

# The Detective That Lies: How Authentication Tools Create Their Own Fiction of Truth

Weekly Analysis — <https://ainews.social>

The pitch is irresistible: a small browser extension, a paid API, a checkbox in a learning-management system, and the question "did a human write this?" gets an answer. A confidence score. A green check or a red flag. The question is ancient — *cui bono*, who really speaks here — but the apparatus is new, and the apparatus, we are told, is impartial. It looks at the text, or the image, or the audio file, and it returns a verdict. It is a detective. It does not get tired, it does not have favorites, and it does not, the marketing materials imply, lie.

It does, though. Not in any dramatic sense — no rogue AI bending the truth — but in the more mundane and more dangerous sense that any classifier lies: it produces an output that looks like a finding when it is actually a guess, dresses statistical correlation in the costume of forensic certainty, and offloads onto the user the burden of knowing when to disbelieve it. Stanford's most recent stocktaking of the field, [15], catalogs the proliferation of detection and provenance tooling alongside steady evidence that the underlying classifiers remain brittle, gameable, and unevenly accurate across the populations they claim to serve. The previous year's edition, [16], pointed at the same gap with a useful pile of empirical studies on detector performance: detectors produce numbers, those numbers look like measurements, and the numbers are wrong often enough that institutional reliance on them is a policy choice masquerading as a technical one.

This essay is about that masquerade. It is about what authentication tools — AI text detectors, image-provenance systems, watermarks, the C2PA standard, the new generation of "deepfake spotters" — actually do versus what their vendors claim they do, and about the strange recursion at the heart of the whole enterprise: we are using AI systems to police the outputs of AI systems, on the theory that the second AI knows something definite about the first. It mostly does not. Meredith Broussard's still-useful term for the genre — what she called [16] — applies here with particular force: a tool can be statistically sophisticated and substantively wrong at the same time, and the sophistication is part of how the wrongness gets through.

[15] The 2026 AI Index Report

[16] HAI\_AI-Index-Report-2024

[16] Artificial Unintelligence

## *What the Detector Is Actually Doing*

Strip the marketing language away from a typical AI-text detector and what remains is a classifier trained to estimate the probability that a given passage was produced by a large language model. The features it relies on are statistical regularities — perplexity, burstiness, token-frequency profiles, stylometric fingerprints — that distinguish the average output of GPT-class models from the average output of human writers. This is real signal, in the aggregate. It is not, however, evidence about any particular document. The detector cannot interview the author. It cannot watch the keystrokes. It can only ask: does this text statistically resemble the texts I was trained to recognize?

The vendors know this. Read carefully through the educator-facing documentation OpenAI maintains and you find the company itself telling teachers, in plain language, that detectors are unreliable for individual cases — see the [5] and the parallel guidance in [3], where the disclaimer is buried but unambiguous. OpenAI itself withdrew its own detector in 2023 because the false-positive rate could not be brought down to a level at which the tool was safe to use as evidence. The third-party detectors that flooded the gap have not solved this problem. They have, in many cases, simply stopped publishing their error rates.

The structural reason is worth dwelling on. A detector trained to recognize "AI-ness" is recognizing a moving target. Every time a model is updated — and at the cadence documented across the corporate-tooling ecosystem, including the rolling release plans visible in [12] and Microsoft's broader [17] overview — the statistical fingerprint shifts. Detectors trained on yesterday's outputs degrade against today's. The arms race is not a metaphor; it is the literal mechanic. Each new model release is, from the detector's point of view, a distribution shift, and distribution shifts are precisely the regime in which classifiers fail silently. You don't get a warning that the tool is now wrong. You get the same confident percentage, just attached to a different ground truth.

And then there is the gameability. The CyberArk research write-up, [6], is about a different attack surface — bypassing safety filters rather than detection — but the underlying lesson generalizes: classifiers that ride on surface features are gameable by anyone willing to perturb the surface. A student who runs AI-generated text through a paraphraser, or simply asks the model to "write in a more bursty, idiosyncratic, error-prone style," collapses the detector's signal. The detector then outputs "human" with high confidence, because that is what the features now look like. The student who wrote their

[5] Educator FAQ | OpenAI Help Center

[3] ChatGPT for Teachers - OpenAI Help Center

[12] Prise en main de la 1re vague de lancement 2026 pour les offres Copilot

[17] What is Microsoft 365 Copilot?

[6] Jailbreaking Every LLM With One Simple Click

paper from scratch in a clean, fluent voice — say, a non-native English speaker who has internalized the rhythms of a textbook — gets flagged. This is not a hypothetical pattern; it is the documented one.

### *The Provenance Stack and Its Weakest Link*

If statistical detection is the weak option, the strong option is supposed to be provenance: don't try to spot AI content after the fact, sign it at birth. Watermark the model's outputs, embed cryptographic content credentials in the file, build a chain of custody from generator to viewer. The C2PA standard, the various invisible-watermark schemes, and the model-card disclosures shipped with open systems like the ones documented across the Hugging Face stack — see, for instance, [14] and its sibling [13] — represent serious technical work. The cryptography is real. The file formats are real. None of that is the problem.

The problem is that provenance is an ecosystem. A signature only matters if everyone signs, every platform preserves the signature, every viewer checks it, and every renderer refuses to strip it. Each of those is-ses is a point of failure. Take a watermarked image, screenshot it, post the screenshot — watermark gone. Take a C2PA-signed video, run it through the standard re-encode that every social platform applies to uploads — manifest gone. Use one of the many open-source generators that do not implement the standard at all — never signed in the first place. The chain is only as strong as its most negligent link, and the link with the most users is almost always the most negligent.

This is the deeper sense in which the whole authentication apparatus depends on adoption ecology, and adoption ecology, for AI tooling, is now thoroughly captured by a handful of vendors. The corporate productivity stack has consolidated around two players, Microsoft and Google, whose deployment guides — see [9] and [7] — describe seamless rollout into hundreds of millions of seats. The decision about whether enterprise-generated content carries content credentials by default is, in practice, a decision that two companies will make, and they will make it on a calendar shaped by their own product roadmaps. The user does not get a vote. Neither does the institution that has to adjudicate disputes about what was generated where.

Kate Crawford's argument in [16] — that the political economy of the technology is inseparable from its technical functioning — is exactly the lens this stack demands. A provenance system administered by the same firms whose models produce the content being authenti-

[14] Stable Diffusion 2 - Hugging Face

[13] Safe Stable Diffusion · Hugging Face

[9] Microsoft 365 Copilot - FastTrack

[7] Las funciones de IA de Gemini ahora están incluidas en las ...

[16] The Atlas of AI

cated is not a neutral check; it is an extension of the firm’s authority into the act of judgment itself. The detector and the generator share a vendor. When the vendor’s classifier confidently says ”this is human,” it is, in the end, the vendor confirming the vendor.

### *The False Positive Has a Name*

Detection error in the abstract is a statistical concept. Detection error in practice is a person. A student handed a zero on a paper they actually wrote. A freelance writer whose contract is terminated because the client ran the deliverable through a tool. A defendant whose alibi video is flagged as synthetic by an ”AI-spotter” that has no published validation against the kind of compression and lighting conditions present in real surveillance footage. These are not edge cases. They are the modal output of a system whose false positive rate, even at vendor-reported figures, multiplied across millions of submissions, produces tens of thousands of accusations per week.

The disparate-impact dimension here is the part the vendors most consistently elide. Detectors trained predominantly on native-English text from particular registers tend to flag non-native writers and writers from outside the training distribution at higher rates. Stanford’s catalog of empirical studies in the responsible-AI appendix to the 2024 edition of [16] makes this explicit, and the pattern reappears across independent evaluations: the tool’s ”objectivity” is calibrated to a center, and writers off-center pay the cost. Ruha Benjamin’s framing in [16] — that systems can encode discrimination through ostensibly neutral mechanisms, what she calls the New Jim Code — is the more general statement of the same phenomenon. A classifier does not need a discriminatory designer to produce discriminatory outcomes; it only needs an unrepresentative training distribution and a deployment context that treats its outputs as evidence.

Virginia Eubanks’ [16] names the further turn: tools of this kind are most consequential not where they are most accurate but where their targets have the least recourse. A wealthy professional flagged by a detector calls a lawyer; a community-college student on financial aid loses the semester. The same percentage error has different gravity at different points in the social structure, and authentication tools are deployed disproportionately at the points where the gravity is highest.

[16] HAI\_AI-Index-Report-2024

[16] Race After Technology

[16] Automating Inequality

### *What the Detector Says It Sees*

There is a second-order problem, harder to see, that runs underneath all of this. The detector is not just sometimes wrong; it is also, in a deeper sense, about a different question than the one its users think they are asking. The user wants to know: did a human do this work? The detector answers: does this text statistically resemble model output? Those are not the same question, and the gap between them is where the fiction lives.

Janelle Shane’s accessible treatment of how AI systems actually behave — in [16] — is full of examples of classifiers confidently labeling things based on features that have nothing to do with the labels users imagine. A model trained to identify tanks learns to identify the time of day; a model trained to detect skin cancer learns to detect rulers in the photo; a model trained to spot AI text learns to spot fluency, polish, and lack of idiosyncrasy, and labels them ”AI.” The detector is honest about its features in a sense — those really are the features it uses. It is dishonest, or at least misleading, about what those features mean. Polish is not authorship. Idiosyncrasy is not honesty. A clean prose register is not a confession.

This is the recursion the marketing language obscures. Authentication tools are AI systems applied to a task that used to be human judgment — the judgment of whether a piece of work was done by the person claiming credit for it. They inherit, as systems of this kind always do, the brittleness, opacity, and confident-seeming wrongness of the larger family. They are not exceptions to the critique of AI; they are instances of it. When a vendor markets a detector as a solution to the AI problem, the vendor is selling more of the problem dressed as the cure.

Broussard’s diagnosis in [16] is the right one to keep close: the systems are not stupid in any simple sense, but they are limited in ways that their interfaces conceal. The interface presents a verdict; the underlying system has a guess. The user reads the verdict; the institution acts on it. By the time anyone goes back to ask what the guess was actually based on, the consequence has already happened.

### *Centralization, Lock-In, and the Vendor as Judge*

Step back from the individual tool to the market it sits in. The detection-and-provenance economy is consolidating along the same lines as the rest of the AI tooling stack. Enterprise educational deployments increasingly route through Microsoft, whose [10] reference

[16] You Look Like a Thing and I Love You

[16] Artificial Unintelligence

[10] Microsoft Copilot in education - M365 Education

documentation describes a baseline that includes generation, agent-building through [11], and the licensing apparatus laid out in [1]. The companion stack from Google, with Gemini features now bundled into Workspace per the support page [8], produces the same effect from the other side: AI generation, AI detection, AI provenance, all administered by the same handful of platforms.

The political consequence of this consolidation is rarely named in the product literature. When the platform that produces the content also adjudicates whether content is authentic, the platform has acquired a power that no individual institution — no university, no court, no newsroom — can effectively contest. The institution that wants to challenge a detector’s verdict cannot, in any practical sense, audit the model. It cannot retrain the classifier. It cannot demand a confidence interval. It can, at best, switch vendors, which means trading one opaque judge for another. The lock-in is not just commercial; it is epistemic.

This is also where the cost structure starts to bite. Detection-as-a-service is priced per query, per seat, per institutional license. The Microsoft FastTrack documentation and the various Copilot SKUs walk through a cost model in which the institution pays for the privilege of accusing its own members. The contract pre-commits the institution to the tool’s framing of the question — there is no line item for “send the case to a human reader who can read.” That option has been priced out of the budget. The institution, having paid, will tend to use what it bought.

There is a wider point about labor here that the productivity-tool documentation makes almost accidentally. The same vendor stack that markets generation as a productivity gain — see the deployment guide for [2] and Amazon’s [4] — also markets detection as a way to police the outputs of that productivity. The gain and the policing are sold to the same buyer, by the same seller, in the same contract. The buyer ends up paying for both sides of a tension the vendor created.

### *The Question the Tool Cannot Answer*

A careful reader at this point will notice that I have not argued for abandoning detection and provenance tooling altogether. That would be the wrong conclusion. There are real uses for these systems — at scale, in aggregate, for triage, for trend-spotting, for forensic context where the result is one input among many. The argument is not that the tools have no signal; it is that the signal is being asked to bear more weight than it can.

[11] Microsoft Copilot Studio

[1] Assign user licenses and manage access - Microsoft Copilot Studio

[8] Lo mejor de la IA de Google ahora se incluye en las suscripciones a ...

[2] Casos de uso de asistentes de IA generativa en el desarrollo de ...

[4] CodeWhisperer Documentation

The crucial move for a careful adopter is to insist on a distinction the marketing language is built to dissolve. There is a difference between a tool that says “this document has features that are statistically associated with model output” and a tool that says “this document was generated by AI.” The first is what the classifier actually computes. The second is what users — and, more importantly, institutions — read off the screen. The translation from the first to the second is the place where the fiction enters, and it is not a translation the vendor performs in the contract; it is performed silently, by the interface, in the act of presenting a percentage as a verdict.

The questions a careful adopter should be asking, and rarely asks, are these: What is the false-positive rate on my population, not on the vendor’s benchmark? What is the disparate-impact analysis across the demographic groups represented in my user base? What is the recourse process when the tool is wrong, and who bears the burden of proof? What happens to my detector’s accuracy when the underlying model — over which I have no control — is updated next quarter? When the contract ends, do I retain the audit logs, or does the vendor? These are not gotcha questions; they are the basic due-diligence questions for any forensic instrument, and the answers, in the current detection market, are almost universally unsatisfying.

The 2026 AI Index data on responsible AI adoption that Stanford documents in [15] shows the predictable pattern: rapid deployment of AI tools, including detection and provenance tools, has run far ahead of the evaluation and governance infrastructure that would let users trust them responsibly. The institutional buyer is being asked to trust the tool because the tool exists, the vendor is reputable, and the alternative — slow human review — has been priced and timed out of viability. None of that is a reason to trust the tool. It is a reason to feel pressed into trusting it, which is a different and more dangerous condition.

[15] The 2026 AI Index Report

### *The Detective and the Story*

There is something worth saying, at the end, about the metaphor that the whole apparatus runs on. The tool is sold as a detective: it looks at evidence, it follows the clues, it solves the case. But a detective in fiction works because the detective is a character — a mind that holds the story together, that knows when a piece of evidence does not fit, that can revise the theory when the new fact arrives. The classifier has none of those capacities. It does not hold a story. It produces a number. The story is supplied by whoever reads the number, and that

person almost always supplies the story the number seems to want.

This is the deepest sense in which the authentication tool creates its own fiction of truth. The fiction is not in the algorithm. It is in the interface — in the moment when the percentage is rendered as a verdict, when the verdict is rendered as evidence, when the evidence is rendered as the basis for an institutional action that has consequences for an actual person. At each of those steps, something is being added that the underlying computation does not contain. The tool is not lying in the sense that it is producing a false output; it is producing exactly the output it was built to produce. The lie is in what we are reading the output to mean.

A pro-reader posture toward this market is, accordingly, neither credulity nor refusal. It is a habit of mind: read every detector verdict as a probability statement about a feature distribution, not as a claim about authorship. Read every provenance check as a claim about file metadata, not about truth. Read every vendor confidence score as a description of the vendor’s confidence, not yours. Treat the tool as one input among many, never as the input. Insist on human review at the points where consequences are highest. And remember, when the interface tells you with great certainty what happened, that the interface is the part the vendor designed to be most persuasive, which means it is the part that least deserves your trust.

Crawford’s image, in [16], of the technology as a planetary infrastructure of extraction is the right one to keep in view here. The detection and provenance economy is a continuation of that infrastructure, not an exception to it. It extracts judgments from people who used to make them; it concentrates the judgment-making capacity in a small set of vendors; it sells the judgment back to the institution that originally had the capacity itself. The detective that lies is not a bad detective. It is a detective whose interests — whose business model, whose product roadmap, whose contractual obligations — are not the user’s interests. Read the verdict accordingly.

[16] The Atlas of AI

The honest version of the authentication tool would come with a warning printed on every output: “This is a statistical estimate. It is wrong some of the time. The rate at which it is wrong on your population is unknown. The model that generated this estimate may be replaced next quarter. Do not use this output as the sole basis for any decision that affects a person.” No vendor will print that warning. The reader who wants to use these tools well will have to print it themselves, and read it every time the green check or the red flag appears on the screen. That mental discipline is, at this point in the technology’s life, the only authentication system that actually works

— the one the user runs in their own head, on the output of the one they bought.

### *References*

1. Assign user licenses and manage access - Microsoft Copilot Studio
2. Casos de uso de asistentes de IA generativa en el desarrollo de ...
3. ChatGPT for Teachers - OpenAI Help Center
4. CodeWhisperer Documentation
5. Educator FAQ | OpenAI Help Center
6. Jailbreaking Every LLM With One Simple Click
7. Las funciones de IA de Gemini ahora están incluidas en las ...
8. Lo mejor de la IA de Google ahora se incluye en las suscripciones a ...
9. Microsoft 365 Copilot - FastTrack
10. Microsoft Copilot in education - M365 Education
11. Microsoft Copilot Studio
12. Prise en main de la 1re vague de lancement 2026 pour les offres Copilot
13. Safe Stable Diffusion · Hugging Face
14. Stable Diffusion 2 - Hugging Face
15. The 2026 AI Index Report
16. The Atlas of AI
17. What is Microsoft 365 Copilot?