



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

Anonymous Authors¹

Abstract

Autoresearch agents are reshaping the research pipeline, but they also let flawed claims enter the literature at scale. Human advisors catch such issues through careful, traceable feedback, yet advisor-style review requires extensive manual effort and does not scale. To shift automated paper assessment from a judge to a diagnostician, we introduce PaperDoctor, an agent framework for pre-submission feedback with three key innovations: **(i) Holistic hierarchical pipeline.** Every paper is reviewed across writing, layout, references, code, theory, prior work, and experiments through three layers: L1 paper-only screening runs cheaply on every submission; L2 verifiers route each claim to the skill that owns its evidence; and L3 reproducers rerun experiments by priority. **(ii) Evidence-grounded actionable feedback.** Each finding is a triple of an observation (Why), a pointer to a specific location such as a sentence, equation, or code line (Where), and a revision suggestion (How), making every critique auditable and actionable. **(iii) Effective experimental reproduction.** Beyond reading the paper, PaperDoctor selectively rebuilds and reruns experiments based on claim importance and compute budget, surfacing reproducibility gaps and quantitative limitations that are invisible from the manuscript alone. We evaluate PaperDoctor on human studies which are junior students on their pre-submission papers, yield 85% accuracy, notably high in claim assumption validation. PaperDoctor reframes automated paper assessment as a diagnostic process rather than a verdict, taking a concrete step toward AI advisors that help authors raise the quality of scientific writing in the autoresearch era.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

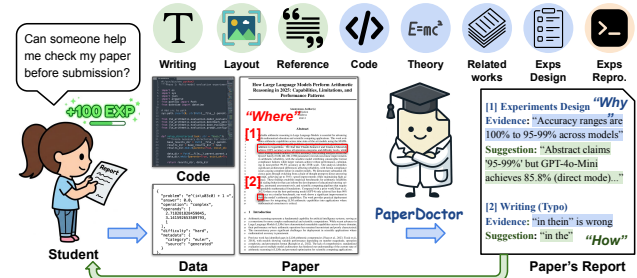


Figure 1. PaperDoctor Agent provides evidence-grounded, actionable feedback. Given a paper along with its code and datasets, PaperDoctor returns a holistic report of findings across multiple dimensions (writing, citations, figures, equations, code, reproduction, and related work). Each finding pairs 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.

1. Introduction

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

Autoresearch agents that turn a proposed idea into a full manuscript that (Lu et al., 2024; Yamada et al., 2025; Schmidgall et al., 2025; Tang et al., 2025; Bianchi et al., 2026) increases the risk of producing work whose quality cannot be guaranteed, such as these claims are often difficult to verify. Moreover, most human-written drafts are now at least partly AI-assisted, whether in writing, figure drawing, or literature survey. This trend raises paper volume while leaving draft quality uncontrolled. Existing automated reviewers (D’Arcy et al., 2024; Jin et al., 2024; Gao et al., 2025; 2024; Zhu et al., 2025; Taechoyotin and Acuna, 2025) do not close this gap. By returning a decision (“accept” or “reject”) with a justification, they primarily serve to filter submissions, which is useful for the reviewer’s side but uninformative for the author.

Human advisors catch exactly these issues during discussion with students (Can and Walker, 2011; Steiss et al., 2024). They go through the draft carefully: circling a claim with no experiments validation, underlining a theorem whose assumption breaks, or flagging a figure that is unclear. Each mark on the page is a small diagnosis. As illustrated in Figure 1, it points to *where* the symptom is in the paper, explains *why* it is a problem, and prescribes *how* to fix it; in short, it is **evidence-grounded** and **actionable**. This is

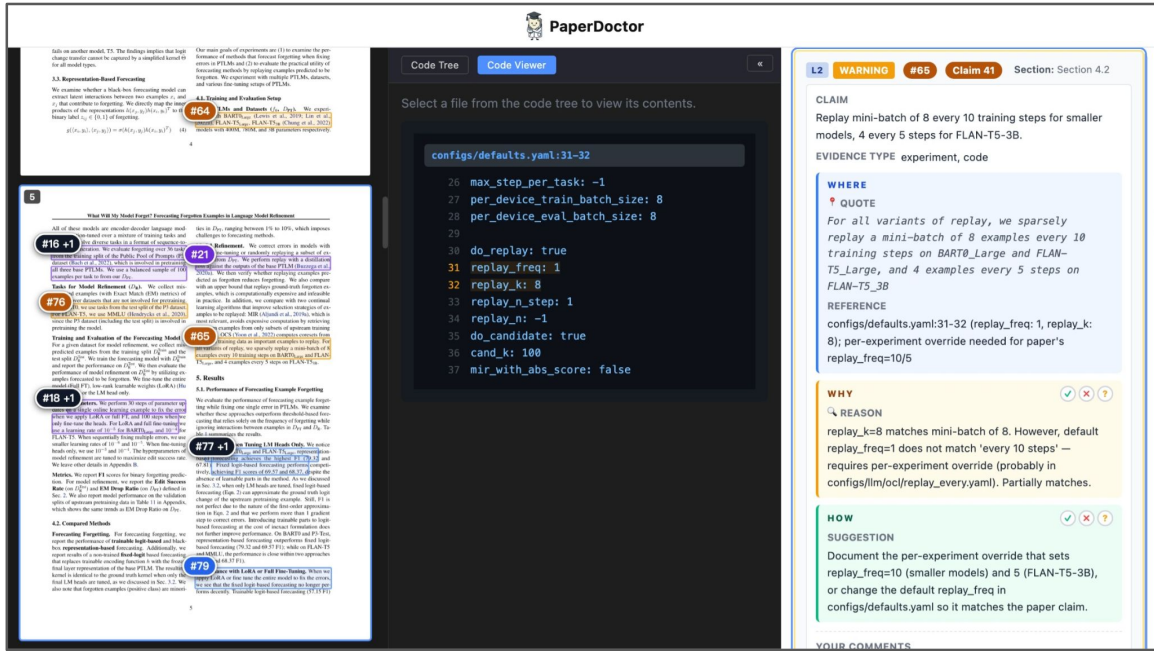


Figure 2. The PaperDoctor Diagnosis Interface. Authors upload their paper and code, and PaperDoctor returns the holistic report. **Left:** the paper, with each finding overlaid as a numbered, color-coded bubble anchored to the exact span it critiques (Where). **Middle:** the code viewer, so findings 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.

how students learn from an advisor’s diagnosis on the page, which is fundamentally different from a reviewer’s judgment. Yet such depth comes at a cost: it takes hours per paper, making it impossible to keep pace with agent-assisted writing (Lu et al., 2024; Yamada et al., 2025; Schmidgall et al., 2025; Tang et al., 2025). Current agents either write the paper for the author (Tang et al., 2025; Bianchi et al., 2026) or grade it at the door (Liang et al., 2024; Rouzrokh et al., 2025; Jin et al., 2024; Tan et al., 2024; Li et al., 2026). Neither is what a student receives from an advisor: a careful diagnosis, like that of a doctor, that points to a specific section and says what to change. However, delivering this level of feedback is non-trivial. A paper is a multi-faceted artifact spanning presentation, implementation, experimental analysis, and more. This is why authors typically rely on feedback from a range of people, such as advisors, peers, and industry mentors, each catching issues the others miss.

Motivated by this, we ask: *can an agent provide diagnosis-level feedback at scale?* We introduce PaperDoctor, a research agent framework that delivers constructive feedback to human authors. Notably, PaperDoctor highlights the following: (i) **Holistic, hierarchical pipeline.** PaperDoctor decomposes a paper, together with its code and datasets, along dimensions ranging from writing to experimental reproduction, and organizes the assessment into three hierarchical levels of increasing cost and subjectivity. **L1 (surface screening)** runs cheaply on every submission and targets

objective issues in writing, layout, references, and figures. **L2 (claim verification)** extracts atomic claims and routes each one to the skill that owns its evidence—web search for prior-work checks, a vision language model (VLM) for figure assessment, code analysis for implementation claims, and theory verifiers for derivations. **L3 (experimental reproduction)** selectively execute experiments, validating reproducibility and correctness at execution time. This design spends effort proportional to the cost of verification and keeps the full pipeline tractable. (ii) **Evidence-grounded, actionable feedback.** Each finding is a triple of an observation with a reason (Why), a pointer to a specific location in the paper (Where, *i.e.*, a sentence, equation, code line, or external URL), and a concrete suggestion for revision (How). This makes every critique both auditable and directly actionable for human authors. (iii) **Effective experimental reproduction.** Beyond reading the manuscript, PaperDoctor selectively prioritize and reruns experiments based on claim importance and compute budget. By executing code rather than only reading it, PaperDoctor surfaces reproducibility gaps and quantitative limitations that are invisible from the manuscript alone.

To validate the effectiveness of PaperDoctor, we first conduct a human study with junior student participants, collecting their in-progress paper drafts and asking them to rate the feedback produced by PaperDoctor. PaperDoctor reaches 85% agreement with their ratings. Moreover, we

evaluate PaperDoctor on 60 manuscripts covering machine learning, the natural sciences, and the social sciences, including both human- and agent-written papers with accompanying code and data. This diverse benchmark allows us to systematically assess PaperDoctor’s robustness across disciplines, writing styles, and varying levels of methodological rigor. We found that (i) PaperDoctor produces more grounded feedback than human and LLM reviewers, its findings track paper quality, and its full coverage surfaces limitations missed from the paper main body. In particular, by executing the accompanying code and re-running key experiments, PaperDoctor uncovers discrepancies between reported and actual results that are otherwise difficult to detect through reading-only review.

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

2. PaperDoctor

2.1. Overview

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

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

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

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

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

where \mathcal{P}_t is the paper as section-organized markdown (via Mathpix¹), \mathcal{P}_v is the same paper rendered page-by-page as images for downstream Vision-Language Model (VLM) use, and \mathcal{C}_t is the code indexed with tree-sitter² into per-

file units. \mathcal{B} denotes the parsed bibliography of the paper (*i.e.*, .bib file). Every downstream skill reads from this shared representation and requests only the section, page image, or code snippet it needs.

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

2.2. L1 – Surface Screening

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

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

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

¹<https://mathpix.com/>

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

165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219

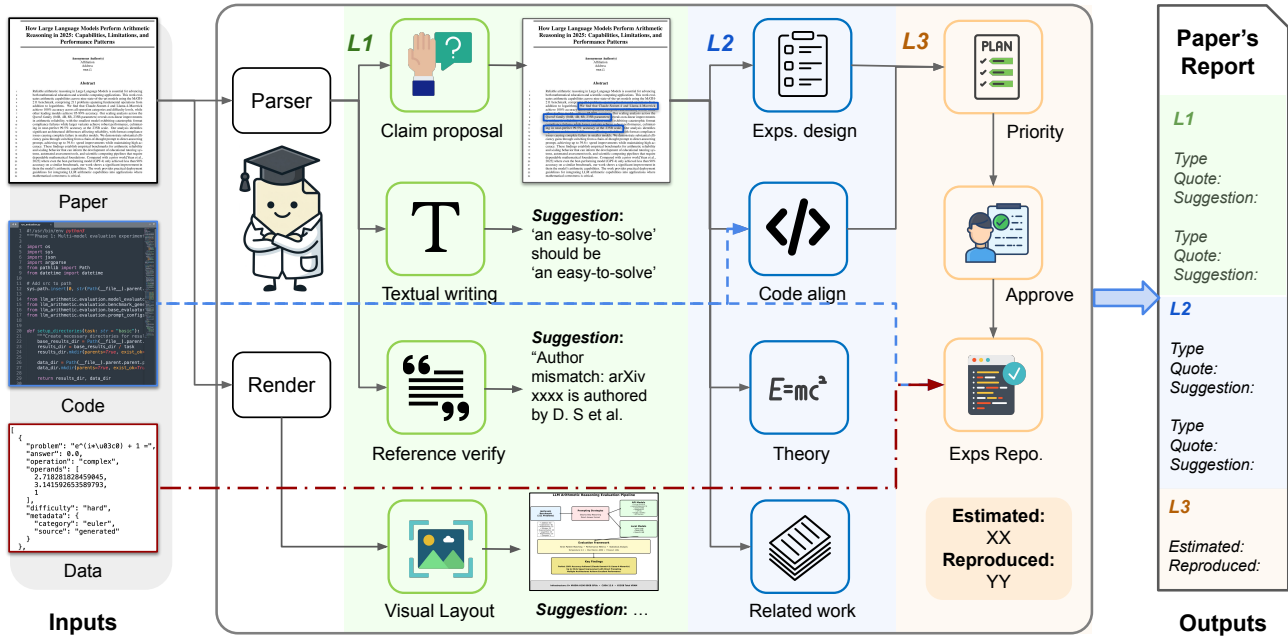


Figure 3. Illustration of PaperDoctor pipeline. A preparation step parses and renders the paper and code into a structured format and page images. (i) The L1 surface-screening stage runs four paper-only skills in parallel, including claim extraction. (ii) The L2 claim-verification stage dispatches each 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.

Citation Check. AI-assisted manuscripts routinely contain references that do not resolve to any real paper, and this is tedious to catch by reading the bibliography alone. Conditioned on the \mathcal{B} , we pair the agent with a web search backend: for each reference, the LLM issues up to three search queries 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.

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

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

2.3. L2 – Claim Verification

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

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

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

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.

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.`

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

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.

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.

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, and storage footprint, for the author to approve. Only then does an LLM-driven executor handle 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.

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

3. Experiments

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

3.1. Settings

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

Furthermore, to study the performance across various paper types, we collect 60 papers from four diverse sources (as summarized in Tab. 1): twenty come from Agents4Science (Bianchi et al., 2026), half of them accepted as orals or spotlights and half rejected, so we could ablate AI-written outcome. Another twenty are from PaperBench (Starace et al., 2025): human-written papers at top AI venues (ICML), all with released code. The remaining twenty are human-written papers from Nature Communications and Nature Human Behaviour (ten each), which pulls the evaluation outside machine learning and into the natural and social sciences. We use ‘Claude-opus-4.7’ as the base model and implement PaperDoctor as skills.

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 | ✓ |

3.2. How do author assess PaperDoctor’s feedback?

As shown in Fig. 4, we use PaperDoctor to generate 349 findings across 7 pre-submission papers, and ask each author to label every finding on their own paper. Authors accept 85% of the findings on average.

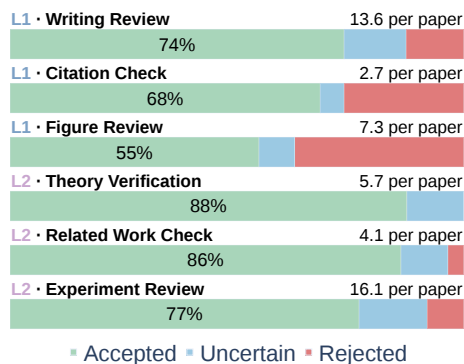


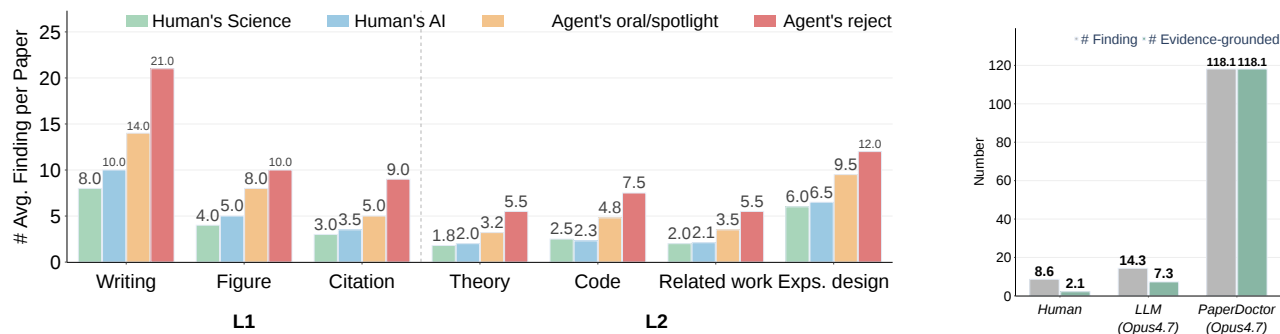
Figure 4. Author votings (Accepted / Uncertain / Rejected) on PaperDoctor’s findings, broken down by individual modules.

Where do authors agree most? As shown in Fig. 4, authors accept the majority of PaperDoctor’s findings across all six dimensions, with acceptance rates ranging from 55% (figure) to 88% (theory) without rejection. L2 dimensions, which target specific claim types such as theory and related work, achieve consistently higher acceptance (77–88%) than L1 surface-level checks (55–74%), owing to their grounding in external evidence (*i.e.*, web search) and full-paper consistency over equations and notation. Among them, Experiment review yields the largest number of findings per paper, providing authors with a systematic reminder to revisit their experimental design.

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

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

Is PaperDoctor share the same role as Humans? Figure 5a breaks PaperDoctor’s findings down by skill across four paper groups: human-written Nature papers, human-



(a) **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.

(b) **Finding number** (with and without evidence grounding), compared across human reviewers, an LLM reviewer, and PaperDoctor.

written AI papers from PaperBench, and agent-written papers that were either accepted (oral/spotlight) or rejected at Agents4Science. Several observations stand out.

Human-written papers have fewer suggestions than Agent’s papers. Across every dimension, agent papers get more findings than human ones, suggesting PaperDoctor’s finding counts track paper quality. Among human papers, Nature submissions sit at the low end on most dimensions, likely because their stricter review process catches surface issues earlier than open conference review. **Agent rejected papers have more suggestions than Agent top papers.** The gap is consistent across both L1 and L2 dimensions. Within L2, Experiment Design draws the most, ahead of theory, code, and related work. Many of these come from PaperDoctor checking the paper’s internal self-consistency, for instance whether a number in the table matches the text, or whether the reported metric matches the one defined in the method.

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

3.4. Comparison with Human Reviewers and LLM Reviewers

We study how human reviewers and LLM reviewers behave on the PaperBench papers (Starace et al., 2025) and obtain reviews directly from their authors. These reviews are not publicly available, ensuring that no current LLM has seen them during training.

How well do humans and LLMs ground their findings?

We first ask whether human and LLM reviewers actually ground their findings in the paper, using an LLM judge to label each finding as evidence-grounded or not. We report

results in Fig. 5b. We observe two gaps between baseline reviewers and PaperDoctor.

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

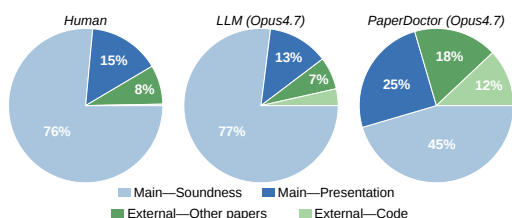
(ii) **Evidence grounding:** Only 25% of human findings and 51% of LLM findings cite a concrete locus (a figure, table, equation, or quoted line), whereas PaperDoctor grounds every one of them (100%). These gaps are not accidental. For human reviewers, writing a precise locus during review is costly: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

What do Humans, LLMs, and PaperDoctor focus on?

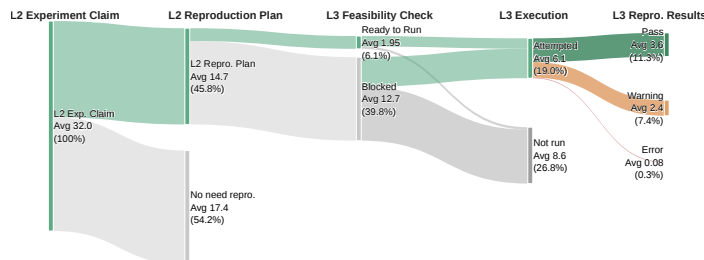
In figure 6a, we classified the review content into four buckets: soundness and presentation within the paper body, and external checks against other literatures or against the released code. Human reviewers focus almost entirely on the paper body, with 76% of their findings on soundness and 15% on presentation;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

PaperDoctor distributes its attention very differently. Only 70% of its findings fall on the paper body, while 18% check other literature works and 12% verify the code. The sharpest gap is on code: very small partition of human and LLM’s touching, and 12% of PaperDoctor’s. This shows that PaperDoctor plays a **complementary** role to human and LLM

385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439



(a) **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.



(b) **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.

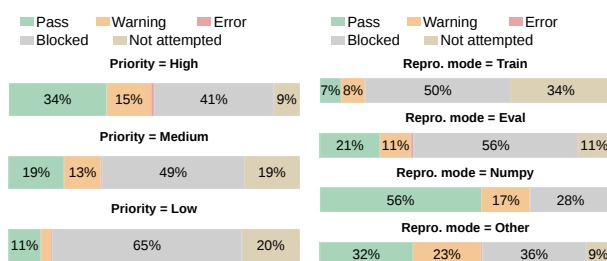
reviewers rather than a competing one: instead of producing a denser version of what they already write, PaperDoctor covers dimension that humans simply cannot reach.

3.5. Experimental Reproduction Analysis

How does an L2 claim transformed into actual execution? Figure 6b shows the end-to-end funnel of L2 experiments claims forwarded to L3. PaperDoctor surfaces 32.0 experimental claims per paper. About half (45.8%) are routed to L3 with a reproduction plan, while the rest are flagged as not requiring re-execution, for example simple inference checks or claims already self-verified within the paper. Of the routed claims, only 1.95 (6.1% of all L2 claims) clear the feasibility check and are ready to run, while 12.7 (39.8%) are blocked, most often by missing model weights, data access, or heavy GPU compute requirements. After human approval, 6.1 are eventually attempted, leading to 3.6 passes, 2.4 warnings per paper. End-to-end, only 11.3% of L2 claims yield a reproduced number, with conditional on execution the success rate (pass plus warning) is 59.5%.

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

Correlation between reproduction mode and success rate? Figure 7b breaks the same outcomes down by reproduction mode. Pass rates span an order of magnitude: 56% for Numpy-style standalone scripts, 32% for Other



(a) By L2-assigned priority. (b) By reproduction mode.

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

(Mainly from the social science papers), 21% for Eval (running released checkpoints), and only 7% for Train. **Training is by far the hardest mode**, with 84% of plans blocked before any code runs, because training at once needs jointly the right environment, the data, and serious compute.

4. Conclusions

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

References

- 440 Kaito Baba, Chaoran Liu, Shuhei Kurita, and Akiyoshi San-
441 nai. Prover agent: An agent-based framework for formal
442 mathematical proofs. *arXiv preprint arXiv:2506.19923*,
443 2025.
- 444 Federico Bianchi, Owen Queen, Nitya Thakkar, Eric Sun,
445 and James Zou. Exploring the use of ai authors and
446 reviewers at agents4science. *Nature Biotechnology*, 44
447 (1):11–14, 2026.
- 448 Bernd Bohnet, Vinh Q Tran, Pat Verga, Roei Aharoni,
449 Daniel Andor, Livio Baldini Soares, Massimiliano Cia-
450 ramita, Jacob Eisenstein, Kuzman Ganchev, Jonathan
451 Herzig, et al. Attributed question answering: Evaluation
452 and modeling for attributed large language models. *arXiv*
453 *preprint arXiv:2212.08037*, 2022.
- 454 Gulfidan Can and Andrew Walker. A model for doctoral stu-
455 dents’ perceptions and attitudes toward written feedback
456 for academic writing. *Research in Higher Education*, 52
457 (5):508–536, 2011.
- 458 Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung,
459 Dane Sherburn, Evan Mays, Giulio Starace, Kevin Liu,
460 Leon Maksin, Tejal Patwardhan, et al. Mle-bench: Evalu-
461 ating machine learning agents on machine learning engi-
462 neering. *arXiv preprint arXiv:2410.07095*, 2024.
- 463 Yuan Chang, Ziyue Li, Hengyuan Zhang, Yuanbo Kong,
464 Yanru Wu, Hayden Kwok-Hay So, Zhijiang Guo, Liya
465 Zhu, and Ngai Wong. Treereview: A dynamic tree of
466 questions framework for deep and efficient llm-based
467 scientific peer review. In *Proceedings of the 2025 Con-*
468 *ference on Empirical Methods in Natural Language Pro-*
469 *cessing*, pages 15662–15693, 2025.
- 470 Ziru Chen, Shijie Chen, Yuting Ning, Qianheng Zhang,
471 Boshi Wang, Botao Yu, Yifei Li, Zeyi Liao, Chen Wei,
472 Zitong Lu, et al. Scienceagentbench: Toward rigorous
473 assessment of language agents for data-driven scientific
474 discovery. *arXiv preprint arXiv:2410.05080*, 2024.
- 475 Mike D’Arcy, Tom Hope, Larry Birnbaum, and Doug
476 Downey. MARG: Multi-agent review generation for sci-
477 entific papers. In *arXiv preprint arXiv:2401.04259*, 2024.
- 478 Christopher Foster. Openness in ai and downstream gov-
479 ernance: A global value chain approach. *arXiv preprint*
480 *arXiv:2509.10220*, 2025.
- 481 Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen,
482 Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao,
483 Ni Lao, Hongrae Lee, Da-Cheng Juan, et al. Rarr: Re-
484 searching and revising what language models say, using
485 language models. In *Proceedings of the 61st Annual*
486 *Meeting of the Association for Computational Linguistics*
487 *(Volume 1: Long Papers)*, pages 16477–16508, 2023a.
- 488 Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen.
489 Enabling large language models to generate text with
490 citations. In *Proceedings of the 2023 Conference on Em-*
491 *pirical Methods in Natural Language Processing*, pages
492 6465–6488, 2023b.
- 493 Xian Gao, Jiacheng Ruan, Zongyun Zhang, Jingsheng Gao,
494 Ting Liu, and Yuzhuo Fu. MMReview: A multidisciplinary
and multimodal benchmark for LLM-based peer
review automation. *arXiv preprint arXiv:2508.14146*,
2025.
- Zhaolin Gao, Kianté Brantley, and Thorsten Joachims. Re-
viewer2: Optimizing review generation through prompt
generation. *arXiv preprint arXiv:2402.10886*, 2024.
- Mengze Hong, Di Jiang, Weiwei Zhao, Yawen Li, Yihang
Wang, Xinyuan Luo, Yanjie Sun, and Chen Jason Zhang.
Multimodal peer review simulation with actionable to-
do recommendations for community-aware manuscript
revisions. *arXiv preprint arXiv:2511.10902*, 2025.
- Tianyu Hua, Harper Hua, Violet Xiang, Benjamin Klieger,
Sang Truong, Weixin Liang, Fan-Yun Sun, and Nick
Haber. Researchcodebench: Benchmarking LLMs on
implementing novel machine learning research code. *Ad-*
vances in Neural Information Processing Systems, 38,
2026.
- Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec.
Mlagentbench: Evaluating language agents on ma-
chine learning experimentation. *arXiv preprint*
arXiv:2310.03302, 2023.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu
Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-
bench: Can language models resolve real-world github
issues? *arXiv preprint arXiv:2310.06770*, 2023.
- Yiqiao Jin, Qinlin Zhao, Yiyang Wang, Hao Chen, Kai-
jie Zhu, Yijia Xiao, and Jindong Wang. Agentreview:
Exploring peer review dynamics with llm agents. In *Pro-*
ceedings of the 2024 Conference on Empirical Methods in
Natural Language Processing, pages 1208–1226, 2024.
- Rui Li, Jia-Chen Gu, Po-Nien Kung, Heming Xia, Xi-
angwen Kong, Zhifang Sui, Nanyun Peng, et al. Llm-
reval: Can we trust llm reviewers yet? *arXiv preprint*
arXiv:2510.12367, 2025.
- Shuaimin Li, Liyang Fan, Yufang Lin, Zeyang Li, Xian
Wei, Shiwen Ni, Hamid Alinejad-Rokny, and Min Yang.
Automatic paper reviewing with heterogeneous graph rea-
soning over llm-simulated reviewer-author debates. In
Proceedings of the AAAI Conference on Artificial Intelli-
gence, volume 40, pages 31717–31725, 2026.

- 495 Weixin Liang, Yuhui Zhang, Hancheng Cao, Binglu Wang,
496 Daisy Yi Ding, Xinyu Yang, Kailas Vodrahalli, Siyu
497 He, Daniel Scott Smith, Yian Yin, et al. Can large lan-
498 guage models provide useful feedback on research pa-
499 pers? a large-scale empirical analysis. *NEJM AI*, 1(8):
500 AIoa2400196, 2024.
- 501 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster,
502 Jeff Clune, and David Ha. The ai scientist: Towards
503 fully automated open-ended scientific discovery. *arXiv*
504 *preprint arXiv:2408.06292*, 2024.
- 505 MZ Naser. How llms cite and why it matters: A cross-
506 model audit of reference fabrication in ai-assisted aca-
507 demic writing and methods to detect phantom citations.
508 *arXiv preprint arXiv:2603.03299*, 2026.
- 509 Azim Ospanov, Zijin Feng, Jiacheng Sun, Haoli Bai, Xin
510 Shen, and Farzan Farnia. Hermes: Towards efficient and
511 verifiable mathematical reasoning in llms. *arXiv preprint*
512 *arXiv:2511.18760*, 2025.
- 513 Abhilasha Ravichander, Shruti Ghela, David Wadden, and
514 Yejin Choi. Halogen: Fantastic llm hallucinations and
515 where to find them. In *Proceedings of the 63rd Annual*
516 *Meeting of the Association for Computational Linguistics*
517 *(Volume 1: Long Papers)*, pages 1402–1425, 2025.
- 518 ZZ Ren, Zhihong Shao, Junxiao Song, Huajian Xin,
519 Haocheng Wang, Wanxia Zhao, Liyue Zhang, Zhe Fu,
520 Qihao Zhu, Dejian Yang, et al. Deepseek-prover-v2:
521 Advancing formal mathematical reasoning via reinforce-
522 ment learning for subgoal decomposition. *arXiv preprint*
523 *arXiv:2504.21801*, 2025.
- 524 Pouria Rouzrokh, Bardia Khosravi, Parsa Rouzrokh, and
525 Moein Shariatnia. Latterreview: a multi-agent framework
526 for systematic review automation using large language
527 models. *arXiv preprint arXiv:2501.05468*, 2025.
- 528 Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun,
529 Jialian Wu, Xiaodong Yu, Jiang Liu, Zicheng Liu, and
530 Emad Barsoum. Agent laboratory: Using LLM agents
531 as research assistants. *arXiv preprint arXiv:2501.04227*,
532 2025.
- 533 Zachary S Siegel, Sayash Kapoor, Nitya Nagdir, Benedikt
534 Stroebel, and Arvind Narayanan. Core-bench: Fostering
535 the credibility of published research through a computa-
536 tional reproducibility agent benchmark. *arXiv preprint*
537 *arXiv:2409.11363*, 2024.
- 538 Giulio Starace, Oliver Jaffe, Dane Sherburn, James Aung,
539 Jun Shern Chan, Leon Maksin, Rachel Dias, Evan Mays,
540 Benjamin Kinsella, Wyatt Thompson, et al. Paperbench:
541 Evaluating ai’s ability to replicate ai research. *arXiv*
542 *preprint arXiv:2504.01848*, 2025.
- 543 Jacob Steiss, Tamara Tate, Steve Graham, Jazmin Cruz,
544 Michael Hebert, Jiali Wang, Youngsun Moon, Waverly
545 Tseng, Mark Warschauer, and Carol Booth Olson. Com-
546 paring the quality of human and chatgpt feedback of
547 students’ writing. *Learning and Instruction*, 91:101894,
548 2024.
- 549 Pawin Taechoyotin and Daniel Acuna. Remor: Auto-
mated peer review generation with llm reasoning and
multi-objective reinforcement learning. *arXiv preprint*
arXiv:2505.11718, 2025.
- Cheng Tan, Dongxin Lyu, Siyuan Li, Zhangyang Gao, Jingx-
uan Wei, Siqi Ma, Zicheng Liu, and Stan Z Li. Peer
review as a multi-turn and long-context dialogue with
role-based interactions. *arXiv preprint arXiv:2406.05688*,
2024.
- Jiabin Tang, Lianghao Xia, Zhonghang Li, and Chao Huang.
Ai-researcher: Autonomous scientific innovation. *arXiv*
preprint arXiv:2505.18705, 2025.
- Xiangru Tang, Yuliang Liu, Zefan Cai, Yanjun Shao, Junjie
Lu, Yichi Zhang, Zexuan Deng, Helan Hu, Kaikai An,
Ruijun Huang, et al. Ml-bench: Evaluating large lan-
guage models and agents for machine learning tasks on
repository-level code. *arXiv preprint arXiv:2311.09835*,
2023.
- Nitya Thakkar, Mert Yuksekogunul, Jake Silberg, Animesh
Garg, Nanyun Peng, Fei Sha, Rose Yu, Carl Vondrick, and
James Zou. Can llm feedback enhance review quality?
a randomized study of 20k reviews at iclr 2025. *arXiv*
preprint arXiv:2504.09737, 2025.
- Nitya Thakkar, Mert Yuksekogunul, Jake Silberg, Animesh
Garg, Nanyun Peng, Fei Sha, Rose Yu, Carl Vondrick,
and James Zou. A large-scale randomized study of large
language model feedback in peer review. *Nature Machine*
Intelligence, pages 1–11, 2026.
- Haoxin Tu, Huan Zhao, Yahui Song, Mehtab Zafar, Rui-
jie Meng, and Abhik Roychoudhury. Agentic program
verification. *arXiv preprint arXiv:2511.17330*, 2025.
- Sumanth Varambally, Thomas Voice, Yanchao Sun, Zhifeng
Chen, Rose Yu, and Ke Ye. Hilbert: Recursively building
formal proofs with informal reasoning. *arXiv preprint*
arXiv:2509.22819, 2025.
- Yixuan Weng, Minjun Zhu, Guangsheng Bao, Hongbo
Zhang, Jindong Wang, Yue Zhang, and Linyi Yang. Cy-
clereviewer: Improving automated research via auto-
mated review. *arXiv preprint arXiv:2411.00816*, 2024.
- Hjalmar Wijk, Tao Lin, Joel Becker, Sami Jawhar, Neev
Parikh, Thomas Broadley, Lawrence Chan, Michael Chen,
Josh Clymer, Jai Dhyani, et al. Re-bench: Evaluating

- 550 frontier ai r&d capabilities of language model agents
551 against human experts. *arXiv preprint arXiv:2411.15114*,
552 2024.
- 553 Kevin Wu, Eric Wu, Kevin Wei, Angela Zhang, Allison
554 Casasola, Teresa Nguyen, Sith Riantawan, Patricia Shi,
555 Daniel Ho, and James Zou. An automated framework for
556 assessing how well llms cite relevant medical references.
557 *Nature Communications*, 16(1):3615, 2025.
- 559 Yanzheng Xiang, Hanqi Yan, Shuyin Ouyang, Lin Gui, and
560 Yulan He. SciReplicate-Bench: Benchmarking LLMs
561 in agent-driven algorithmic reproduction from research
562 papers. *arXiv preprint arXiv:2504.00255*, 2025.
- 564 Zuyao Xu, Yuqi Qiu, Lu Sun, FaSheng Miao, Fubin Wu,
565 Xinyi Wang, Xiang Li, Haozhe Lu, ZhengZe Zhang,
566 Yuxin Hu, et al. Ghostcite: A large-scale analysis of cita-
567 tion validity in the age of large language models. *arXiv*
568 *preprint arXiv:2602.06718*, 2026.
- 569 Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran
570 Hu, Chris Lu, Jakob Foerster, Jeff Clune, and David
571 Ha. The ai scientist-v2: Workshop-level automated sci-
572 entific discovery via agentic tree search. *arXiv preprint*
573 *arXiv:2504.08066*, 2025.
- 575 Shuo Yan, Ruochen Li, Ziming Luo, Zimu Wang, Daoyang
576 Li, Liqiang Jing, Kaiyu He, Peilin Wu, George
577 Michalopoulos, Yue Zhang, et al. Lmr-bench: Evaluating
578 llm agent’s ability on reproducing language modeling
579 research, 2025. URL <https://arxiv.org/abs/2506.17335>.
- 581 Zhengqing Yuan, Kaiwen Shi, Zheyuan Zhang, Lichao Sun,
582 Nitesh V Chawla, and Yanfang Ye. Citeaudit: You
583 cited it, but did you read it? a benchmark for verify-
584 ing scientific references in the llm era. *arXiv preprint*
585 *arXiv:2602.23452*, 2026.
- 586 Xuanle Zhao, Zilin Sang, Yuxuan Li, Qi Shi, Weilun Zhao,
587 Shuo Wang, Duzhen Zhang, Xu Han, Zhiyuan Liu, and
588 Maosong Sun. Autoreproduce: Automatic ai experi-
589 ment reproduction with paper lineage. *arXiv preprint*
590 *arXiv:2505.20662*, 2025.
- 592 Minjun Zhu, Yixuan Weng, Linyi Yang, and Yue Zhang.
593 Deepreview: Improving llm-based paper review with
594 human-like deep thinking process. In *Proceedings of*
595 *the 63rd Annual Meeting of the Association for Com-*
596 *putational Linguistics (Volume 1: Long Papers)*, pages
597 29330–29355, 2025.
- 599 Zhenzhen Zhuang, Jiandong Chen, Hongfeng Xu, Yuwen
600 Jiang, and Jialiang Lin. Large language models for au-
601 tomated scholarly paper review: A survey. *Information*
602 *Fusion*, 124:103332, 2025.
- 603
604

Appendix

A. Related Work

Autoresearch for Papers End-to-end autoresearch pipelines chain ideation, experimentation, and drafting into a single agentic loop (Lu et al., 2024; Yamada et al., 2025; Schmidgall et al., 2025; Tang et al., 2025; Bianchi et al., 2026), increasingly supported by language and coding agents (Hua et al., 2026; Yan et al.; Xiang et al., 2025) that probe whether agents can implement and run published methods. Recent variants further explore tool-augmented, multi-agent collaboration, and long-horizon execution, pushing manuscripts toward greater autonomy. A shared trait, however, is that manuscripts are produced without an internal quality-control stage, leaving characteristic failure modes such as hallucinated citations, inflated novelty, mismatches between claims and the code that implements them, and unverified empirical statements—silently propagated into the final draft. PaperDoctor is complementary to this line of work: rather than producing papers, it consumes a draft together with its accompanying code and data, decomposes it across writing, references, theory, and experiments, and localises the artefacts that require revision before submission. In this sense, PaperDoctor acts as an internal advisor that closes the quality-control gap left open by current autoresearch agents.

Paper Peer Review A rapidly growing line casts LLMs as peer reviewers. Static prompting or fine-tuning yields a full review per paper (Liang et al., 2024; D’Arcy et al., 2024; Gao et al., 2024; Zhu et al., 2025; Taechoyotin and Acuna, 2025; Rouzrokh et al., 2025); multi-agent variants simulate the review cycle with role-specialised agents (Jin et al., 2024; Tan et al., 2024; Li et al., 2026); decomposition-based TreeReview (Chang et al., 2025) recursively asks sub-questions; and CycleResearcher (Weng et al., 2024) closes the loop via iterative preference optimization. Multimodal and multidisciplinary variants (Gao et al., 2025; Hong et al., 2025) read figures and tables; pre-submission assistance has been proposed as an ethically cleaner deployment (Foster, 2025). At scale, the Review Feedback Agent (Thakkar et al., 2025) was deployed on 20k ICLR-2025 reviews. Parallel audits document persistent failure modes: prompt-injection susceptibility, sycophancy, and poor novelty calibration (Li et al., 2025; Zhuang et al., 2025). These systems share a judge-oriented output contract, optimised against held-out reviewer opinions rather than author utility, and under-serve revision along three axes that PaperDoctor inverts. (i) Evidence-grounded. Verdicts are rarely tied to a specific sentence, equation, or code line; PaperDoctor emits (finding, evidence, suggestion) triples attached to concrete artefacts. (ii) Claim-level auditability. Holistic judgements let individual unverifiable assertions pass silently; PaperDoctor

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

Paper Verification For rigorous assessment, a paper can be treated as a verifiable system whose claims must be checked along multiple dimensions. (i) Citation verification. LLM drafts routinely contain fabricated references, with audits reporting hallucination rates significantly (Xu et al., 2026; Naser, 2026); dedicated verifiers (Yuan et al., 2026; Wu et al., 2025) and attributed-generation frameworks (Gao et al., 2023b;a; Bohnet et al., 2022; Ravichander et al., 2025) supply the primitives for PaperDoctor’s reference verifier. (ii) Theoretical-claim verification. LLM-based provers (Ren et al., 2025; Baba et al., 2025; Varambally et al., 2025; Ospanov et al., 2025) and agentic program verifiers (Tu et al., 2025) combine informal reasoning with Lean and Coq checking. PaperDoctor incorporates this perspective: it densely extract informal argument and check whether it matches formal or executable content. Moreover, the reproducibility is considered (iii) Experiment reproduction. Beyond textual verification, reproducibility benchmarks such as (Chan et al., 2024; Huang et al., 2023; Tang et al., 2023) evaluate ML engineering on Kaggle- and repo-scale tasks, while (Jimenez et al., 2023; Siegel et al., 2024; Chen et al., 2024; Wijk et al., 2024) extend this to software, computational, scientific, and frontier-R&D settings, with top agents still below 40% execution accuracy (Hua et al., 2026; Yan et al.; Xiang et al., 2025). PaperDoctor internalises the lesson that full reproduction is expensive and noisy: L3 issues a *prioritised* plan from the paper’s own claims and executes only with author approval, spending compute where evidence, not bandwidth, is the constraint.

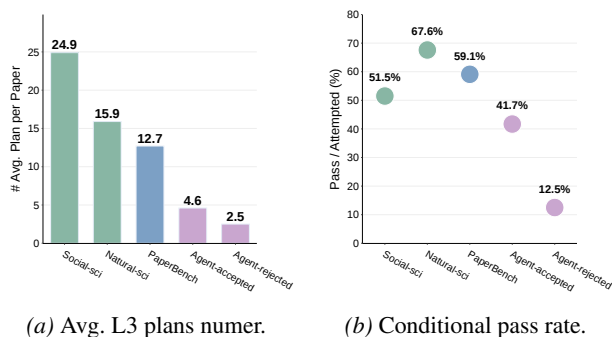


Figure 8. **L3 reproduction breakdowns** (a) Average L3 plans per paper, by paper source. (b) Conditional pass rate (pass/attempted), by paper source.

How does reproduction differ across paper sources? Figure 8a breaks reproduction down across five paper sources. Plan density varies by an order of magnitude, from 24.9

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 (Yamada et al., 2025) | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ |
| Peer Review | | | | | | | | |
| MARG (D’Arcy et al., 2024) | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| AgentReview (Jin et al., 2024) | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| Reviewer2 (Gao et al., 2024) | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| DeepReview (Zhu et al., 2025) | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| TreeReview (Chang et al., 2025) | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| CycleReviewer (Weng et al., 2024) | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| MMReview (Gao et al., 2025) | ✓ | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Verification | | | | | | | | |
| CiteAudit (Yuan et al., 2026) | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| PaperBench (Starace et al., 2025) | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ |
| AutoReproduce (Zhao et al., 2025) | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ |
| Feedbacks | | | | | | | | |
| Review feedback (Thakkar et al., 2026) | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| Human advisor | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ |
| PaperDoctor (ours) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

plans per paper for Social-sci down to 4.6 and 2.5 for Agent4Science accepted and rejected, reflecting how much of each domain’s claims map to concrete numerical targets: **social-science papers are dominated by statistical tests**, while **agent-written papers contain very few experiments** overall and are often short-paper format. In Fig. 8b, conditional pass rates (defined as Pass / Attempted) tell a different story: Nature science leads at 67.6%, while Agent-rejected drops to 12.5%. Two failure modes are visible. PaperBench struggles before execution, with many plans blocked by environment or compute issues, but the experiments that do run usually match. Agent-rejected struggles after execution: most of its plans run, but the numbers rarely line up with what the paper reports. Notably, **“Pass / Attempted”** measures how often a paper’s experiments hold up *once execution is feasible*. By this metric, peer-reviewed papers from established venues (Nature, PaperBench oral/spotlight) sit at the top, while agent-written rejects sit at the bottom, suggesting that conditional pass rate can serve as a useful indicator of empirical reliability.

B. Additional Experiments

What kinds of failures dominate each domain? Figure 9 shows the failure reasons by paper source, and the patterns are very different. Social Science is mostly stuck on restricted-access data and R/Stan tooling. Natural Sci-

ence runs into platform-specific binaries and bioinformatics pipelines. PaperBench is dominated by infrastructure cost: missing weights, gated APIs, and multi-GPU training. Agent4Science fails in a different way altogether: its training environment is not hard to satisfy, but the numbers disagree with what the paper reports. Together, these patterns show that feasibility remains a major bottleneck for reproduction, and underline the need for more reproducible and openly documented research artifacts.

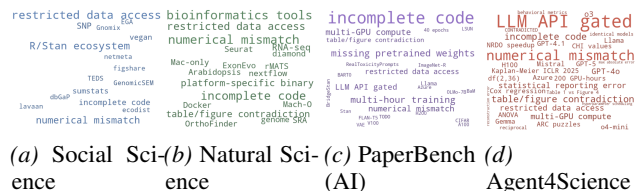


Figure 9. Distribution of reasons that why PaperDoctor fail at reproduction stage.

C. Future works

Future directions stand out: (i) *Full paper-to-code reproduction*. Even when authors do not release a runnable codebase, an agent could synthesise a reference implementation directly from the paper itself, as PaperBench (Starace et al., 2025) does. One challenge is that this setting will yield even lower reproduction rates than current system. (ii) *Human-*

715 *in-the-loop collaboration*. Our analysis (Fig. 6a) shows that
716 PaperDoctor and human reviewers cover complementary as-
717 pects, motivating an interactive setup where authors accept
718 or revise each suggestion and trigger affected skills to re-run,
719 so that human judgement and agent coverage compound.

720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769