



Description

The arrival of powerful large language models (LLMs) has changed scholarly writing and posed new risks to research integrity. Evidence from large-scale studies suggests that a non-trivial share of recent biomedical abstracts show stylistic signals consistent with LLM intervention one [analysis](#) estimated at least 13.5% of 2024 biomedical abstracts were processed with LLMs. This dual reality widespread utility and emerging misuse begs the question why some AI-generated fraudulent papers are quickly exposed and retracted while others remain undetected for longer. This article explains why detection is inconsistent, what factors determine exposure, and practical steps researchers and research managers can take to reduce risk and preserve trust in scholarship.

Why Some AI-Generated Papers Get Exposed

Detection often hinges on a combination of telltale linguistic patterns, editorial scrutiny, and contextual red flags. Editors and reviewers spot anomalies such as unnatural phrasing, inconsistent terminology, or references that cannot be verified; these cues can trigger closer checks that reveal AI-generated passages or fabricated citations. Some journals have also added automated screening to editorial triage; combined human review and technical checks increase the likelihood that AI-origin content will be flagged early. High-profile publisher investigations have led to mass retractions when clusters of submissions share similar stylistic fingerprints or originate from the same institutions. For example, [an investigation](#) into a Springer Nature journal in 2025 resulted in scores of retractions after editors concluded many commentaries and letters showed strong indications of large language model (LLM) generation without disclosure.

In the teaching context, detection vendors report large volumes of student submissions with probable AI content. [Turnitin](#) has stated that its tools reviewed hundreds of millions of student papers and flagged a substantial share as containing AI-generated content, a figure that helped spark institutional responses and policy changes. Such large-scale scanning, when combined with human follow-up, explains many exposures outside research publishing.

Why Some AI-Generated Papers Slip Through the Cracks

Detection tools and workflows are far from foolproof. [Independent evaluations](#) show wide variance in detector accuracy and significant vulnerability to simple evasive techniques. Systematic testing of multiple detectors reported mixed results, with many tools scoring below high reliability thresholds and performance degrading when AI-generated text was paraphrased or edited by humans. This inconsistency means some altered or carefully post-edited AI drafts evade automated flags, and if editors or reviewers do not notice linguistic or citation anomalies, the manuscript proceeds to publication.

Other factors that allow AI-generated content to pass include discipline-specific writing conventions (which can mask AI style), limited time for peer reviewers to perform deep verification, and the difficulty of spotting factual hallucinations in long, domain-specific texts. Additionally, authors can use tools designed to obfuscate machine origin (for example, paraphrasing networks or text “humanizers”), reducing detector scores without necessarily improving factual accuracy. Emerging detection approaches can sometimes be bypassed at modest cost in time or resources.

The Technical and Methodological Landscape: Detectors, Evasion, and Specialized Classifiers

Detection methods range from simple linguistic-feature classifiers to more advanced watermarking proposals and specialized machine-learning classifiers trained on journal-specific corpora. A 2023 [study](#) demonstrated a specialized classifier that distinguished ChatGPT-generated chemistry introductions from human-authored ones with very high accuracy in that narrow domain; however, its success depended on domain-specific training and may not generalize across disciplines. At the same time, research shows that paraphrasing or minimal human edits can drastically reduce the detection scores of general-purpose detectors, and new methods such as prompting an LLM to rewrite a text and measuring editing distance are under development to improve robustness. These findings illustrate a cat-and-mouse dynamic: specialized detectors may perform well for certain journal styles, but general detectors remain vulnerable to obfuscation.

Editorial Practices and Contextual Signals That Matter

Journals and editors rarely rely on a single signal. Exposure is most likely when multiple red flags align: unusual submission volume from the same affiliation, repetitive or mechanical language across different manuscripts, unverifiable references, inconsistent author contributions, and reviewer reports that raise methodological questions. Policies that require explicit disclosure of AI assistance (and name which tools were used and for what purpose) make it easier to identify undisclosed reliance. In contrast, inconsistent disclosure expectations across journals and disciplines produce gaps that allow undisclosed AI use to go unnoticed. Publisher-level audits or [whistleblower reports](#) also play a role in uncovering patterns of misuse.

Practical Steps for Researchers: What to Do and What to Avoid

Researchers can reduce the risk of exposure and retraction by adopting transparent, verifiable practices. The following checklist provides immediate action items that fit most disciplines:

1. **Disclose AI assistance:** If an LLM or other generative tool contributed to drafting, editing, or data handling, state the tool, version, and nature of assistance (for example, “language editing and phrasing suggestions only”). Place this statement in the Methods or Acknowledgements section as per journal guidance.
2. **Verify every citation and factual claim:** Never accept AI-suggested references at face value check that each source exists and supports the point made.
3. **Preserve human accountability:** Ensure authors can explain and defend key conceptual choices, analyses, and conclusions during [peer review](#). If AI produced a draft, authors should substantially rewrite and contextualize it to reflect original reasoning.
4. **Keep revision logs:** Maintain internal version control showing human edits and decision points to evidence authorship and contribution.
5. **Use AI for low-risk tasks:** Limit generative AI to language polishing, grammar checks, or formatting, and avoid relying on it for interpretation, data analysis, or synthesis without rigorous human oversight.

Tips for Institutions and Journals

- **Make disclosure mandatory** and define acceptable vs. unacceptable AI uses in clear, discipline-sensitive language.
- **Train editors and reviewers** to recognize linguistic and citation anomalies and to verify references as part of the review workflow.
- **Use detection tools as a triage step** never as definitive evidence and pair automated flags with human inspection.
- **Foster transparent processes** for investigating suspected misuse that protect due process for authors and minimize harm from false positives. Recent university and publisher reversals of [detector-driven accusations](#) illustrate the risk of over-reliance on imperfect tools.

How Detection Strategies Are Evolving

Detection is becoming more sophisticated and contextual. Domain-specific classifiers trained on journal text, methods that measure how an LLM itself rewrites content, and proposals for cryptographic or embedded watermarks are part of a multi-pronged approach. However, as detection tools evolve, so do techniques for evasion, especially when human editing is combined with AI output. No single technical solution will be definitive soon: effective governance will pair detection with training, disclosure requirements, and editorial judgment to sustain trust while allowing legitimate, responsible use of AI tools.

Common Mistakes to Avoid

Relying solely on an AI-detector score as proof of misconduct, failing to verify references, and not documenting the role of AI in manuscript preparation are frequent errors that lead either to wrongful accusations or to avoidable retractions. Non-native English authors can be disproportionately affected by false positives; equitable policies must account for these biases in detector performance.

Conclusion and Next Steps

AI will continue to change scholarly workflows. Exposure of AI-origin content depends less on a single tool and more on an ecosystem: the combination of detector technologies, editorial policies, human review, and author transparency. Researchers should treat generative AI as a powerful drafting aid that requires verification and explicit disclosure. Editors and institutions should deploy detectors thoughtfully, pair them with human checks, and adopt fair investigation procedures.

For authors seeking practical support, professional [manuscript editing](#) can help ensure language clarity while documenting human revision and accountability; Enago's [manuscript editing](#) and [Responsible AI resources](#) provide guidance on disclosure and ethical use in academic writing. These services can help researchers present manuscripts that meet journal expectations and reduce the risk of procedural issues that can lead to retraction. Consider using such support to align submissions with publisher policies and to strengthen the human-authorship record.

Category

1. Reporting Research

Date Created

2026/01/09

Author

editor