

# Beyond the Label: Why Provenance Literacy Demands a Renaissance of Skepticism

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

A small blue checkmark, a tiny "Generated by AI" tag in the corner of an image, a shimmer of invisible watermark embedded in pixels: this is the iconography our institutions have settled on to inoculate the public against synthetic media. The label is the talisman. The reasoning, when articulated, is touchingly nineteenth-century — a faith that disclosure, properly mandated, will produce informed citizens, the way a nutrition panel was supposed to produce slimmer ones. It has not. It will not. And the reason is not that the labels are poorly designed, though many are. The reason is that the problem we have decided to call a labeling problem is in fact a problem of skepticism, of the cognitive habits with which an ordinary person greets a piece of information arriving on their screen — and skepticism is not something a sticker can transmit.

Across the last two years, "AI literacy" has hardened into a policy keyword, deployed by ministries, philanthropies, and platforms as the answer to questions ranging from labor market disruption to electoral integrity. The phrase has the comforting quality of inevitability — of course we should teach people about the systems shaping their lives — but it conceals a battle over definition. Is AI literacy the ability to write better prompts for a chatbot? Is it the ability to distinguish a real photograph from a generated one? Is it the capacity to ask whether the question itself was the right one? UNESCO, in framing the synthetic media crisis, has begun to argue for the third reading, treating deepfakes not as a content-moderation problem but as a knowledge crisis that destabilizes the conditions of public reason [3]. That is the right altitude. But altitude is not curriculum, and the gap between the diagnosis and what is actually being taught — when literacy is taught at all — is the subject of this essay.

[3] Deepfakes and the crisis of knowing - UNESCO

The argument here is unfashionable in its modesty. Labels do not change behavior. Detection tools age out within months of release. Watermarks fracture under the most ordinary forms of manipulation — a screenshot, a recompression, a translation through a second model. The legal scaffolding being erected in the United States and

Europe to compel disclosure is necessary but not sufficient, and the most candid voices in that conversation will say so [19]. What citizens need is not a better label but an older virtue: the disciplined refusal to take any digitally mediated artifact at face value, regardless of whether a platform has blessed it as authentic. Call it provenance literacy. It is not a skill in the workshop sense. It is a posture, a habit of mind, a renaissance of the skepticism that pre-digital readers once brought to a pamphlet handed to them on the street.

### *The Label Fallacy*

Begin with what labels actually do. A content-provenance tag — whether the C2PA cryptographic standard adopted by major camera manufacturers and platforms, or the simpler “AI-generated” watermarks now mandated for political advertising in several jurisdictions — encodes a claim about origin that the viewer is invited to trust. The claim is brittle in three directions. It can be stripped: a screenshot of a labeled image carries none of its metadata forward. It can be spoofed: an authentic-looking provenance signal can be attached to a fabricated artifact by an adversary with access to the right keys. And, most damagingly, it can be ignored: even when present and accurate, the label competes for attention with the affective pull of the image itself, and the image almost always wins.

Microsoft’s documentation for its Azure AI Content Safety product is unintentionally clarifying on this point. The tooling is sophisticated — it scans for policy violations, jailbreaks, prompt injections, and a growing taxonomy of harms — but its design assumes a downstream developer making decisions about what to surface, not a downstream citizen making decisions about what to believe [20]. The locus of judgment has been quietly relocated from the reader to the platform. This is the architecture of managed information, and it is the opposite of what a literate public requires.

The labeling regime fails in a more subtle way as well. Synthetic content is now so pervasive — from generated stock imagery on corporate websites to hallucinated citations in chatbot outputs — that the absence of a label has stopped meaning what it used to mean. A study examining ChatGPT’s behavior on bibliographic tasks found that more than half of the references the system produced were either incorrect or wholly fabricated, presented in the confident grammar of a source that exists [2]. None of these fabricated citations carried a “this is invented” tag, because the system did not know it was inventing. Provenance, in such cases, is a category mistake — there is no source

[19] We Need Laws to Stop AI-Generated Deepfakes | Scientific American

[20] What is Azure AI Content Safety? - Azure AI services

[2] ChatGPT’s Hallucination Problem: Study Finds More Than Half Of AI’s ...

from which the artifact derives, only a probability distribution that the model has sampled.

This is why a reader trained only to look for the label is more dangerous than a reader trained to look for nothing at all. The first reader has learned to outsource judgment; the second still has it.

### *Competing Definitions, Competing Stakes*

The phrase "AI literacy" papers over at least three distinct projects, each with different beneficiaries. The first is operational fluency: the ability to use AI tools efficiently, write prompts, understand context windows, navigate the affordances of a chatbot. GitHub's own developer documentation on prompt engineering for Copilot Chat is a perfectly executed example of this register — practical, instrumental, oriented toward productivity [12]. The audience is the worker, and the goal is throughput. There is nothing wrong with this, but it is not literacy in any sense that would have been recognizable to Erasmus or to John Dewey. It is training.

[12] Prompt engineering for GitHub Copilot Chat

The second project is consumer protection: equipping people to recognize manipulated content, avoid scams, and refuse harmful interactions. UNESCO's communications on AI and disinformation belong here, framing literacy as the public's first line of defense in an environment where verification has become asymmetrically expensive — the lie travels in seconds, the correction takes hours and reaches a tenth of the audience [6]. This is closer to the older sense of media literacy, and it is genuinely necessary. But its theory of the citizen is largely defensive. The literate person, in this telling, is one who has not been deceived. That is a low bar.

[6] Inteligencia artificial y desinformación - UNESCO

The third project — the one that animates the strongest current scholarship and the one that ought to organize public investment — is the cultivation of a critical orientation toward all mediated information, AI-generated or not. UNESCO's *Think Critically Click Wisely* materials begin to gesture at this register, treating deepfake detection not as a forensic skill but as a window onto the broader epistemic conditions of digital life [14]. The literate person, here, is not the one who has not been deceived but the one who can articulate the conditions under which they would expect to be — who can model their own credulity. That is the renaissance posture: not certainty, but disciplined doubt.

[14] UNESCO Think Critically Click Wisely

These three projects do not coexist comfortably. A government framework that prioritizes operational fluency, as the recently finalized U.S. Department of Education priority on AI in schools largely does,

will produce a population well-equipped to deploy systems but poorly equipped to interrogate them [15]. A framework that emphasizes only consumer protection will produce a defensive citizenry, watching for fakes, missing the deeper questions about who is producing the genuine articles and to what end. Only the third register approaches what democratic life actually requires.

[15] The US Department of Education Just Finalized Its AI in Education ... - AEI

### *The Fragility of Technical Fixes*

Watermarking is the technocrat’s preferred answer, and its enthusiasts are not foolish. A robust, model-level watermark — a statistical fingerprint embedded in the token-level distribution of an LLM’s outputs, or in the high-frequency components of a generated image — has the elegant quality of being detectable by tools the public does not need to own. The user need not be skeptical; the system is skeptical on their behalf. This is the fantasy.

In practice, watermarks degrade. Recent work on frontier image generation models documents how trivially synthetic visual provenance signals can be removed, attenuated, or spoofed by adversaries with modest resources [4]. Translation through a second generative model — re-rendering an image at slightly different parameters — destroys the signal altogether. And the watermarks that survive such manipulation are precisely the ones that are most aggressive and therefore most likely to degrade output quality, creating a competitive disincentive for any single provider to adopt them robustly. The economic gradient runs the wrong way.

[4] Frontier Image Generation Models, Synthetic Visual ...

There is a second fragility, less often discussed: the watermark is only as trustworthy as the institutional chain that produces it. When a platform certifies that an image is “authentic,” what is being certified? That a particular camera produced it? That the metadata has not been altered since capture? That the depicted scene corresponds to a real-world event? These are radically different claims. The C2PA specification is technically rigorous about the first, partially rigorous about the second, and silent on the third — and it is the third that ordinary people care about. An image can be cryptographically authentic and substantively false. The watermark certifies the bytes, not the world.

Even the most cautious institutional voices acknowledge that detection is a losing race. A guide produced for K-12 schools by MIT Teaching Systems Lab notes that detection tools should be treated as supplementary at best, and that students must be taught to triangulate across multiple kinds of evidence rather than rely on any single

technical signal [9]. This is the right instinct, and it is not confined to schools. It is the instinct any reader of any age now requires.

The point is not that watermarking should be abandoned. It should not. The point is that watermarking is infrastructure, not literacy, and confusing the two has produced a generation of policy documents that mistake the existence of a technical countermeasure for the resolution of the social problem it was meant to address. Meredith Broussard's term for this category error — that we systematically overestimate what computational systems can do, and underestimate the human judgment that real situations require — applies directly here [14]. The label is technochauvinism in miniature.

### *Skepticism as Civic Muscle*

What, then, would a literacy actually adequate to this moment look like? It would begin by treating skepticism as a muscle rather than a skill — something that strengthens with regular use and atrophies in its absence, and that cannot be installed in a single training session.

A piece written for educators in Quebec on equipping young people to confront deepfakes makes this argument with unusual clarity, insisting that exposure to manipulated media must be repeated, contextualized, and discussed across years rather than dispatched in a single workshop [7]. The same logic applies to adults, and the absence of any analogous infrastructure for adult provenance literacy — outside the narrow precincts of journalism schools and a few civil society organizations — is one of the larger failures of the present moment.

A muscular skepticism has specific components. It includes asking, as a near-reflex, whether a piece of content can be triangulated against an independent source, and treating the absence of triangulation as informationally meaningful. It includes a working understanding of metadata — what kinds of claims metadata can and cannot support, how trivially it can be altered, and where the chain of custody breaks. It includes familiarity with adversarial techniques: how prompt injection works, how a model can be manipulated by the framing of a question, how the cheapness of producing plausible-sounding content shifts the burden of verification onto the reader. None of this is exotic. All of it is teachable. Almost none of it is being taught at any scale.

The asymmetry between what we know about AI's effects and how we are responding is itself part of the story. A recent national study of youth and AI found that the impact of these tools on mental health depends acutely on context — who is using them, in what relational frame, with what alternatives available — and that broad

[9] PDF A GUIDE TO AI IN SCHOOLS - tsl.mit.edu

[14] Artificial Unintelligence

[7] Les « deepfakes » : Comment donner aux jeunes les moyens de lutter ...

claims about AI being good or bad for young people obscure more than they reveal [8]. Context is the precise thing that labels strip away. A disclosure tag tells you that a piece of content was generated; it does not tell you by whom, for whom, with what intent, against what alternatives. Skepticism is the practice of asking those questions even — especially — when the system has politely declined to answer them.

[8] New National Research Reveals How Context Shapes AI's Impact on Youth ...

### *Who Defines Literacy, and Whose Literacy Counts*

Definitions are political. The framework that gets adopted at the federal or supranational level shapes the curricula that follow, the funding that flows, and the metrics by which success is measured. It also determines whose competence is treated as the reference standard and whose is treated as deficit. This is not a neutral question, and the answers being offered are not neutral.

Consider how the discourse around AI in special education has been organized. Recent reporting documents that AI tools are being deployed at scale to draft Individualized Education Programs, with advocates raising sharp concerns about whether the parents and students most affected understand what these systems can and cannot do, what data they consume, and how their outputs should be weighted against human judgment [13]. The literacy required to participate meaningfully in such a process — to push back on a generated recommendation, to ask for the reasoning, to demand human review — is not the literacy taught by a one-hour module on prompt engineering. It is closer to the literacy a tenant needs to read a lease, or a patient needs to read a consent form. It is adversarial.

[13] Teachers Are Using AI to Help Write IEPs. Advocates Have Concerns

The same structural pattern recurs in legal aid, in benefits adjudication, in immigration proceedings, in content moderation appeals. Reporting from Spain on digital violence, including AI-generated abusive imagery and grooming, documents how the burden of recognition falls heaviest on those least equipped to bear it, and how the institutional response has lagged the threat by years [17]. In each of these domains, the stakes of provenance literacy are not abstract. They are the difference between an individual being heard and an individual being processed.

[17] Una violencia que nunca acaba: ciberacoso, 'grooming', contenido sexual ...

A human-rights-grounded analysis from Latin America makes this point in its sharpest form, arguing that the asymmetry between those who deploy AI systems and those who are subject to them generates a literacy gap that is itself a form of structural injustice [11]. The vendor and the regulator share, broadly, a vocabulary; the citizen does

[11] PDF Derechos Humanos E Inteligencia Artificial: Una Mirada Desde Los ...

not. Closing that gap is not a matter of producing better explainer videos. It is a matter of redistributing interpretive authority — of insisting that the standards by which a system’s claims are evaluated are accessible to the people the system is acting upon. Kate Crawford’s mapping of the supply chains, labor regimes, and energy flows that constitute AI underwrites this redistribution: literacy that does not include the political economy of the systems being labeled is literacy in name only [14].

[14] The Atlas of AI

### *The Cross-Domain Stakes*

Provenance literacy does not stay in its lane. It sits underneath every other category of AI-related public concern, because every other concern depends on the public’s ability to evaluate claims made by, about, and through AI systems.

Consider election integrity. Recent analysis of AI’s role in governing disinformation across democracies identifies a recurring pattern: the technical interventions that platforms deploy in response to electoral pressure (detection systems, labeling regimes, content takedowns) tend to migrate, over the course of a single election cycle, from emergency measures to permanent infrastructure — without the public deliberation that ordinarily accompanies the construction of permanent infrastructure [21]. A citizenry that cannot interrogate the provenance of these governance moves — who proposed the system, who audits it, who benefits from its silences — has effectively delegated the management of its public sphere to whichever coalition of platforms and ministries was in the room when the decision was made. Disclosure labels on individual posts are a sideshow. The deeper provenance question is about the labeling regime itself.

[21] When AI Governs  
(Dis)information: Five Lessons for  
...

Consider fraud at scale. A recent investigation traced a network of fake college websites — full domain spoofs, generated faculty photographs, fabricated course catalogues — operating to harvest tuition payments and student data from international applicants [5]. The defenses that prevent such fraud are not technical detectors; the spoofs are good enough that automated detection lags. The defense is a public habituated to verifying institutional identity through multiple independent channels — registrations, accreditations, third-party reporting — before parting with money. That habit is provenance literacy in its civic form, and the ease with which such networks proliferate suggests how thinly the habit is currently distributed.

[5] Inside a Network of Fake College  
Websites

Consider the routine epistemic environment of work. A growing body of research on responsible AI in educational and professional

settings argues that the most consequential failures occur not at the dramatic edges (the deepfake, the hallucinated citation in a high-stakes filing) but in the accumulated drift of small unverified claims: the chatbot summary trusted because it was convenient, the generated image used because it was free, the model output cited because the model was confident [16]. The drift is the danger. It is also the place where literacy must be most embodied, because no individual instance of drift seems worth the friction of verification. Skepticism is the willingness to apply that friction anyway.

The Brookings analysis of AI's future for students argues, in a similar register, that the most meaningful interventions are not technological but cultural — sustained shifts in how questions are asked, how sources are evaluated, how confidence is calibrated [1]. This is correct, and it generalizes well beyond the cohort of students. The cultural shift is the literacy. Everything else is scaffolding.

### *A Renaissance of Skepticism*

The word "renaissance" is ambitious, and it is meant to be. The original European renaissance was not, despite its self-mythology, a discovery of new texts; it was a recovery of habits of reading that had survived in attenuated form through the medieval period and that, when reactivated, made possible a new relationship to authority. The reader of a Latin manuscript in 1480 began to ask, with new seriousness, who copied this, from what exemplar, with what corruptions, in whose service. The questions were not novel. The systematic application of them was.

The provenance literacy this moment requires is structurally similar. The questions — who made this artifact, from what materials, in whose service, with what occlusions — are old questions. They are the questions a careful reader has always brought to a pamphlet, a newspaper, a photograph, a film. What has changed is that the volume and verisimilitude of artifacts has overwhelmed the casual versions of these habits, and the casual habits are what most readers have. The labels and watermarks now being deployed at scale are an attempt to substitute for those habits with a technical signal. The substitution will not hold.

What will hold — what will scale, what will compound — is a public practice of asking the questions out loud, in the company of others, repeatedly, in the contexts in which the artifacts actually appear. A guide produced for evaluating evidence on AI in school settings makes a related point about practice as a precondition for judgment, arguing

[16] Towards responsible artificial intelligence in education: a systematic ...

[1] AI's future for students is in our hands - Brookings

that decisions about AI's role cannot be delegated to evidence-grading rubrics applied once and then filed away [18]. The judgment must be made and remade. The same is true at the level of the citizen.

This essay has avoided the temptation to issue a checklist, because checklists are precisely the form provenance literacy has been failing to take. A checklist gets memorized and then ignored; a habit gets internalized and applied without remembering the rule. What can be said, by way of orientation, is that the literate person of the next decade will be the one who treats every interface as a claim, every artifact as an argument, and every label as a piece of evidence to be evaluated rather than a verdict to be accepted. They will know that detection lags generation, that watermarks are infrastructure, and that the institutions producing both have interests not identical to their own. They will recognize, as the most candid literacy frameworks now do, that the boundary between AI-generated and human-generated content is becoming less informationally useful than the older boundary between content that has been verified and content that has not [10].

The label fades into the background; the verifying intelligence comes forward. That is the shift. It is not glamorous, it cannot be deployed by a vendor, and it does not lend itself to keynote announcements. It is, however, the only literacy that has ever worked, and there is no reason to think the present moment will produce a shortcut where prior moments produced none. The work is the work. A renaissance, properly understood, is the moment a culture remembers this.

## *References*

1. AI's future for students is in our hands - Brookings
2. ChatGPT's Hallucination Problem: Study Finds More Than Half Of AI's ...
3. Deepfakes and the crisis of knowing - UNESCO
4. Frontier Image Generation Models, Synthetic Visual ...
5. Inside a Network of Fake College Websites
6. Inteligencia artificial y desinformación - UNESCO
7. Les « deepfakes » : Comment donner aux jeunes les moyens de lutter ...
8. New National Research Reveals How Context Shapes AI's Impact on Youth ...

[18] Understanding the Evidence Base on AI in K-12 Education | SCALE Initiative

[10] PDF AI Child Safety Final Report - d19ob9sqegt2wc.cloudfront.net

9. PDF A GUIDE TO AI IN SCHOOLS - [tsl.mit.edu](https://tsl.mit.edu)
10. PDF AI Child Safety Final Report - [d19ob9sqegt2wc.cloudfront.net](https://d19ob9sqegt2wc.cloudfront.net)
11. PDF Derechos Humanos E Inteligencia Artificial: Una Mirada Desde Los ...
12. Prompt engineering for GitHub Copilot Chat
13. Teachers Are Using AI to Help Write IEPs. Advocates Have Concerns
14. The Atlas of AI
15. The US Department of Education Just Finalized Its AI in Education ... - AEI
16. Towards responsible artificial intelligence in education: a systematic ...
17. Una violencia que nunca acaba: ciberacoso, 'grooming', contenido sexual ...
18. Understanding the Evidence Base on AI in K-12 Education | SCALE Initiative
19. We Need Laws to Stop AI-Generated Deepfakes | Scientific American
20. What is Azure AI Content Safety? - Azure AI services
21. When AI Governs (Dis)information: Five Lessons for ...