument for granted, PaperDoctor surfaces every claim that could
 284 potentially require verification.

285 **(ii) Evidence grounding:** Only 25% of human findings and 51% of LLM findings cite a concrete
 286 locus (a figure, table, equation, or quoted line), whereas PaperDoctor grounds every one of them
 287 (100%). These gaps are not accidental. For human reviewers, writing a precise locus during review is
 288 costly: it requires flipping through the paper, copying the exact span, and cross-checking line numbers,
 289 all competing with limited reviewing time. LLM reviewers fail differently: with the full paper in
 290 context, they treat every comment as self-evident thus omitting explicit citations that triggered the
 291 issue.



(a) **Finding number** (with and without evidence grounding), compared across human reviewers, an LLM reviewer, and PaperDoctor. (b) **Distribution of finding types.** *Main* refers to findings grounded in the paper body, *External* to those grounded in supplementary material, code, or other literature. *Soundness* covers method, theory, and experiments; *Presentation* covers writing and figures.

292 **What do Humans, LLMs, and PaperDoctor focus on?** In figure 6b, we classified the review
 293 content into four buckets: soundness and presentation within the paper body, and external checks
 294 against other literatures or against the released code. Human reviewers focus almost entirely on the
 295 paper body, with 76% of their findings on soundness and 15% on presentation; they rarely venture
 296 beyond the paper itself (8% on prior work, under 1% on code). The LLM reviewer (Opus 4.7) behaves
 297 almost identically, concentrating 90% of its attention on the paper body.

298 PaperDoctor distributes its attention very differently. Only 70% of its findings fall on the paper body,
 299 while 18% check other literature works and 12% verify the released code. The sharpest gap is on
 300 code: very small partition of human and LLM’s touching, and 12% of PaperDoctor’s. This shows that
 301 PaperDoctor plays a **complementary** role to human and LLM reviewers rather than a competing one:
 302 instead of producing a denser version of what they already write, PaperDoctor covers dimensions that
 303 a reader simply cannot reach.

304 **Experimental Reproduction Analysis** Due to the limited page space, please see Section D for a
 305 detailed reproduction analysis covering 60 papers, such as “How does an L2 claim transformed into
 306 actual execution?”, “Correlation between L2 priority and L3 reproduction success rate”, “Correla-
 307 tion between reproduction mode and success rate?”, “How does reproduction differ across paper
 308 sources?”.

309 4 Conclusions

310 We introduced PaperDoctor, an agent framework that shifts automated paper assessment from a
 311 judge to a diagnostician. Given a paper with the code, PaperDoctor produce each finding, which is
 312 a triple of an observation (Why), a pointer to a specific location (Where), and a concrete revision
 313 suggestion (How). PaperDoctor’s findings are produced through a holistic hierarchical pipeline:
 314 L1 cheap surface checks run on paper surface, L2 typed verifiers route each claim to the skill that
 315 owns its evidence, and a L3 experimental-reproduction stage selectively reruns experiments under
 316 a priority budget. This pipeline scales effort to the cost of verification. Across a human study and
 317 diverse manuscripts spanning ML, the natural sciences, and the social sciences, PaperDoctor produces
 318 substantially more grounded feedback than human or LLM reviewers, its findings track paper quality,
 319 and its reproduction stage surfaces inconsistencies invisible from the manuscript alone. We release
 320 PaperDoctor together with a demo interface, in the hope of supporting more rigorous and reproducible
 321 research in the community.

322 References

- 323 [1] Kaito Baba, Chaoran Liu, Shuhei Kurita, and Akiyoshi Sannai. Prover agent: An agent-based
324 framework for formal mathematical proofs. *arXiv preprint arXiv:2506.19923*, 2025.
- 325 [2] Federico Bianchi, Owen Queen, Nitya Thakkar, Eric Sun, and James Zou. Exploring the use of
326 ai authors and reviewers at agents4science. *Nature Biotechnology*, 44(1):11–14, 2026.
- 327 [3] Bernd Bohnet, Vinh Q Tran, Pat Verga, Roe Aharoni, Daniel Andor, Livio Baldini Soares,
328 Massimiliano Ciaramita, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, et al. Attributed
329 question answering: Evaluation and modeling for attributed large language models. *arXiv
330 preprint arXiv:2212.08037*, 2022.
- 331 [4] Gulfidan Can and Andrew Walker. A model for doctoral students’ perceptions and attitudes
332 toward written feedback for academic writing. *Research in Higher Education*, 52(5):508–536,
333 2011.
- 334 [5] Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio
335 Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, et al. Mle-bench: Evaluating machine
336 learning agents on machine learning engineering. *arXiv preprint arXiv:2410.07095*, 2024.
- 337 [6] Yuan Chang, Ziyue Li, Hengyuan Zhang, Yuanbo Kong, Yanru Wu, Hayden Kwok-Hay So,
338 Zhijiang Guo, Liya Zhu, and Ngai Wong. Treereview: A dynamic tree of questions framework
339 for deep and efficient llm-based scientific peer review. In *Proceedings of the 2025 Conference
340 on Empirical Methods in Natural Language Processing*, pages 15662–15693, 2025.
- 341 [7] Ziru Chen, Shijie Chen, Yuting Ning, Qianheng Zhang, Boshi Wang, Botao Yu, Yifei Li, Zeyi
342 Liao, Chen Wei, Zitong Lu, et al. Scienceagentbench: Toward rigorous assessment of language
343 agents for data-driven scientific discovery. *arXiv preprint arXiv:2410.05080*, 2024.
- 344 [8] Mike D’Arcy, Tom Hope, Larry Birnbaum, and Doug Downey. MARG: Multi-agent review
345 generation for scientific papers. In *arXiv preprint arXiv:2401.04259*, 2024.
- 346 [9] Christopher Foster. Openness in ai and downstream governance: A global value chain approach.
347 *arXiv preprint arXiv:2509.10220*, 2025.
- 348 [10] Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng
349 Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, et al. Rarr: Researching and revising
350 what language models say, using language models. In *Proceedings of the 61st Annual Meeting
351 of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16477–16508,
352 2023.
- 353 [11] Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. Enabling large language models to
354 generate text with citations. In *Proceedings of the 2023 Conference on Empirical Methods in
355 Natural Language Processing*, pages 6465–6488, 2023.
- 356 [12] Xian Gao, Jiacheng Ruan, Zongyun Zhang, Jingsheng Gao, Ting Liu, and Yuzhuo Fu. MMRe-
357 view: A multidisciplinary and multimodal benchmark for LLM-based peer review automation.
358 *arXiv preprint arXiv:2508.14146*, 2025.
- 359 [13] Zhaolin Gao, Kianté Brantley, and Thorsten Joachims. Reviewer2: Optimizing review genera-
360 tion through prompt generation. *arXiv preprint arXiv:2402.10886*, 2024.
- 361 [14] Mengze Hong, Di Jiang, Weiwei Zhao, Yawen Li, Yihang Wang, Xinyuan Luo, Yanjie Sun, and
362 Chen Jason Zhang. Multimodal peer review simulation with actionable to-do recommendations
363 for community-aware manuscript revisions. *arXiv preprint arXiv:2511.10902*, 2025.
- 364 [15] Tianyu Hua, Harper Hua, Violet Xiang, Benjamin Klieger, Sang Truong, Weixin Liang, Fan-
365 Yun Sun, and Nick Haber. Researchcodebench: Benchmarking LLMs on implementing novel
366 machine learning research code. *Advances in Neural Information Processing Systems*, 38, 2026.
- 367 [16] Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. Mlagentbench: Evaluating language
368 agents on machine learning experimentation. *arXiv preprint arXiv:2310.03302*, 2023.

- 369 [17] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and
370 Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv*
371 *preprint arXiv:2310.06770*, 2023.
- 372 [18] Yiqiao Jin, Qinlin Zhao, Yiyang Wang, Hao Chen, Kaijie Zhu, Yijia Xiao, and Jindong Wang.
373 Agentreview: Exploring peer review dynamics with llm agents. In *Proceedings of the 2024*
374 *Conference on Empirical Methods in Natural Language Processing*, pages 1208–1226, 2024.
- 375 [19] Rui Li, Jia-Chen Gu, Po-Nien Kung, Heming Xia, Xiangwen Kong, Zhifang Sui, Nanyun Peng,
376 et al. Llm-reval: Can we trust llm reviewers yet? *arXiv preprint arXiv:2510.12367*, 2025.
- 377 [20] Shuaimin Li, Liyang Fan, Yufang Lin, Zeyang Li, Xian Wei, Shiwen Ni, Hamid Alinejad-
378 Rokny, and Min Yang. Automatic paper reviewing with heterogeneous graph reasoning over
379 llm-simulated reviewer-author debates. In *Proceedings of the AAAI Conference on Artificial*
380 *Intelligence*, volume 40, pages 31717–31725, 2026.
- 381 [21] Weixin Liang, Yuhui Zhang, Hancheng Cao, Binglu Wang, Daisy Yi Ding, Xinyu Yang,
382 Kailas Vodrahalli, Siyu He, Daniel Scott Smith, Yian Yin, et al. Can large language models
383 provide useful feedback on research papers? a large-scale empirical analysis. *NEJM AI*,
384 1(8):AIoa2400196, 2024.
- 385 [22] Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scien-
386 tist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*,
387 2024.
- 388 [23] MZ Naser. How llms cite and why it matters: A cross-model audit of reference fabrication
389 in ai-assisted academic writing and methods to detect phantom citations. *arXiv preprint*
390 *arXiv:2603.03299*, 2026.
- 391 [24] Azim Ospanov, Zijin Feng, Jiacheng Sun, Haoli Bai, Xin Shen, and Farzan Farnia. Hermes: To-
392 wards efficient and verifiable mathematical reasoning in llms. *arXiv preprint arXiv:2511.18760*,
393 2025.
- 394 [25] Abhilasha Ravichander, Shruti Ghela, David Wadden, and Yejin Choi. Halogen: Fantastic
395 llm hallucinations and where to find them. In *Proceedings of the 63rd Annual Meeting of the*
396 *Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1402–1425, 2025.
- 397 [26] ZZ Ren, Zhihong Shao, Junxiao Song, Huajian Xin, Haocheng Wang, Wanxia Zhao, Liyue
398 Zhang, Zhe Fu, Qihao Zhu, Dejian Yang, et al. Deepseek-prover-v2: Advancing formal
399 mathematical reasoning via reinforcement learning for subgoal decomposition. *arXiv preprint*
400 *arXiv:2504.21801*, 2025.
- 401 [27] Pouria Rouzrokh, Bardia Khosravi, Parsa Rouzrokh, and Moein Shariatnia. Latterreview: a
402 multi-agent framework for systematic review automation using large language models. *arXiv*
403 *preprint arXiv:2501.05468*, 2025.
- 404 [28] Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu,
405 Zicheng Liu, and Emad Barsoum. Agent laboratory: Using LLM agents as research assistants.
406 *arXiv preprint arXiv:2501.04227*, 2025.
- 407 [29] Zachary S Siegel, Sayash Kapoor, Nitya Nagdir, Benedikt Stroebel, and Arvind Narayanan. Core-
408 bench: Fostering the credibility of published research through a computational reproducibility
409 agent benchmark. *arXiv preprint arXiv:2409.11363*, 2024.
- 410 [30] Giulio Starace, Oliver Jaffe, Dane Sherburn, James Aung, Jun Shern Chan, Leon Maksin,
411 Rachel Dias, Evan Mays, Benjamin Kinsella, Wyatt Thompson, et al. Paperbench: Evaluating
412 ai’s ability to replicate ai research. *arXiv preprint arXiv:2504.01848*, 2025.
- 413 [31] Jacob Steiss, Tamara Tate, Steve Graham, Jazmin Cruz, Michael Hebert, Jiali Wang, Youngsun
414 Moon, Waverly Tseng, Mark Warschauer, and Carol Booth Olson. Comparing the quality of
415 human and chatgpt feedback of students’ writing. *Learning and Instruction*, 91:101894, 2024.

- 416 [32] Pawin Taechoyotin and Daniel Acuna. Remor: Automated peer review generation with llm
417 reasoning and multi-objective reinforcement learning. *arXiv preprint arXiv:2505.11718*, 2025.
- 418 [33] Cheng Tan, Dongxin Lyu, Siyuan Li, Zhangyang Gao, Jingxuan Wei, Siqi Ma, Zicheng Liu, and
419 Stan Z Li. Peer review as a multi-turn and long-context dialogue with role-based interactions.
420 *arXiv preprint arXiv:2406.05688*, 2024.
- 421 [34] Jiabin Tang, Lianghao Xia, Zhonghang Li, and Chao Huang. Ai-researcher: Autonomous
422 scientific innovation. *arXiv preprint arXiv:2505.18705*, 2025.
- 423 [35] Xiangru Tang, Yuliang Liu, Zefan Cai, Yanjun Shao, Junjie Lu, Yichi Zhang, Zexuan Deng,
424 Helan Hu, Kaikai An, Ruijun Huang, et al. Ml-bench: Evaluating large language models and
425 agents for machine learning tasks on repository-level code. *arXiv preprint arXiv:2311.09835*,
426 2023.
- 427 [36] Nitya Thakkar, Mert Yuksekgonul, Jake Silberg, Animesh Garg, Nanyun Peng, Fei Sha, Rose
428 Yu, Carl Vondrick, and James Zou. Can llm feedback enhance review quality? a randomized
429 study of 20k reviews at iclr 2025. *arXiv preprint arXiv:2504.09737*, 2025.
- 430 [37] Nitya Thakkar, Mert Yuksekgonul, Jake Silberg, Animesh Garg, Nanyun Peng, Fei Sha, Rose
431 Yu, Carl Vondrick, and James Zou. A large-scale randomized study of large language model
432 feedback in peer review. *Nature Machine Intelligence*, pages 1–11, 2026.
- 433 [38] Haoxin Tu, Huan Zhao, Yahui Song, Mehtab Zafar, Ruijie Meng, and Abhik Roychoudhury.
434 Agentic program verification. *arXiv preprint arXiv:2511.17330*, 2025.
- 435 [39] Sumanth Varambally, Thomas Voice, Yanchao Sun, Zhifeng Chen, Rose Yu, and Ke Ye. Hilbert:
436 Recursively building formal proofs with informal reasoning. *arXiv preprint arXiv:2509.22819*,
437 2025.
- 438 [40] Yixuan Weng, Minjun Zhu, Guangsheng Bao, Hongbo Zhang, Jindong Wang, Yue Zhang, and
439 Linyi Yang. Cycleresearcher: Improving automated research via automated review. *arXiv*
440 *preprint arXiv:2411.00816*, 2024.
- 441 [41] Hjalmar Wijk, Tao Lin, Joel Becker, Sami Jawhar, Neev Parikh, Thomas Broadley, Lawrence
442 Chan, Michael Chen, Josh Clymer, Jai Dhyani, et al. Re-bench: Evaluating frontier ai r&d
443 capabilities of language model agents against human experts. *arXiv preprint arXiv:2411.15114*,
444 2024.
- 445 [42] Kevin Wu, Eric Wu, Kevin Wei, Angela Zhang, Allison Casasola, Teresa Nguyen, Sith Ri-
446 antawan, Patricia Shi, Daniel Ho, and James Zou. An automated framework for assessing how
447 well llms cite relevant medical references. *Nature Communications*, 16(1):3615, 2025.
- 448 [43] Yanzheng Xiang, Hanqi Yan, Shuyin Ouyang, Lin Gui, and Yulan He. SciReplicate-Bench:
449 Benchmarking LLMs in agent-driven algorithmic reproduction from research papers. *arXiv*
450 *preprint arXiv:2504.00255*, 2025.
- 451 [44] Zuyao Xu, Yuqi Qiu, Lu Sun, FaSheng Miao, Fubin Wu, Xinyi Wang, Xiang Li, Haozhe Lu,
452 ZhengZe Zhang, Yuxin Hu, et al. Ghostcite: A large-scale analysis of citation validity in the
453 age of large language models. *arXiv preprint arXiv:2602.06718*, 2026.
- 454 [45] Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Foerster, Jeff
455 Clune, and David Ha. The ai scientist-v2: Workshop-level automated scientific discovery via
456 agentic tree search. *arXiv preprint arXiv:2504.08066*, 2025.
- 457 [46] Shuo Yan, Ruochen Li, Ziming Luo, Zimu Wang, Daoyang Li, Liqiang Jing, Kaiyu He, Peilin
458 Wu, George Michalopoulos, Yue Zhang, et al. Lmr-bench: Evaluating llm agent’s ability on
459 reproducing language modeling research, 2025. URL <https://arxiv.org/abs/2506.17335>.
- 460 [47] Zhengqing Yuan, Kaiwen Shi, Zheyuan Zhang, Lichao Sun, Nitesh V Chawla, and Yanfang Ye.
461 Citeaudit: You cited it, but did you read it? a benchmark for verifying scientific references in
462 the llm era. *arXiv preprint arXiv:2602.23452*, 2026.

- 463 [48] Xuanle Zhao, Zilin Sang, Yuxuan Li, Qi Shi, Weilun Zhao, Shuo Wang, Duzhen Zhang, Xu Han,
464 Zhiyuan Liu, and Maosong Sun. Autoreproduce: Automatic ai experiment reproduction with
465 paper lineage. *arXiv preprint arXiv:2505.20662*, 2025.
- 466 [49] Minjun Zhu, Yixuan Weng, Linyi Yang, and Yue Zhang. Deepreview: Improving llm-based
467 paper review with human-like deep thinking process. In *Proceedings of the 63rd Annual Meeting*
468 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 29330–29355,
469 2025.
- 470 [50] Zhenzhen Zhuang, Jiandong Chen, Hongfeng Xu, Yuwen Jiang, and Jialiang Lin. Large lan-
471 guage models for automated scholarly paper review: A survey. *Information Fusion*, 124:103332,
472 2025.

Table 2: **PaperDoctor vs. representative research agents.** **Reviewer Report:** produces reviewer-style prose, not only a score. **Grounded Evidence:** every finding is anchored to a concrete span/figure/equation/code line. **Revision Suggestion:** output specifies *what to change*, not only *what is wrong*. **Text (LLM):** reads body text, tables, and document structure. **Visual (VLM):** reads figures as images. **Code Audit:** cross-checks paper claims against released code. **Exp. Reproduction:** re-executes experiments. **In-completed papers:** can the agent support or focus on these in-completed manuscripts. ✓=full, ✓=partial, ✗=absent.

System	Feedback Form			Assessment Coverage				Focusness
	Reviewer Report	Grounded Evidence	Revision Suggestion	Text (LLM)	Visual (VLM)	Code Imple.	Exp. Repro.	In-progress Papers
Autoresearch								
AI Scientist [45]	✓	✗	✗	✓	✗	✓	✗	✗
Peer Review								
MARG [8]	✓	✗	✓	✓	✗	✗	✗	✗
AgentReview [18]	✓	✗	✓	✓	✗	✗	✗	✗
Reviewer2 [13]	✓	✗	✓	✓	✗	✗	✗	✗
DeepReview [49]	✓	✗	✓	✓	✗	✗	✗	✗
TreeReview [6]	✓	✓	✓	✓	✗	✗	✗	✗
CycleReviewer [40]	✓	✗	✗	✓	✗	✗	✗	✗
MMReview [12]	✓	✓	✗	✓	✓	✗	✗	✗
Verification								
CiteAudit [47]	✗	✓	✗	✓	✗	✗	✗	✗
PaperBench [30]	✗	✓	✗	✓	✗	✓	✗	✗
AutoReproduce [48]	✓	✓	✗	✗	✗	✗	✓	✗
Feedbacks								
Review feedback [37]	✓	✓	✗	✓	✗	✗	✗	✗
Human advisor	✓	✓	✓	✓	✓	✓	✗	✓
PaperDoctor (ours)	✓	✓	✓	✓	✓	✓	✓	✓

473 A Related Work

474 **Autoresearch for Papers** End-to-end autoresearch pipelines chain ideation, experimentation, and
475 drafting into a single agentic loop [22, 45, 28, 34, 2], increasingly supported by language and coding
476 agents [15, 46, 43] that probe whether agents can implement and run published methods. Recent
477 variants further explore tool-augmented, multi-agent collaboration, and long-horizon execution,
478 pushing manuscripts toward greater autonomy. A shared trait, however, is that manuscripts are
479 produced without an internal quality-control stage, leaving characteristic failure modes such as
480 hallucinated citations, inflated novelty, mismatches between claims and the code that implements
481 them, and unverified empirical statements—silently propagated into the final draft. PaperDoctor is
482 complementary to this line of work: rather than producing papers, it consumes a draft together with
483 its accompanying code and data, decomposes it across writing, references, theory, and experiments,
484 and localises the artefacts that require revision before submission. In this sense, PaperDoctor acts as
485 an internal advisor that closes the quality-control gap left open by current autoresearch agents.

486 **Paper Peer Review** A rapidly growing line casts LLMs as peer reviewers. Static prompting or fine-
487 tuning yields a full review per paper [21, 8, 13, 49, 32, 27]; multi-agent variants simulate the review
488 cycle with role-specialised agents [18, 33, 20]; decomposition-based TreeReview [6] recursively
489 asks sub-questions; and CycleResearcher [40] closes the loop via iterative preference optimization.
490 Multimodal and multidisciplinary variants [12, 14] read figures and tables; pre-submission assistance
491 has been proposed as an ethically cleaner deployment [9]. At scale, the Review Feedback Agent [36]
492 was deployed on 20k ICLR-2025 reviews. Parallel audits document persistent failure modes: prompt-
493 injection susceptibility, sycophancy, and poor novelty calibration [19, 50]. These systems share a
494 judge-oriented output contract, optimised against held-out reviewer opinions rather than author utility,
495 and under-serve revision along three axes that PaperDoctor inverts. (i) Evidence-grounded. Verdicts
496 are rarely tied to a specific sentence, equation, or code line; PaperDoctor emits (finding, evidence,
497 suggestion) triples attached to concrete artefacts. (ii) Claim-level auditability. Holistic judgements let

498 individual unverifiable assertions pass silently; PaperDoctor extracts claims and verifies each one.
 499 (iii) Cost tiering. Review pipelines are binary: every check always runs or never runs; PaperDoctor
 500 exposes a three-tier pipeline so cheap screens are default, typed verifiers are routed by evidence, and
 501 full reproduction is gated on explicit author approval.

502 **Paper Verification** For rigorous assessment, a paper can be treated as a verifiable system whose
 503 claims must be checked along multiple dimensions. (i) Citation verification. LLM drafts routinely
 504 contain fabricated references, with audits reporting hallucination rates significantly [44, 23]; dedicated
 505 verifiers [47, 42] and attributed-generation frameworks [11, 10, 3, 25] supply the primitives for
 506 PaperDoctor’s reference verifier. (ii) Theoretical-claim verification. LLM-based provers [26, 1, 39,
 507 24] and agentic program verifiers [38] combine informal reasoning with Lean and Coq checking.
 508 PaperDoctor incorporates this perspective: it densely extract informal argument and check whether it
 509 matches formal or executable content. Moreover, the reproducibility is considered (iii) Experiment
 510 reproduction. Beyond textual verification, reproducibility benchmarks such as [5, 16, 35] evaluate
 511 ML engineering on Kaggle- and repo-scale tasks, while [17, 29, 7, 41] extend this to software,
 512 computational, scientific, and frontier-R&D settings, with top agents still below 40% execution
 513 accuracy [15, 46, 43]. PaperDoctor internalises the lesson that full reproduction is expensive and
 514 noisy: L3 issues a *prioritised* plan from the paper’s own claims and executes only with author
 515 approval, spending compute where evidence, not bandwidth, is the constraint.

516 B Experiment Details

517 **Implementation.** PaperDoctor is implemented as a set of Claude-skills orchestrated by the tiered
 518 pipeline of Section 2. L1 and L2 skills use CLAUDE-OPUS-4.7 for text reasoning and the same model
 519 with vision enabled for read-vis. Web lookups in read-bib and read-prior use the built-in
 520 WebSearch backend, bounded to at most three queries per claim. L3 reproduction runs on a single
 521 node with an NVIDIA A100 (80 GB). All 60 papers go through the same pipeline with the same skill
 522 configuration; the L3 reproduction gate is enabled uniformly so that read-exp plans are dispatched
 523 to run-exp under a fixed 30-minute per-claim budget.

524 C Pipeline Algorithm

525 Algorithm 1 gives the full pseudocode for the three-tier PaperDoctor pipeline summarised in Sec-
 526 tion 2.1. Stage 0 (PREPARE) is executed once per paper; L1 and L2 skills run in parallel over the
 527 shared artefact bundle \mathcal{A} ; the L3 RUNEXP stage is only executed when the user explicitly approves
 528 the reproduction plan produced by read-exp.

Algorithm 1 PaperDoctor review pipeline

Require: paper \mathcal{P} , code \mathcal{C} , data \mathcal{D}
 1: $\mathcal{A} \leftarrow \text{PREPARE}(\mathcal{P}, \mathcal{C})$
 2: $\{\mathcal{Q}, \mathcal{F}_{\text{txt}}, \mathcal{F}_{\text{vis}}, \mathcal{F}_{\text{bib}}\} \leftarrow \text{L1}(\mathcal{A})$ // parallel
 3: $\mathcal{F}_{\text{L2}} \leftarrow \emptyset$
 4: **for all** $c_i \in \mathcal{Q}$ **in parallel do**
 5: **for all** $t \in T_i$ **do**
 6: $\mathcal{F}_{\text{L2}} \leftarrow \mathcal{F}_{\text{L2}} \cup \{\text{VERIFY}_t(c_i, \mathcal{A})\}$
 7: **end for**
 8: **end for**
 9: **if** user approves reproduction **then**
 10: $\mathcal{F}_{\text{L3}} \leftarrow \text{RUNEXP}(\mathcal{F}_{\text{L2}}^{\text{exp}}, \mathcal{C}, \mathcal{D})$
 11: **else**
 12: $\mathcal{F}_{\text{L3}} \leftarrow \emptyset$
 13: **end if**
 14: **return** REPORT($\mathcal{P}, \mathcal{F}_{\text{txt}} \cup \mathcal{F}_{\text{vis}} \cup \mathcal{F}_{\text{bib}} \cup \mathcal{F}_{\text{L2}} \cup \mathcal{F}_{\text{L3}}$)

529 **D Reproduction Analysis**

530 **How does an L2 claim transformed into actual execution?** Figure 8 shows the end-to-end funnel
 531 of L2 experiments claims forwarded to L3. PaperDoctor surfaces 32.0 experimental claims per paper.
 532 About half (45.8%) are routed to L3 with a reproduction plan, while the rest are flagged as not
 533 requiring re-execution, for example simple inference checks or claims already self-verified within
 534 the paper. Of the routed claims, only 1.95 (6.1% of all L2 claims) clear the feasibility check and are
 535 ready to run, while 12.7 (39.8%) are blocked, most often by missing model weights, data access, or
 536 heavy GPU compute requirements. After human approval, 6.1 are eventually attempted, leading to
 537 3.6 passes, 2.4 warnings per paper. End-to-end, only 11.3% of L2 claims yield a reproduced number,
 538 with conditional on execution the success rate (pass plus warning) is 59.5%.

539 **Correlation between L2 priority and L3 reproduction success rate?** Figure 7a breaks down L3
 540 outcomes by L2-assigned priority. Pass rates drop monotonically from 34% for High-priority claims
 541 to 19% for Medium and 11% for Low, validating L2’s prioritisation. Higher-priority claims tend
 542 to come with richer reproduction specifications, including explicit numeric targets, well-described
 543 commands, and clearly identified entry points, and therefore more often clear the feasibility check.
 544 Lower-priority claims, by contrast, are typically auxiliary numbers such as parameter sweeps or
 545 sensitivity studies, for which authors often do not release a runnable script in the first place. As a
 546 result, most are blocked rather than executed.

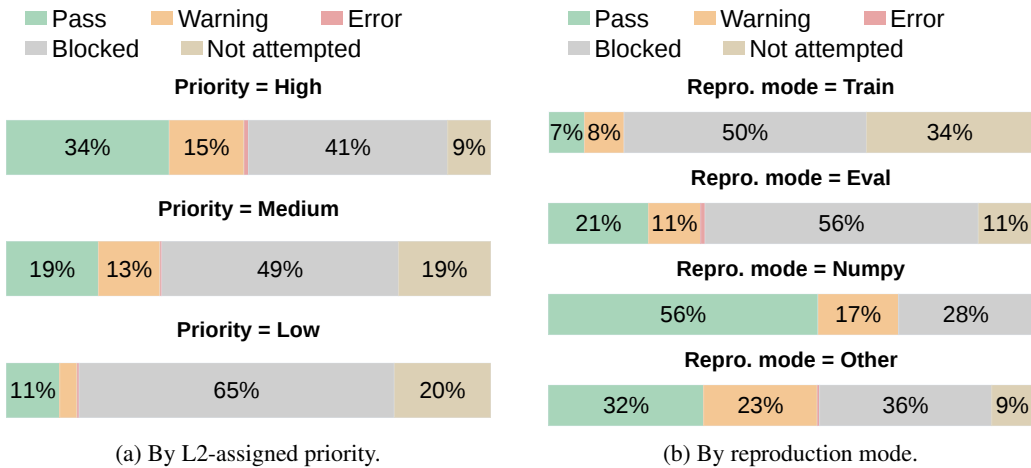


Figure 7: L3 reproduction breakdowns. (a) Pass rate by L2 priority. (b) Pass rate by reproduction mode.

547 **Correlation between reproduction mode and success rate?** Figure 7b breaks the same outcomes
 548 down by reproduction mode. Pass rates span an order of magnitude: 56% for Numpy-style standalone
 549 scripts, 32% for Other (Mainly from the social science papers), 21% for Eval (running released
 550 checkpoints), and only 7% for Train. **Training is by far the hardest mode**, with 84% of plans
 551 blocked before any code runs, because training at once needs jointly the right environment, the data,
 552 and serious compute.

553 **How does reproduction differ across paper sources?** Figure 9a breaks reproduction down across
 554 five paper sources. Plan density varies by an order of magnitude, from 24.9 plans per paper for
 555 Social-sci down to 4.6 and 2.5 for Agent4Science accepted and rejected, reflecting how much of
 556 each domain’s claims map to concrete numerical targets: **social-science papers are dominated by
 557 statistical tests**, while **agent-written papers contain very few experiments** overall and are often
 558 short-paper format. In Fig. 9b, conditional pass rates (defined as Pass / Attempted) tell a different
 559 story: Nature science leads at 67.6%, while Agent-rejected drops to 12.5%. Two failure modes are
 560 visible. PaperBench struggles before execution, with many plans blocked by environment or compute
 561 issues, but the experiments that do run usually match. Agent-rejected struggles after execution:
 562 most of its plans run, but the numbers rarely line up with what the paper reports. Notably, **“Pass /
 563 Attempted”** measures how often a paper’s experiments hold up *once execution is feasible*. By this
 564 metric, peer-reviewed papers from established venues (Nature, PaperBench oral/spotlight) sit at the

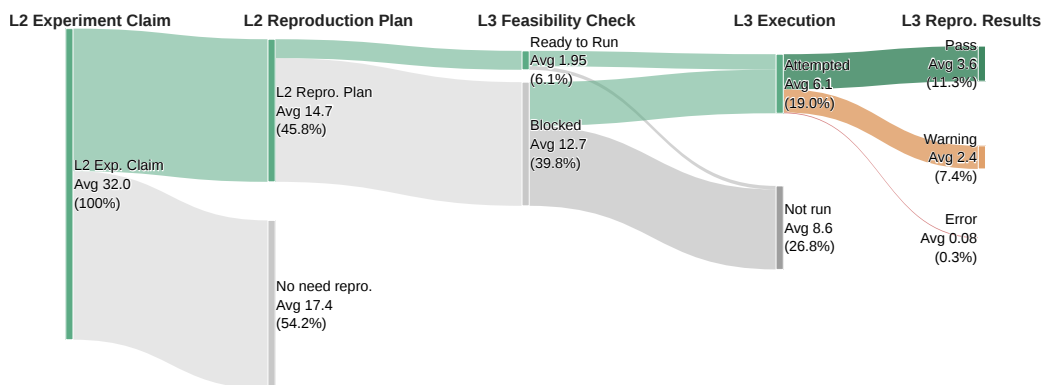


Figure 8: **End-to-end reproduction funnel across 60 papers.** PaperDoctor surfaces 32.0 experimental claims per paper on average; 14.7 receive a reproduction plan, 1.95 clear the feasibility check, and 6.1 ultimately execute. Of those that run, 3.6 pass and 2.4 yield warnings. The end-to-end yield is 11.3% reproduced claims per paper, or 59.5% conditional on execution.

565 top, while agent-written rejects sit at the bottom, suggesting that conditional pass rate can serve as a
 566 useful indicator of empirical reliability.

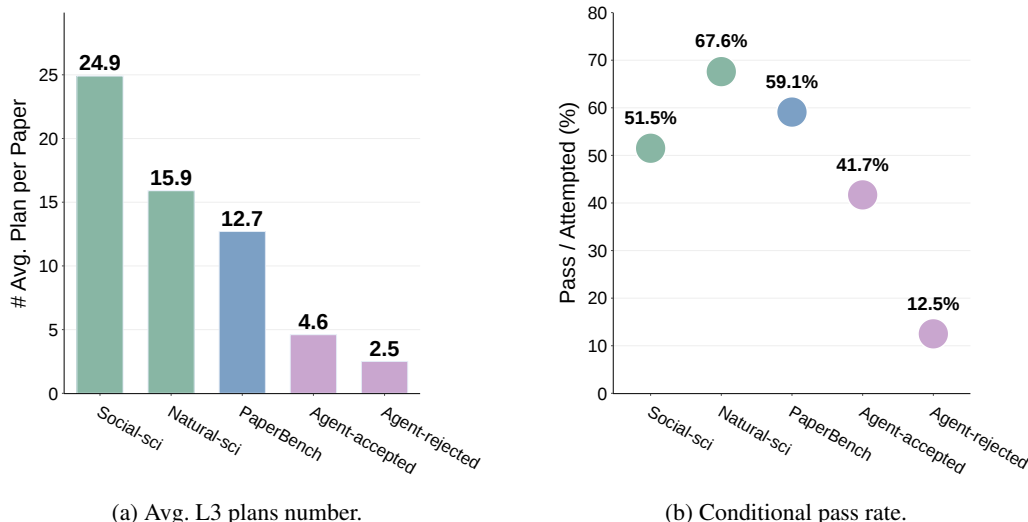


Figure 9: **L3 reproduction breakdowns by paper source.** (a) Average L3 plans per paper. (b) Conditional pass rate (pass/attempted).

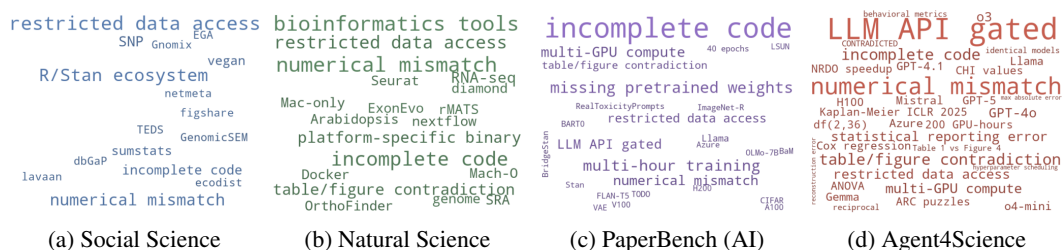


Figure 10: **Distribution of reasons that why PaperDoctor fail at reproduction stage.**

567 **What kinds of failures dominate each domain?** Figure 10 shows the failure reasons by paper
568 source, and the patterns are very different. Social Science is mostly stuck on restricted-access data
569 and R/Stan tooling. Natural Science runs into platform-specific binaries and bioinformatics pipelines.
570 PaperBench is dominated by infrastructure cost: missing weights, gated APIs, and multi-GPU training.
571 Agent4Science fails in a different way altogether: its training environment is not hard to satisfy,
572 but the numbers disagree with what the paper reports. Together, these patterns show that feasibility
573 remains a major bottleneck for reproduction, and underline the need for more reproducible and openly
574 documented research artifacts.

575 **E Case Studies**

576 The aggregate numbers in Section 3 describe PaperDoctor’s behaviour at scale, but they do not
577 show what an individual finding looks like to an author opening the report. We walk through five
578 representative cases drawn from our 60-paper benchmark, one for an L1 surface skill, three for
579 distinct L2 typed verifiers, and one for an L3 experimental reproduction. Each case follows the same
580 triple of *Why*, *Where*, and *How*, and was selected because the issue would be invisible (or at least
581 very hard to spot) under conventional reading-only review.

582 **L1 Citation Check: a fabricated future date.**

ERROR #59 ref-005

CITATION [5] OpenAI. Introducing gpt-4.5, 2024. Research preview model.

WHERE
QUOTE <https://openai.com/index/introducing-gpt-4-5/>

WHY ✓ ✗ ?
REASON Wrong year: GPT-4.5 was announced and released as a research preview on February 27, 2025, not 2024. The citation year is incorrect.

HOW ✓ ✗ ?
SUGGESTION Update the citation year from 2024 to 2025.

YOUR COMMENTS
Add your comments here...

583

584 In *Reasoning Models Outperform Standard Language Models in De Novo Protein Design*
585 (Agents4Science accepted), the bibliography contains the entry “[5] OpenAI. Introducing GPT-
586 4.5, 2024. Research preview model.” PaperDoctor’s citation verifier issues a web search for this
587 reference and finds that *GPT-4.5* was announced and released on February 27, 2025, not in 2024.
588 Although the manuscript itself appeared in late 2024, the cited release date is therefore chronologically
589 impossible. Standard citation checkers accept any correctly-formatted entry; only a search-grounded
590 verifier flags this kind of temporally inconsistent metadata, which is a recurring failure mode of
591 AI-assisted drafts that synthesise plausible-sounding venues and years without grounding them in
592 real release notes.

ERROR
#64
Claim 15
Section: method

CLAIM Materials are encoded using SciBERT into 768-dim vectors

EVIDENCE TYPE code

WHERE

QUOTE `code_data/scripts/embedding_indexing.py:21 – model_name='all-MiniLM-L6-v2'; pipeline_config.yaml:9 – embedding_model='all-MiniLM-L6-v2'; rag_retrieval.py:33 – model_name='all-MiniLM-L6-v2' encoded using SciBERT into 768-dimensional vectors`

WHY ✓ ✗ ?

REASON Paper claims SciBERT (768-dim) but code consistently uses all-MiniLM-L6-v2 (384-dim) across all files. The embedding model and dimension both differ from what the paper describes. SciBERT is never imported or used in the actual code.

HOW ✓ ✗ ?

SUGGESTION Update the paper text to describe the actual encoder (all-MiniLM-L6-v2, 384-dim), or replace the embedding implementation with SciBERT to match the claim.

YOUR COMMENTS

Add your comments here...

594

595 *LLM-Driven Discovery of High-Entropy Catalysts via Retrieval-Augmented Generation* (a second
 596 Agents4Science accepted paper) states in its Methods section that materials are “encoded using
 597 *SciBERT* into 768-dimensional vectors”. PaperDoctor’s code verifier dispatches the embedding-
 598 model claim to the released repository and finds that all three relevant files, `code_data/scripts/`
 599 `embedding_indexing.py` (line 21), `code_data/scripts/rag_retrieval.py` (line 33), and
 600 `pipeline_config.yaml` (line 9), set `model_name='all-MiniLM-L6-v2'`. *SciBERT* is never
 601 imported anywhere in the codebase, and the produced vectors are 384-dimensional rather than 768.
 602 The discrepancy spans both the model identity and its dimensionality. This is a paradigmatic case of a
 603 mismatch that is invisible to reading-only review: the paper text reads cleanly, but the implementation
 604 differs substantively from what is described.

ERROR
#27
Claim 7
Section: Section 4

CLAIM Approximation error bound: $\|G - G_\theta\|_{L^2} \leq C_1 W^{-\alpha/d} + C_2 L^{-\beta} + C_3 \|D - D_{\text{approx}}\|_\infty$
 – $D_{\text{approx}}\|_\infty$ holds for the heterogeneous solution operator

EVIDENCE TYPE theoretical

WHERE

📌 **QUOTE** Section 4 there exists a neural operator G_θ such that:
 $\|G - G_\theta\|_{L^2} \leq C_1 W^{-\alpha/d} + C_2 L^{-\beta} + C_3 \|D - D_{\text{approx}}\|_\infty$

WHY ✔ ✖ ?

🔍 **REASON** The bound is stated without proof, without stated assumptions, and without defining five key quantities: (1) W is not defined (network width? number of parameters?); (2) L is not defined (network depth? number of layers?); (3) the exponents α and β are not derived or specified; (4) the constants C_1, C_2, C_3 are not characterized; (5) it is not stated what function space G maps between. The bound cannot be verified, checked for tightness, or used to guide architecture design.

HOW ✔ ✖ ?

✍ **SUGGESTION** Define W, L, α, β , and C_1-C_3 in the surrounding text, and supply a proof or appendix reference before stating the bound.

YOUR COMMENTS

Add your comments here...

606

607 In *Neural Reaction-Diffusion Operators for Spatially Heterogeneous Tumor Modeling*
 608 (Agents4Science rejected), Section 4 states the approximation bound $\|G - G_\theta\|_{L^2} \leq C_1 W^{-\alpha/d} +$
 609 $C_2 L^{-\beta} + C_3 \|D - D_{\text{approx}}\|_\infty$ as if it were established. PaperDoctor’s step by step re-derivation finds
 610 that none of the symbols composing the bound are introduced anywhere in the paper: W (presumably
 611 network width), L (presumably depth), the exponents α and β , and the constants C_1, C_2, C_3 all
 612 appear without definition, and no proof or pointer to an appendix is given. Without these, the bound
 613 cannot be checked for tightness or used to guide architecture design. On a first pass, such a gap looks
 614 like a finished theorem; only step by step re-derivation surfaces it.

ERROR
#38
Claim 5
Section: Introduction, Section 1

CLAIM DeepONet and Fourier Neural Operators are primarily designed for homogeneous PDEs with constant coefficients and struggle with spatially heterogeneous systems

EVIDENCE TYPE related_work

WHERE

QUOTE *DeepONet [12] and Fourier Neural Operators (FNO) [10] are primarily designed for homogeneous PDEs with constant coefficients and struggle with spatially heterogeneous systems prevalent in biological applications*

WHY ✓ ✗ ?

REASON This characterization of FNO [10] and DeepONet [12] is inaccurate as a blanket statement. Li et al. (FNO, 2020) showed the method applies to parametric PDEs with varying initial conditions and forcing terms. DeepONet (Lu et al., 2021) was explicitly demonstrated on nonlinear operators including those with variable coefficients. Neither paper restricts to 'constant coefficients'. The paper provides no ablation or citation demonstrating FNO/DeepONet failure on heterogeneous settings.

HOW ✓ ✗ ?

SUGGESTION Restate the comparison to reflect the variable-coefficient settings already covered by Li 2020 and Lu 2021, or add an empirical ablation demonstrating FNO/DeepONet failure on heterogeneous PDEs.

YOUR COMMENTS

Add your comments here...

616

617 The same paper claims that “*DeepONet [12] and Fourier Neural Operators (FNO) [10]* are primarily
 618 designed for homogeneous PDEs with constant coefficients and struggle with spatially heterogeneous
 619 systems prevalent in biological applications.” PaperDoctor dispatches this novelty assertion as
 620 a web search query, and the original method papers refute the framing directly: Li et al. (FNO,
 621 2020, [10]) explicitly demonstrates the method on parametric PDEs with varying initial conditions
 622 and forcing terms, and Lu et al. (DeepONet, 2021, [12]) demonstrates nonlinear operators with
 623 variable coefficients. Neither prior work is restricted to the homogeneous setting the paper attributes
 624 to it. This pattern, where an inflated novelty claim is contradicted by the very references the paper
 625 cites, is exactly what a search-grounded literature check is built to catch.

ERROR
#49
Claim 16
Section: results

CLAIM Using multi-time-point input with soft focal loss roughly doubles recall compared to single-time-point input.

EVIDENCE TYPE experiment

WHERE

📍 **QUOTE** *Fig. 2d, data/Source_Data.xlsx With soft focal loss, this increased the recall of the best performing models by almost a factor of 2*

WHY ✔ ✖ ?

🔍 **REASON** The recall column of the shipped Figure 2d data gives 0.43 and 0.49 for the two soft-focal configurations, an actual best-pair ratio of 1.14× and a mean of 1.04× — far from a factor of 2. Both numbers are derivable from artefacts the authors themselves shipped, so the discrepancy is not environmental drift but a prose summary that overstates the effect.

HOW ✔ ✖ ?

✍ **SUGGESTION** Revise the prose to reflect the measured ~1.1× improvement, e.g. ‘soft focal loss yielded a moderate ~14% improvement in recall on the best-performing models’.

YOUR COMMENTS

Add your comments here...

627

628 In *Smart hybrid microscopy for cell-friendly detection of rare events*, a Nature Communications
 629 paper on mitochondrial imaging, the Results section claims that “soft focal loss . . . increased the
 630 recall of the best performing models by almost a factor of 2”. The L3 stage rebuilds the analysis from
 631 the source data shipped with the paper (`data/Source_Data.xlsx`) and reads the recall column
 632 for Figure 2d directly. The two soft-focal configurations have contact-task recalls of 0.43 and 0.49,
 633 giving an actual best-pair ratio of 1.14× and a mean ratio of 1.04×, far from the claimed factor of
 634 two. Both the claimed and reproduced numbers are derivable from artefacts the authors themselves
 635 shipped, so the discrepancy is not a matter of environmental drift but of how a quantitative claim was
 636 summarised in the text. The aggregate reproduction analysis (Sec. ??) shows this kind of text-vs-data
 637 mismatch concentrated on Agents4Science papers; the case here shows that even peer-reviewed
 638 venues are not immune, and that reading the paper alone cannot easily catch it.

639 **Takeaways.** Across the five cases, two patterns recur. First, every finding is grounded in a concrete
640 locus: a bibliography line, a config file, an equation, a cited paper’s abstract, or a column of a shipped
641 spreadsheet. The author can therefore audit each critique without re-reading the paper. Second,
642 PaperDoctor catches the failure modes that reading alone misses: paper text that contradicts the
643 released code or shipped data, undefined symbols hidden inside a clean-looking derivation, and
644 novelty claims that a single web search refutes. Aggregate metrics show that PaperDoctor produces
645 denser and better-grounded feedback than human or LLM reviewers; the cases here illustrate what
646 PaperDoctor’s feedback looks like in practice.

647 **F Future works**

648 Future directions stand out: (i) *Full paper-to-code reproduction.* Even when authors do not release a
649 runnable codebase, an agent could synthesise a reference implementation directly from the paper
650 itself, as PaperBench [30] does. One challenge is that this setting will yield even lower reproduction
651 rates than current system. (ii) *Human-in-the-loop collaboration.* Our analysis (Fig. 6b) shows that
652 PaperDoctor and human reviewers cover complementary aspects, motivating an interactive setup
653 where authors accept or revise each suggestion and trigger affected skills to re-run, so that human
654 judgement and agent coverage compound.