

On the Impossibility of Artifact-Based Truth: Indistinguishability, Undecidability, and the Probabilistic Geometry of Plausibility

Indistinguishability, Undecidability, and the Limits of Artifact-Based Truth

Author: Marcos Eduardo Elias

Affiliation: The Ramanujan Institute - Brazil (Mathematics, Computer Science, Cryptography)

Contact: available upon request

Abstract

We prove that perfect authenticity verification of digital or photographed artifacts is impossible under adversarial conditions commonly encountered in real-world human-machine systems. We present a unified impossibility theorem combining (i) observational indistinguishability (in the spirit of FLP-style arguments from distributed computing) with (ii) algorithmic undecidability (via a reduction to the Halting Problem and Gödelian diagonalization).

We show that when an adversary can replicate all computable observables of an artifact—including typography, interface layout, optical capture, and metadata—no verifier operating solely on the artifact can decide authenticity with certainty. We further demonstrate that any claim of a universal, perfect verifier implies the decidability of the Halting Problem, yielding a contradiction.

Finally, we propose a reconciliatory framework (Theorem 3), grounded in von Neumann’s asymmetry principles and Shannon’s separation of signal and semantics, establishing that authenticity can be bounded probabilistically only by introducing external asymmetries (trusted channels, cost, time, or intersubjective corroboration), never by artifact analysis alone.

Preliminaries and Model

Artifacts and Worlds

Let Σ^* denote the set of all finite binary strings.

An artifact $A \in \Sigma^*$ represents any digital object or capture (e.g., image bytes of a photographed screen).

Define two generative processes (worlds):

- R: the real world generator, producing artifacts according to distribution D_R over Σ^* .
- F: the fake world generator, producing artifacts according to distribution D_F over Σ^* .

A verifier is an algorithm $V: \Sigma^* \rightarrow \{\text{REAL}, \text{FAKE}\}$.

V may be deterministic or randomized.

Theorem 1 — Indistinguishability (Distributed-Systems Style)

Definitions

Let $\text{Feat}(A)$ be the (finite) set of computable observables extracted by V from A (e.g., pixels, noise statistics, typography, EXIF, optical artifacts).

We say R and F are observationally indistinguishable for V if:

For every artifact A in the support of D_R , there exists an artifact A' in the support of D_F such that

$$\text{Feat}(A') = \text{Feat}(A).$$

Equivalently, the induced distributions over Feat are equal:

$$D_R \circ \text{Feat} = D_F \circ \text{Feat}.$$

Statement of Theorem 1

Theorem 1 (Indistinguishability).

If $D_R \circ \text{Feat} = D_F \circ \text{Feat}$, then no verifier V can distinguish REAL from FAKE with certainty.

Formally, for any V ,

$$P_{\{AD_R\}}[V(A) = \text{REAL}] = P_{\{AD_F\}}[V(A) = \text{REAL}].$$

In particular, no perfect verifier exists.

Proof

1. Assume, for contradiction, that there exists a perfect verifier V .
2. Then $P_{\{AD_R\}}[V(A) = \text{REAL}] = 1$ and $P_{\{AD_F\}}[V(A) = \text{REAL}] = 0$.
3. Since $D_R \circ \text{Feat} = D_F \circ \text{Feat}$, V sees identically distributed inputs in both worlds.
4. Therefore $P_{\{AD_R\}}[V(A) = \text{REAL}] = P_{\{AD_F\}}[V(A) = \text{REAL}]$.
5. This contradicts step (2).
6. Hence, no such V exists. QED.

Interpretation

This mirrors classical indistinguishability proofs in distributed computing:

if two global states induce the same local view, no process can decide differently without error.

Theorem 2 — Undecidability (Gödel / Turing Style)

Computational Foundation

We work within Peano Arithmetic (PA), a recursively axiomatizable, consistent system sufficient to encode Turing computation.

Let $\langle M, x \rangle$ be a Gödel encoding of a Turing machine M with input x .

Define $\text{Halt}(e) =$ “the machine encoded by e halts”.

It is known that $\text{Halt}(e)$ is undecidable.

Authenticity as a Decidable Predicate

Assume a property $P(A)$ meaning “ A is an authentic artifact of event E ”.

A universal perfect verifier would be a total algorithm V such that:

For all $A \in \Sigma^*$:

$V(A) = \text{REAL}$ if and only if $P(A)$ holds.

Reduction to the Halting Problem

Construct, for each e , an event E_e defined as:

“Event E_e occurs if and only if the machine encoded by e halts.”

Define an artifact A_e as “a purported capture of event E_e ”.

Then:

- If $\text{Halt}(e)$ is true, E_e occurs, and there exists a real artifact A_e .
- If $\text{Halt}(e)$ is false, E_e never occurs, so any A_e is fake.

If V exists, we can decide $\text{Halt}(e)$ by computing $V(A_e)$.

This contradicts the undecidability of Halt .

Statement of Theorem 2

Theorem 2 (Undecidability).

No total algorithm can perfectly decide authenticity for all possible artifact-defined events.

Gödelian Diagonalization (Intuition)

Define an artifact A^* that asserts: “ $V(A^*) = \text{FAKE}$ ”.

If $V(A^*) = \text{REAL}$, A^* lies; if $V(A^*) = \text{FAKE}$, A^* tells the truth.

Thus, V cannot be consistent on all inputs.

Unified Impossibility Theorem

Theorem (Unified Impossibility of Perfect Authenticity Verification).

For any verifier V operating solely on an artifact $A \in \Sigma^*$:

1. If the adversary can replicate all computable observables, authenticity is indistinguishable (Theorem 1).
2. If V is required to be universal and perfect, authenticity is undecidable (Theorem 2).

Therefore, perfect authenticity verification from artifacts alone is impossible.

Theorem 3 — Reconciliation via Asymmetry (von Neumann–Shannon)

Motivation

The impossibility above does not imply nihilism. It implies that verification must introduce asymmetry.

von Neumann’s Principle (Asymmetry)

John von Neumann emphasized that computation requires irreversibility or asymmetry.

A symmetric system cannot generate information gain.

Applied here: if the verifier and adversary share identical capabilities over the artifact, no decision is possible.

Shannon’s Separation

Claude Shannon showed that signal and semantics are separable:

information theory alone cannot certify meaning or truth.

Thus, authenticity cannot be inferred from signal statistics alone.

Statement of Theorem 3

Theorem 3 (Asymmetry-Bounded Verification).

Authenticity can be bounded probabilistically only by introducing at least one external asymmetry, such as:

- a trusted external channel,
- a cost asymmetry (economic, temporal, reputational),
- a cryptographic commitment outside the artifact, or
- intersubjective corroboration across independent agents.

Without such asymmetry, artifact-based verification collapses to Theorems 1 and 2.

Consequence

Verification shifts from a binary decision problem to a field of plausibility over time, cost, and social consistency.

Implications

- No cryptographic primitive can rescue authenticity if the adversary controls the artifact pipeline.
- Human–machine hybrids defeat purely technical detectors.
- The correct object of verification is not the artifact, but the event and its consequences.

Conclusion

Authenticity is not a property of artifacts alone.

It is an emergent property of asymmetry, cost, time, and intersubjective validation.

Any system claiming perfect artifact-based verification contradicts either indistinguishability or undecidability.

Glossary

- Artifact: Any finite digital representation (Σ^*).
- Indistinguishability: Equality of observable distributions.
- Verifier: Algorithm deciding REAL vs FAKE.
- Halting Problem: The undecidable problem of determining whether a computation halts.
- Asymmetry: Any resource or constraint not shared symmetrically by adversary and verifier.

A Mainstream Architecture for Authenticity as Time-Evolving Plausibility

Bayesian Filtering + Likelihood Ratios under Adversarial Indistinguishability

Core objects

Let E be a binary event:

$E = 1$ means “the claim is true” (e.g., “the SEI memo existed and was routed at time t_0 ”)

$E = 0$ means “the claim is false”

You observe evidence over time. At discrete times $t = 1, 2, 3, \dots$ you receive observation o_t .

Define the observation history up to time t :

$$O_{1:t} = (o_1, o_2, \dots, o_t)$$

The plausibility field is the posterior:

$$p_t = P(E=1 \mid O_{1:t})$$

This is the single canonical quantity reviewers respect.

The impossibility shows up as symmetry in likelihoods

Theorem 1's content, in Bayesian language:

If the adversary can make the distribution of observations identical under $E=1$ and $E=0$, i.e.

$$P(o_t \mid E=1, O_{1:t-1}) = P(o_t \mid E=0, O_{1:t-1}) \text{ for all } t,$$

then the observations carry zero information and the posterior never moves:

$$p_t = p_0 \text{ for all } t.$$

So "symmetry" here means:

the likelihood ratio is 1 forever.

Define the (conditional) likelihood ratio at time t :

$$LR_t = P(o_t \mid E=1, O_{1:t-1}) / P(o_t \mid E=0, O_{1:t-1})$$

If $LR_t = 1$ always, you cannot learn.

This is the cleanest mainstream restatement of the impossibility.

Where asymmetry is injected

You inject asymmetry by adding observation channels whose conditional distributions differ under $E=1$ vs $E=0$. That is, you force $LR_t \neq 1$ sometimes.

A mainstream way to model this is to decompose each observation into multiple channels:

$$o_t = (x_t^{(1)}, x_t^{(2)}, \dots, x_t^{(k)})$$

Examples (purely abstract categories):

- $x_t^{(A)}$: physical/optical trace (camera noise/moiré etc.)
- $x_t^{(B)}$: narrative-temporal consistency (calendar-time constraints)
- $x_t^{(C)}$: intersubjective corroboration (independent reports)
- $x_t^{(D)}$: system-side attestation (server log / signature) if available

You are not claiming any new crypto; you are saying: only channels with different likelihoods create learning.

If the adversary can fully simulate a channel, then that channel has $LR \approx 1$ and is informationally useless. If the adversary cannot simulate it without high cost or coordination, LR deviates from 1 and you can update.

Sequential Bayes update (time-evolving plausibility)

Start with a prior $p_0 = P(E=1)$.

The odds form is simplest and copyable:

$$\text{Odds}_t = p_t / (1 - p_t)$$

Sequential update:

$$\text{Odds}_t = \text{Odds}_{\{t-1\}} * LR_t$$

Equivalently, in log-odds:

$$L_t = \log(\text{Odds}_t)$$

$$L_t = L_{\{t-1\}} + \log(LR_t)$$

So the plausibility field is a cumulative log-likelihood ratio process.

This is the exact mainstream substrate used in:

- sequential probability ratio tests (SPRT, Wald)
- Bayesian filtering
- online detection

A robust mainstream specialization: Adversarial mixture model

Reviewers will immediately ask: “But under fake, the adversary adapts.”

You handle this with a standard robust Bayesian approach:

Under $E=0$ (fake), observations are generated from a mixture:

$$P(o_t | E=0) = (1 - \alpha) * P_{\text{benign}}(o_t | E=0) + \alpha * P_{\text{adversarial}}(o_t)$$

where $\alpha \in [0,1]$ is the probability the adversary actively crafts the artifact.

This is mainstream (robust statistics, Huber-style contamination, adversarial priors).

Then LR_t becomes:

$$LR_t = P(o_t | E=1) / [(1 - \alpha) P_{\text{benign}}(o_t | E=0) + \alpha P_{\text{adversarial}}(o_t)]$$

You do not need to know $P_{\text{adversarial}}$ exactly; you can bound it. Which leads to:

- worst-case LR bounds
- conservative posterior bounds

Computing plausibility fields when you can't model everything

A practical, still mainstream approach is to define evidence features $\varphi_i(o_t)$ and use a log-linear likelihood ratio model:

$$\log(LR_t) = \sum_i w_i * \varphi_i(o_t)$$

This is equivalent to using a generalized likelihood ratio with features; it's also equivalent to logistic regression under certain assumptions, which is mainstream enough.

Key: reviewers accept this if you are explicit: it's a model, not a proof of reality.

The process becomes:

$$L_t = L_{t-1} + \sum_i w_i \varphi_i(o_t)$$

And p_t is recovered from odds:

$$p_t = \text{Odds}_t / (1 + \text{Odds}_t) = 1 / (1 + \exp(-L_t))$$

Decision policy: Wald-style stopping (no binary truth claims)

You don't output "REAL" as a metaphysical truth. You output actions based on thresholds:

Choose two thresholds $A > 0$, $B < 0$ on log-odds L_t :

- if $L_t \geq A$: act as if E is true (accept / escalate)
- if $L_t \leq B$: act as if E is false (reject / de-escalate)
- else: wait for more evidence

This is exactly SPRT logic, a classical result (Wald). It's "sacred" mainstream.

Mapping the impossibility theorem into the architecture

- Theorem 1 corresponds to the case where for every channel, $LR_t = 1$ (or very close), so L_t is a martingale with zero drift and never separates.
- Theorem 2 corresponds to demanding a universal verifier that decides E for arbitrary event specifications; the architecture refuses that demand and instead (i) restricts E to operational schemas, and (ii) provides probabilistic bounds and stopping rules.

So the architecture is not trying to refute the theorems; it is the constructive layer above them.

What counts as "asymmetry" in this mainstream framing?

Asymmetry is any evidence stream where:

$$P(o_t | E=1) \text{ is materially different from } P(o_t | E=0)$$

and where the adversary cannot force them equal without incurring cost or coordination that itself becomes observable.

In short:

Asymmetry \Leftrightarrow expected log-likelihood ratio drift is nonzero.

Define drift:

$$\mu = E[\log(LR_t) \mid E=1] \text{ and } \nu = E[\log(LR_t) \mid E=0]$$

Learning happens when $\mu > 0$ and $\nu < 0$ (on average).

Perfect indistinguishability implies $\mu = \nu = 0$.

This drift criterion is classical and reviewer-proof.

Minimal, reviewer-safe “pipeline” description

1. Define event schema E (operational, testable).
2. Set prior p_0 (explicit).
3. For each evidence arrival o_t , compute feature vector $\varphi(o_t)$.
4. Compute incremental log-likelihood $\Delta L_t = \sum_i w_i \varphi_i(o_t)$.
5. Update $L_t = L_{\{t-1\}} + \Delta L_t$, hence p_t .
6. Apply Wald thresholds to decide accept/reject/continue.

That’s it. No tables, no magic.

Now I’ll do the step in a way that is (1) technically mainstream, (2) operational, (3) explicitly not a recipe for forging, and (4) uses only copy-paste symbols.

We’ll take a concrete claim and show how $p(t)$ evolves as evidence arrives, using a small set of features φ_i that (a) do not rely on typos, and (b) are grounded in physics, timing, and corroboration.

Concrete example: the memo photo claim

Event (binary):

$E = 1$: “A SEI/BNDES internal memo ‘Nota à Diretoria nº 05/2026 – PRESI’ existed in SEI, was signed by the Presidency, and routed to directors at time $t_0 = 2026-02-03$ (local), and the photo A is a capture of that real screen session.”

$E = 0$: “The photo A does not correspond to such a real SEI session (fabricated, reenacted, composited, or otherwise not genuine to the claimed event).”

Observation stream:

At $t=1$ we receive the artifact A (a photo of a screen).

Then over time we receive additional observations o_2, o_3, \dots (corroboration, context consistency, etc.)

We compute the plausibility field:

$$p_t = P(E=1 \mid O_{1:t})$$

Use log-odds:

$$L_t = \log(p_t / (1 - p_t))$$

Update rule:

$$L_t = L_{t-1} + \Delta L_t$$

where ΔL_t is the incremental log-likelihood evidence.

The evidence model (minimal but rigorous)

We'll define a feature vector $\varphi(o_t) = (\varphi_1, \varphi_2, \dots, \varphi_m)$ for each observation.

Then model:

$$\Delta L_t = \sum_i w_i * \varphi_i(o_t)$$

Interpretation:

- If φ_i supports authenticity, it is positive and pushes L_t up.
- If φ_i supports fabrication, it is negative and pushes L_t down.
- Weights w_i encode how diagnostic that feature is (in log-likelihood units).

This is a standard log-linear approximation to $\log(LR_t)$. It is exactly the kind of thing a reviewer accepts as a model.

Features for the artifact photo (t=1), ignoring typos completely

We define $\varphi_i(A)$ from physics and interface behavior. None of these are “spellcheck.”

φ_1 : Optical moiré plausibility score

Does the photo display a moiré/subpixel pattern consistent with a camera photographing a real display, with anisotropy and nonstationary interference?

- $\varphi_1 = +1$ if strongly consistent
- $\varphi_1 = 0$ if ambiguous
- $\varphi_1 = -1$ if inconsistent (too clean / too synthetic / pattern repeats unnaturally)

φ_2 : Sensor noise consistency score

Does noise behave like real phone sensor noise (chromatic noise in shadows, nonuniform variance, rolling-shutter micro artifacts)?

$\varphi_2 \in \{+1, 0, -1\}$ similarly.

φ_3 : Perspective + lens distortion coherence

Do straight UI lines show plausible slight barrel/pincushion distortion and perspective geometry consistent with a handheld capture?

$\varphi_3 \in \{+1, 0, -1\}$.

φ_4 : UI micro-consistency score (non-textual)

This is about “tiny UI truths,” e.g., scroll position, icon states, spacing of SEI elements, typical UI density.

Importantly: not “can a human copy the UI,” but “are there internal micro-constraints that are typically violated when reenacting.”

$\varphi_4 \in \{+1, 0, -1\}$.

φ_5 : “Effort plausibility” prior (cost-of-fake)

A Bayesian prior feature: given this artifact’s complexity and the claimant’s environment, is high-effort fabrication plausible?

This is not moral judgment; it’s adversary modeling.

$\varphi_5 \in \{+1, 0, -1\}$.

Now, define the first update:

$$\Delta L_{-1} = w_1 \varphi_1 + w_2 \varphi_2 + w_3 \varphi_3 + w_4 \varphi_4 + w_5 \varphi_5$$

Example weights (order-of-magnitude, not sacred):

$$w_1 = 0.7 \text{ (moiré is fairly diagnostic)}$$

$$w_2 = 0.5$$

$$w_3 = 0.4$$

$$w_4 = 0.6$$

$$w_5 = 0.3 \text{ (prior/cost cue, weaker)}$$

These are log-likelihood units; 0.7 means $LR \approx \exp(0.7) \approx 2.0$.

Start with a prior

Pick a conservative prior, since single-photo claims are common:

$$p_{-0} = 0.20$$

Then:

$$L_{-0} = \log(0.20 / 0.80) = \log(0.25) \approx -1.386$$

5. Apply $t=1$ update (artifact photo arrives)

Suppose a neutral-to-positive forensic read of the photo yields:

$$\varphi_1 = +1 \text{ (moiré looks real)}$$

$$\varphi_2 = +1 \text{ (noise looks sensor-real)}$$

$\varphi_3 = +1$ (perspective coherent)

$\varphi_4 = 0$ (UI micro-consistency ambiguous)

$\varphi_5 = 0$ (we don't assume high/low effort yet)

Then:

$$\Delta L_1 = 0.7*(+1) + 0.5*(+1) + 0.4*(+1) + 0.6*(0) + 0.3*(0)$$

$$\Delta L_1 = 1.6$$

So:

$$L_1 = L_0 + 1.6 = -1.386 + 1.6 = 0.214$$

Convert back:

$$p_1 = 1 / (1 + \exp(-L_1)) = 1 / (1 + \exp(-0.214)) \approx 0.553$$

Interpretation: after a single strong-looking photo, we move from 20% to ~55%.

This is exactly how a sober system should behave: “interesting, not proven.”

Time-evolving evidence beyond the photo

We now define additional evidence channels that create asymmetry. Again, not typos.

At $t=2$: narrative-temporal coherence observation o_2

Example: does the claimed time “2026-02-03 18:47” match plausible internal decision cadence (end-of-day signature), time zone, and the SEI notice format?

Define feature:

φ_6 : Temporal plausibility

+1 if coherent, 0 if unclear, -1 if strongly incoherent.

Weight $w_6 = 0.4$ (useful but not decisive; humans can imitate time).

Suppose $\varphi_6 = +1$. Then $\Delta L_2 = 0.4$ and:

$$L_2 = 0.214 + 0.4 = 0.614$$

$$p_2 \approx 0.649$$

At $t=3$: independent corroboration arrives o_3

This is the big one. Define:

φ_7 : Independent corroboration count (bounded)

Let k be the number of mutually independent sources that report the same memo existence, with independence assessed structurally (different org units, different access paths).

We can model $\varphi_7 = k$, but cap it to avoid overconfidence.

Use a diminishing returns function:

$$\varphi_7 = \min(k, 3)$$

Weight $w_7 = 1.0$ per independent corroboration unit (very strong).

If we get $k=1$ (one credible independent confirmation), then:

$$\Delta L_3 = 1.0 * (1) = 1.0$$

$$L_3 = 1.614$$

$$p_3 \approx 0.834$$

Now it's "very likely," but still not absolute.

At $t=4$: contradiction arrives o_4

Example: a credible insider says: "Nota 05/2026 exists but the subject line differs materially," or "the numbering scheme that week was different."

Define:

φ_8 : High-grade contradiction flag

$\varphi_8 = -1$ if present, 0 otherwise.

Weight $w_8 = 1.2$ (contradictions from high-trust channels are strong).

If $\varphi_8 = -1$:

$$\Delta L_4 = 1.2 * (-1) = -1.2$$

$$L_4 = 0.414$$

$$p_4 \approx 0.602$$

Notice what happened: one strong contradiction can pull you back from 83% to ~60%.

That's a feature, not a bug. This is how you avoid being "narrative-captured."

At $t=5$: system-side attestation arrives o_5 (optional, if you ever have it)

For example: a verifiable SEI hash or server-side log snippet. This is the canonical asymmetry injection.

Define:

φ_9 : Verifiable attestation

$\varphi_9 = +1$ if present and independently verifiable, 0 if absent.

Weight $w_9 = 3.0$ (huge; this is what breaks Theorem 1 in practice).

If $\varphi_9 = +1$:

$$\Delta L_{.5} = +3.0$$

$$L_{.5} = 3.414$$

$$p_{.5} \approx 0.968$$

Now you are in “near-certain” territory.

The plausibility field as a function of time

We’ve computed p_t at discrete t . To make it a field over continuous time, define:

$p(t)$ piecewise constant between evidence arrivals or use smoothing.

But the point is conceptual:

- Without asymmetry channels, ΔL_t hovers near 0 (symmetry), $p(t)$ does not move.
- With at least one channel whose likelihoods differ, ΔL_t has drift, and $p(t)$ moves.

The drift condition (mainstream, clean):

Learning requires $E[\Delta L_t \mid E=1] > 0$ and $E[\Delta L_t \mid E=0] < 0$ for some evidence schedule.

Perfect fake-symmetry makes both expectations ≈ 0 .

Decision policy (Wald thresholds, no metaphysical “truth”)

Choose thresholds A and B on L_t .

Example:

$$A = +2.0 \text{ (accept/act)} \Rightarrow p \approx 0.881$$

$$B = -2.0 \text{ (reject)} \Rightarrow p \approx 0.119$$

Policy:

- If $L_t \geq A$: act as if E true (escalate, allocate resources, publish cautiously).
- If $L_t \leq B$: act as if E false (de-escalate, treat as fabrication).
- Else: keep gathering evidence.

This is literally classical SPRT logic.

What this shows, “seminally,” disregarding typos

A photo can move you from prior skepticism to moderate plausibility (20% \rightarrow 55%).

The decisive movement comes from asymmetrical channels, especially independent corroboration and verifiable attestations.

A single high-grade contradiction can undo a lot, which prevents “one-shot” capture.

The model is reviewer-proof: it admits uncertainty, quantifies updates, and embeds the impossibility theorem as the symmetry case $LR=1$.

I'll now:

- Assume all typos are corrected (completely ignore orthography).
- Work only with the weak material you actually have (full screenshot + photo + crops).
- Perform the MIT-style Bayesian plausibility computation, step by step.
- Do not give any instructions that help fabrication.
- Act as a hostile-but-honest reviewer.

I will explicitly flag what is observable, weakly observable, and not observable, because that discipline is what makes this credible.

Restating the claim (precisely)

Event E:

“A SEI/BNDES internal memorandum titled ‘Nota à Diretoria nº 05/2026 – PRESI’ existed, was signed by the Presidency, and was visible in SEI around $t_0 = 2026-02-03$; the photographed screen corresponds to that real session.”

We compute a plausibility field $p(t) = P(E = 1 \mid \text{evidence up to time } t)$.

Prior (reviewer-grade conservative)

Single-photo institutional leaks are common but often fake.

Set:

$$p_0 = 0.20$$

$$L_0 = \log(p_0 / (1 - p_0)) = \log(0.25) \approx -1.386$$

This is deliberately skeptical.

Feature set (typos ignored entirely)

I will score only features that a reviewer would accept as legitimate.

φ_1 — Optical moiré / subpixel interaction

Observability: weak but present.

From the phone photo (not the PDF mock):

- There is nonuniform aliasing across text lines.
- Moiré is directionally inconsistent, not periodic.
- JPEG compression overlays but does not erase it.

This is consistent with camera→display capture, not decisive but positive.

Score: $\varphi_1 = +0.5$ (not +1; weak evidence)

Weight: $w_1 = 0.7$

Contribution: +0.35

φ_2 — Sensor noise consistency

Observability: weak-to-moderate.

- Noise is chromatic, uneven, stronger in gray regions.
- No obvious GAN-style “pleasant noise”.
- Compression artifacts are plausible for WhatsApp/Telegram forwarding.

Score: $\varphi_2 = +0.5$

Weight: $w_2 = 0.5$

Contribution: +0.25

φ_3 — Perspective & lens distortion

Observability: moderate (full phone photo helps).

- Slight keystone distortion vertically.
- Top bar curvature consistent with phone lens.
- Not perfectly rectified.

Score: $\varphi_3 = +1$

Weight: $w_3 = 0.4$

Contribution: +0.4

φ_4 — UI micro-consistency (non-textual)

Observability: moderate.

Key points:

- SEI header spacing and breadcrumb layout are internally consistent.
- Button placement (“Assinar Documento”, “Voltar”, “Baixar”) is plausible.
- No impossible UI state is visible.

Important: this does not prove reality, but it removes a common failure mode.

Score: $\varphi_4 = +0.5$

Weight: $w_4 = 0.6$

Contribution: +0.3

φ_5 — Effort / adversary cost prior

Observability: weak, but not zero.

Given the level of UI fidelity + photo capture + distribution:

- This is not trivial, but not implausible either.
- A hostile reviewer will not give much credit here.

Score: $\varphi_5 = 0$

Weight: $w_5 = 0.3$

Contribution: 0

First update: artifact-only plausibility

Sum contributions:

$$\Delta L_1 = 0.35 + 0.25 + 0.4 + 0.3 = 1.30$$

Update log-odds:

$$L_1 = L_0 + \Delta L_1 = -1.386 + 1.30 = -0.086$$

Convert to probability:

$$p_1 = 1 / (1 + \exp(0.086)) \approx 0.478$$

Interpretation:

From a single photographed screen that looks physically real,

plausibility moves from 20% \rightarrow ~48%.

That is exactly the right order of magnitude.

Not belief. Not dismissal. Suspended plausibility.

Any system that jumps to 90% here is unserious.

Time evolution with weak but mainstream asymmetries

Now we simulate plausibility over time, not certainty.

t_2 — Temporal / procedural coherence (φ_6)

Evidence: the memo structure, time-of-day plausibility, routing order.

This is easy to fake, but mistakes are common.

Here, nothing jumps out as impossible.

Score: $\varphi_6 = +0.5$

Weight: $w_6 = 0.4$

$\Delta L_2 = +0.2$

$L_2 = -0.086 + 0.2 = 0.114$

$p_2 \approx 0.528$

t_3 — Independent corroboration (single source)

Suppose one structurally independent confirmation appears
(e.g., another insider acknowledges memo existence, not content).

Score: $\varphi_7 = +1$

Weight: $w_7 = 1.0$

$\Delta L_3 = +1.0$

$L_3 = 1.114$

$p_3 \approx 0.753$

This is the first moment where a reviewer would say:

“Likely true, but still reversible.”

t_4 — High-grade contradiction (possible!)

Suppose a credible contradiction appears
(e.g., numbering scheme mismatch, or timing inconsistency).

Score: $\varphi_8 = -1$

Weight: $w_8 = 1.2$

$\Delta L_4 = -1.2$

$L_4 = -0.086$

$p_4 \approx 0.478$

One contradiction almost cancels everything.

This is a feature, not a flaw.

What breaks symmetry (and what does not)

From this exercise:

- The photo alone cannot break symmetry (Theorem 1).
- Optical realism gives $\sim +1.3$ log-odds max, not more.
- Independent corroboration is the dominant asymmetry.
- Server-side attestations (if ever available) would dominate everything.

No typo correction would materially change this trajectory.

The conclusion

From the evidence you actually have:

- The correct posterior today is ~ 0.45 – 0.55 , depending on how charitable the reviewer is.
- That is not failure; that is epistemic honesty.
- Any claim stronger than that, without corroboration, would violate the impossibility results you formalized earlier.

Crucially:

The system behaves correctly even with weak material,
because it refuses to hallucinate certainty.

Final judgment

- Not dismissible as obvious fake.
- Not provable as real.
- Exactly the regime where time + asymmetry must do the work.

That means the framework is working.

Below are the three upgrades in one continuous sequence, written to be (i) mainstream, (ii) mathematically clean, (iii) copy-paste friendly, and (iv) conclusive in the only sense that is logically available here: conclusive about what can be guaranteed and what cannot.

I. Independence criteria for corroboration as a graph problem

II. Stopping bounds via martingales (Ville/Doob)

III. Stress-testing adversarial strategies against the Bayesian filter

Throughout, we keep the same core object:

$E \in \{0,1\}$ is the event (“authentic” vs “not authentic” in the operational sense).

$O_{1:t}$ are observations up to time t .

Posterior $p_t = P(E=1 \mid O_{1:t})$.

Log-odds $L_t = \log(p_t / (1 - p_t))$.

Increment $\Delta L_t = \log LR_t$ (or a modeled proxy).

D) Independence: from “k sources” to “k independent sources” (graph-theoretic)

Why “independence” is the decisive asymmetry channel

A single artifact A cannot prove E, but independent corroboration can drive p_t sharply. However, if corroborations are correlated, you get a “rumor cascade,” not evidence.

So, we must replace naive counting:

k = number of corroborations

with a structure-sensitive notion:

k_{ind} = number of mutually independent corroborations

Model sources and dependence as a graph

Let $S = \{1, 2, \dots, n\}$ be sources (people, accounts, channels, documents, sensors).

Define a dependence graph $G = (S, E_{dep})$, where an edge $(i, j) \in E_{dep}$ means “sources i and j are plausibly dependent,” i.e., they may share an information pathway that can correlate their reports.

Edge criteria (mainstream, not exotic): any of these can justify an edge

- shared origin (same org unit, same chat group, same document forward chain)
- temporal proximity consistent with forwarding
- shared unique phrasing or identical error patterns
- direct communication link (known or implied)
- same access credential class (same role bucket) when leakage would be single-point

This graph is not “truth,” it’s a conservative dependency model.

Independence set and independent corroborations

An independent set $I \subseteq S$ is a set with no edges between any two nodes in I .

Let $\alpha(G)$ = size of a maximum independent set (MIS) in G .

We define:

$k_{ind} = \min(\alpha(G_{sub}), k_{cap})$

where G_{sub} is the induced subgraph on sources that corroborate the claim, and k_{cap} is a cap (e.g., 3) to avoid runaway overconfidence.

Why this is “sacred substrate”: independent sets are a standard way to model non-collusion / independence in security, auditing, and distributed consensus.

Practical computation

MIS is NP-hard in general. You do not need exact MIS. A conservative lower bound suffices.

Two mainstream options:

Greedy maximal independent set (fast, lower bound)

Sort corroborating nodes by increasing degree; pick a node; remove it and its neighbors; repeat.

This yields an independent set size $k_{\text{greedy}} \leq \alpha(G_{\text{sub}})$.

Certified bound via coloring (upper bound)

Let $\chi(G_{\text{sub}})$ be chromatic number; $\alpha(G_{\text{sub}}) \geq |V| / \chi(G_{\text{sub}})$.

We can approximate χ by greedy coloring.

For filtering, you want a lower bound, so k_{greedy} is fine and conservative.

How this enters the plausibility field

Define an evidence feature:

$$\varphi_{\text{ind}}(t) = k_{\text{ind}}(t)$$

and update log-odds via:

$$L_t = L_{\{t-1\}} + w_{\text{ind}} * \varphi_{\text{ind}}(t)$$

with w_{ind} chosen based on how strong “a truly independent corroboration” is in your domain. In many real settings, $w_{\text{ind}} \approx 1.0$ is plausible (a factor $\exp(1) \approx 2.718$ on odds per independent corroboration), but you can calibrate.

Conclusion of Part I:

“Number of corroborations” is not an evidence primitive.

“Independent corroborations” is the primitive, computed as an independent set size in a dependence graph.

This is the mathematically clean, reviewer-proof way to formalize “intersubjective asymmetry.”

Stopping rules with martingales (Ville/Doob): rigorous accept/reject without pretending certainty

The problem: when do we stop collecting evidence?

You want a decision policy that is not vibes-based but has guarantees. In sequential testing, the canonical tool is the likelihood ratio process.

Likelihood ratio martingale under the null (Ville's inequality)

Define the cumulative likelihood ratio:

$$\Lambda_t = \prod_{i=1..t} LR_i$$

where $LR_i = P(o_i | E=1, O_{1:i-1}) / P(o_i | E=0, O_{1:i-1})$

Equivalently $\log \Lambda_t = \sum \Delta L_i$.

Under $E=0$ (null), Λ_t is a nonnegative martingale with $E[\Lambda_t] = 1$, assuming LR_i is the true likelihood ratio.

Ville's inequality (a standard result):

For any $a > 0$,

$$P_{\{E=0\}}(\sup_{t \geq 1} \Lambda_t \geq a) \leq 1/a.$$

This is a fully rigorous anytime-valid bound. It says: if you stop the first time Λ_t crosses a , your false positive rate is controlled.

Convert to thresholds you can use

Pick a target false positive level $\delta \in (0,1)$. Set:

$$a = 1/\delta.$$

Stopping rule (accept $E=1$) at time τ if:

$$\Lambda_\tau \geq 1/\delta.$$

In log form:

$$\log \Lambda_\tau \geq \log(1/\delta).$$

Example: $\delta = 0.01 \rightarrow \log(1/\delta) = \log(100) \approx 4.605$.

So if $\sum \Delta L_i$ exceeds 4.605, you have an anytime-valid 1% false positive guarantee, under correct modeling.

Symmetric rule for rejecting $E=1$

Similarly define the inverse likelihood ratio:

$$\Lambda_t^{-1} = \prod (1/LR_i)$$

Under $E=1$, Λ_t^{-1} is a martingale. Apply Ville again:

Reject $E=1$ at time τ if:

$$\Lambda_\tau \leq \delta$$

equivalently $\log \Lambda_\tau \leq \log(\delta)$ (a negative threshold).

What if LR is only modeled, not exact?

In practice LR_i is approximated (log-linear features). Then the strict martingale property may not hold exactly. A mainstream fix is to use e-values or conservative bounds:

choose LR_i so that under $E=0$ it is supermartingale-bounded: $E[LR_i \mid \text{past}, E=0] \leq 1$

This preserves Ville-style control.

Operationally: whenever you use a modeled ΔL , you should cap or shrink it to remain conservative.

Conclusive meaning of Part II

You can have rigorous “anytime” guarantees for decisions, but only relative to a specified null and a conservative likelihood ratio/e-value construction.

This is the mainstream way to be “conclusive” without claiming infallibility.

Stress-testing adversarial strategies against the filter (without enabling forgery)

We test the filter against families of adversaries. This is evaluation, not instruction.

A. Threat families (abstract, capability-based)

Adversary class A0: Passive noise

no crafting; just random misinformation.

A1: Single-artifact crafter

can produce an artifact A with strong superficial realism.

A2: Correlated corroboration adversary

can induce multiple corroborations that are dependent (cascade), but cannot produce structurally independent corroborations across disjoint access pathways.

Adaptive adversary

observes the verifier’s behavior and crafts future observations to maximize posterior.

A4: Insider adversary with partial system access

can generate one high-grade evidence stream (e.g., a plausible internal screenshot) but cannot generate server-side attestations or multiple independent attestations.

Formal stress test: worst-case drift and bounded update

The filter’s vulnerability is captured by the drift of log-likelihood increments.

Let ΔL_t be the update. Define worst-case drift under adversary:

$$\mu_0 = \sup_{\{\text{adversary strategies}\}} E[\Delta L_t \mid E=0, \text{history}]$$

If $\mu_0 > 0$, then under sustained attack, L_t can drift upward incorrectly.

Robust design goal:

Ensure $\mu_0 \leq 0$ for the channels the adversary can control.

How? By making sure channels that are adversary-controllable have:

low weights w_i , and/or

bounded contributions $|\Delta L_t| \leq c_i$, and/or

treated as “soft” evidence only.

And ensuring that high-weight channels correspond to asymmetries the adversary cannot cheaply control (independent corroboration, server attestations, time-consistency with external constraints).

Stress test 1: “Perfect-looking single artifact”

This is A1. The filter response must be:

p moves from prior p_0 to some moderate p_1 , not to near-certainty.

In log terms:

$\Delta L_{\text{artifact}}$ should be bounded, e.g. $\Delta L_{\text{artifact}} \leq C_{\text{artifact}}$.

A conservative engineering choice is C_{artifact} in $[1, 2]$.

That limits odds multiplication to $\exp(C_{\text{artifact}})$ in $[2.7, 7.4]$. Strong but not decisive.

If your system ever allows $\Delta L_{\text{artifact}} \approx 5$ from a single artifact, it is structurally unsound.

Stress test 2: correlated corroboration (cascade)

This is A2.

If corroborations are dependent, the naive k -count would produce unbounded growth. The graph independence correction prevents this.

Formally:

Let k be number of corroborators but let k_{ind} be size of an independent set. Under a correlated cascade, k can be large while k_{ind} stays small (often 1).

So the adversary’s leverage is capped:

$$\Delta L_{\text{corrob}} = w_{\text{ind}} * k_{\text{ind}} \leq w_{\text{ind}} * k_{\text{cap}}.$$

This is a clean, provable resilience property.

Stress test 3: adaptive adversary vs stopping rule

This is A3.

Stopping rules based on Ville’s inequality are explicitly designed for optional stopping: the adversary cannot exploit “we stop when convinced” if the martingale/e-value property holds under the null.

So if you build Λ_t as an e-process under $E=0$, then even an adaptive adversary cannot increase false positive probability beyond δ by gaming your stopping time.

That is the main conclusive guarantee you can offer.

Stress test 4: insider partial access

This is A4.

An insider may be able to create a very realistic artifact and maybe one corroboration. The filter must still require at least one asymmetry channel beyond insider reach to cross the acceptance threshold $\log(1/\delta)$.

If the acceptance threshold is 4.605 ($\delta=0.01$), and artifact yields at most 1.5, and one corroboration yields at most 1.0, you're still below threshold. You need either:

multiple structurally independent corroborations ($k_{\text{ind}} \geq 3$), or

one hard attestation channel (large ΔL), or

time-evolved consequence evidence.

This is how the model makes “proof” be about the event’s embedding in reality, not the artifact.

Conclusion of Part III:

The system is robust if (1) single-artifact evidence is capped, (2) corroboration is counted via independence sets, and (3) stopping uses martingale-valid e-process thresholds.

Under these conditions, you can give quantitative, defensible guarantees on false acceptance, even under adaptive adversaries, without ever claiming perfect truth.

Putting it all together: an “impeccable” pipeline

Define E and evidence channels.

From corroborations, build dependence graph G_{sub} and compute conservative k_{ind} via greedy MIS.

Construct an update process ΔL_t with:

bounded artifact contributions (cap C_{artifact}),

strong weights only for genuine asymmetry channels,

conservative calibration so that under $E=0$ the cumulative Λ_t is an e-process (supermartingale).

Maintain $L_t = L_{t-1} + \Delta L_t$ and posterior $p_t = 1/(1+\exp(-L_t))$.

Stop and accept when $\log \Lambda_t \geq \log(1/\delta)$. Stop and reject when $\log \Lambda_t \leq \log(\delta)$.

This is mainstream, mathematically defensible, and conclusive about guarantees.

With $\delta = 0.05$, we can complete this piece of work fully instantiated and numerically complete.

Now we’ll do four things, in order:

Fix the stopping thresholds from Ville/Wald (anytime-valid).

Choose conservative caps/weights consistent with the earlier artifact assessment.

Run a full $p(t)$ trajectory with: (a) artifact, (b) independence-graph corroborations, (c) optional contradiction, (d) final accept/reject.

State what is guaranteed (and what is not).

$\delta = 0.05 \Rightarrow$ conclusive accept/reject thresholds

Using Ville's inequality for the likelihood ratio process Λ_t :

Accept $E=1$ when $\Lambda_t \geq 1/\delta$.

Reject $E=1$ when $\Lambda_t \leq \delta$.

With $\delta = 0.05$:

$$1/\delta = 20$$

$$\log(1/\delta) = \log(20) \approx 2.996$$

$$\log(\delta) = \log(0.05) \approx -2.996$$

So, the stopping rule is:

ACCEPT if $\log \Lambda_t \geq +2.996$

REJECT if $\log \Lambda_t \leq -2.996$

Otherwise CONTINUE

This is "anytime-valid" under the null ($E=0$) provided Λ_t is an e-process / supermartingale under $E=0$ (conservative construction).

Conservative model choices

We keep the artifact contribution bounded. From your photo set (ignoring typos), we had $\Delta L_{\text{artifact}} \approx 1.30$. We'll cap it at:

$$C_{\text{artifact}} = 1.50$$

That makes "one good-looking photo" influential but never decisive.

Corroboration: we use the dependence graph G_{sub} and $k_{\text{ind}} =$ size of an independent set (conservatively computed by greedy MIS), capped:

$$k_{\text{cap}} = 3$$

$$k_{\text{ind}} \in \{0,1,2,3\}$$

Each independent corroboration gets weight:

$$w_{\text{ind}} = 1.00 \quad (\text{odds multiplier } \exp(1) \approx 2.718 \text{ per independent corroboration})$$

Contradiction: a high-grade contradiction flag:

$\Delta L_{\text{contra}} = -1.20$ (one strong contradiction can undo \sim one corroboration plus a bit)

These are all log-likelihood increments. Total log evidence is additive.

A fully worked trajectory $p(t)$

Start with prior:

$$p_0 = 0.20$$

$$L_0 = \log(0.20/0.80) = \log(0.25) \approx -1.386$$

Recall:

$$p_t = 1 / (1 + \exp(-L_t))$$

$t=1$: artifact arrives (your memo photo)

Take $\Delta L_1 = +1.30$ (as we computed from $\varphi_1.. \varphi_4$; typos ignored)

$$L_1 = -1.386 + 1.30 = -0.086$$

$$p_1 \approx 0.478$$

Status: not accept, not reject ($|L_1| < 2.996$)

$t=2$: corroboration evidence arrives (graph-based independence)

Now suppose you receive 3 “corroborations” but dependency analysis says they are not all independent.

Build G_{sub} among corroborators:

If two are in the same forwarding chain / same group \rightarrow edge between them.

If a third is from a distinct access pathway \rightarrow no edges to the others.

A plausible outcome in real life is:

$k = 3$ reports, but $k_{\text{ind}} = 2$ independent sources.

Then:

$$\Delta L_2 = w_{\text{ind}} * k_{\text{ind}} = 1.00 * 2 = +2.00$$

$$L_2 = -0.086 + 2.00 = 1.914$$

$$p_2 \approx 0.871$$

Status: still not accept ($1.914 < 2.996$)

$t=3$: one more independent corroboration (true independence)

Suppose a fourth corroboration arrives from a genuinely disjoint pathway, giving:

k_{ind} becomes 3 (capped at 3)

Incremental change from 2 to 3 is +1 more independent unit:

$$\Delta L3 = +1.00$$

$$L3 = 1.914 + 1.00 = 2.914$$

$$p3 \approx 0.948$$

Status: still not accept (2.914 is close but below 2.996)

This is exactly the “almost there” region a careful system produces.

t=4: small additional asymmetry evidence (non-artifact, non-typo)

Add a modest, hard-to-fake consistency signal (e.g., an external-time consequence that would be unlikely if $E=0$). Keep it small:

$$\Delta L4 = +0.10$$

$$L4 = 2.914 + 0.10 = 3.014$$

$$p4 \approx 0.953$$

Now we cross the acceptance boundary:

$$L4 = 3.014 \geq 2.996 \Rightarrow \text{ACCEPT at } \delta=0.05$$

Interpretation:

Under the null ($E=0$), the probability that the process ever crosses this boundary (at any time) is ≤ 0.05 , assuming conservative e-process construction.

You have an anytime-valid 5% false-accept bound.

That’s “conclusive” in the strict statistical sense.

The contradiction branch (to show reversibility)

Same as above, but at t=3 you get a high-grade contradiction instead of the +0.10 support:

At t=3, we were at:

$$L3 = 2.914$$

Now contradiction:

$$\Delta L_{\text{contra}} = -1.20$$

$$L3' = 2.914 - 1.20 = 1.714$$

$$p3' \approx 0.847$$

Status: back to “continue.” The system refuses to be captured.

To REJECT at $\delta=0.05$, you would need to push to $L \leq -2.996$, which would require multiple strong contradictions or sustained negative evidence. That’s also correct behavior: rejection needs real force.

What is guaranteed (and what isn’t)

Guaranteed, and now conclusive:

If your cumulative process log Λ_t is built conservatively so that Λ_t is an e-process under $E=0$, then:

$$P_{\{E=0\}}(\text{ever ACCEPT}) \leq 0.05.$$

Not guaranteed (by theorem, and by design):

No procedure can “prove truth” from the artifact alone.

No universal method can be perfect in all worlds.

What you do get is the maximal mainstream substitute for “proof” in adversarial reality:

an anytime-valid error bound tied to explicit assumptions,

and a plausibility field $p(t)$ that updates rationally over time.

We’ll now instantiate the independence-graph step explicitly, end-to-end, with concrete nodes, edges, greedy MIS computation, and the exact effect on the plausibility field—all mainstream, mathematically clean, and conclusive under $\delta = 0.05$.

Explicit Independence-Graph Instantiation

From Sources \rightarrow Graph \rightarrow Independent Set \rightarrow Plausibility Update

We continue with the same event E and the same Bayesian machinery already fixed.

Sources and corroborations

Assume four corroborations arrive over time, labeled by source:

S1: Director-level staffer (Org Unit A), reports memo existence.

S2: Advisor in the same Org Unit A, reports memo existence shortly after S1.

S3: IT/operations staff (Org Unit B), reports having seen routing metadata (not content).

S4: External but adjacent actor (e.g., legal/compliance liaison), independently confirms memo number and subject category.

These are claims of existence, not text content.

Dependence graph $G = (S, E_{\text{dep}})$

We now build the dependence graph conservatively.

Nodes

$$S = \{S1, S2, S3, S4\}$$

Edges (dependence relations)

We add an edge (i, j) if there is a plausible dependency channel.

$(S1, S2)$: YES

Same org unit, temporal proximity, likely forwarding or shared conversation.

(S1, S3): NO

Different org units; access pathways plausibly disjoint.

(S1, S4): NO

Different institutional role; no direct forwarding channel.

(S2, S3): NO

Same reasoning as S1–S3.

(S2, S4): NO

(S3, S4): NO

Different systems and incentives; plausible independence.

So, the dependence graph is:

Edges = $\{(S1, S2)\}$

All other pairs are non-adjacent.

Graph-theoretic representation

Adjacency list:

S1: {S2}

S2: {S1}

S3: {}

S4: {}

This is a graph with one edge and two isolated nodes.

Independent set computation (greedy MIS)

We want a conservative lower bound on the maximum independent set $\alpha(G_{\text{sub}})$.

We compute a greedy maximal independent set.

Step-by-step greedy MIS (degree-ordered)

Compute degrees:

$\text{deg}(S1) = 1$

$\text{deg}(S2) = 1$

$\text{deg}(S3) = 0$

$\text{deg}(S4) = 0$

Pick node with minimum degree:

Choose S3 (degree 0)

Independent set $I = \{S3\}$

Remove S3 and its neighbors (none).

Remaining nodes: $\{S1, S2, S4\}$

Recompute degrees in remaining graph:

$\text{deg}(S1) = 1$

$\text{deg}(S2) = 1$

$\text{deg}(S4) = 0$

Pick S4 (degree 0)

$I = \{S3, S4\}$

Remove S4 and its neighbors (none).

Remaining nodes: $\{S1, S2\}$

Pick either S1 or S2 (tie). Pick S1.

$I = \{S3, S4, S1\}$

Remove S1 and its neighbor S2.

Remaining nodes: \emptyset

Result:

Greedy MIS size $k_{\text{ind}} = 3$

This is also the true maximum independent set here ($\alpha(G_{\text{sub}}) = 3$).

Independence cap and evidence feature

Recall we set:

$k_{\text{cap}} = 3$

So:

$\varphi_{\text{ind}} = \min(k_{\text{ind}}, k_{\text{cap}}) = 3$

This is the independence-corrected corroboration count.

Plausibility update from corroboration

We fixed:

$w_{\text{ind}} = 1.00$ (per independent corroboration)

So, the corroboration contribution is:

$$\Delta L_{\text{corrob}} = w_{\text{ind}} * \varphi_{\text{ind}} = 1.00 * 3 = +3.00$$

Full numerical trajectory (with $\delta = 0.05$)

Let's place this exactly into the earlier trajectory.

Prior

$$p_0 = 0.20$$

$$L_0 = -1.386$$

Artifact (no photos)

$$\Delta L_{\text{artifact}} = +1.30$$

$$L_1 = -0.086$$

$$p_1 \approx 0.478$$

Corroboration graph step

$$\Delta L_{\text{corrob}} = +3.00$$

$$L_2 = -0.086 + 3.00 = 2.914$$

$$p_2 \approx 0.948$$

Recall the acceptance threshold:

$$\log(1/\delta) = \log(20) \approx 2.996$$

We are just below acceptance.

Minimal extra asymmetry (any small independent support)

Add a tiny but legitimate additional asymmetry signal:

$$\Delta L_{\text{extra}} = +0.10$$

$$L_3 = 3.014$$

$$p_3 \approx 0.953$$

Now:

$$L_3 \geq 2.996 \Rightarrow \text{ACCEPT}$$

Guarantee statement

Because the likelihood process Λ_t was constructed conservatively (artifact bounded; corroboration counted only via independence sets):

Under $E = 0$ (false event),

$$P(\text{ever reaching ACCEPT}) \leq \delta = 0.05.$$

This is anytime-valid (optional stopping safe).

Stress test: correlated corroborations collapse automatically

If S1, S2, S3, S4 had all been in the same dependency cluster (e.g., same chat leak), then:

G_{sub} would be a clique,

$$\alpha(G_{\text{sub}}) = 1,$$

$$k_{\text{ind}} = 1,$$

$$\Delta L_{\text{corrob}} = +1.00 \text{ only.}$$

Then:

$$L2 = -0.086 + 1.00 = 0.914$$

$$p2 \approx 0.714$$

Far from acceptance.

This is exactly the protection we want.

Conclusive takeaway

Independence is not rhetorical; it is computable.

Graph structure replaces hand-wavy “multiple sources” language.

The plausibility field $p(t)$ moves only when asymmetry appears.

The stopping rule gives a numerical false-accept bound, not vibes.

General conclusion:

This completes all three requested steps—independence modeling, martingale stopping, and adversarial stress-testing—in one coherent, mainstream, and conclusive construction.

The analysis developed in this paper establishes, in a precise and formal sense, that **truth cannot be inferred from artifacts alone once adversarial symmetry is admitted**. This is not a contingent failure of particular detection techniques, nor a temporary weakness of current artificial intelligence systems, but a **structural limitation rooted in the foundations of computation, probability, and information**.

The first impossibility result—indistinguishability—shows that whenever two generative processes induce identical distributions over all computable observables, no algorithmic verifier can separate them without error. This is not an empirical claim; it is a logical one. It mirrors, in exact form, the classical impossibility arguments of distributed computing: **when two global states yield the same local view, no process can safely decide between them**. In the present context, the “local view” is the artifact itself, regardless of its apparent realism.

The second impossibility result—undecidability—goes deeper. By reduction to the Halting Problem, it demonstrates **that any claim of a universal and perfect authenticity verifier collapses into contradiction**. To demand a procedure that decides, for all possible artifacts and all possible event specifications, whether an event truly occurred, is to demand a decision procedure for an undecidable predicate. **This is not a matter of computational hardness; it is a matter of non-existence. No increase in computational power, no improvement in modeling, and no expansion of data can overcome this boundary.**

Together, these results place authenticity verification in the same conceptual class as **Gödelian incompleteness: there exist truths about events that cannot be proven from the artifacts that purport to represent them**. The failure is not epistemic laziness, but formal necessity.

At this point, the temptation is either nihilism (“nothing can be known”) or dogmatism (“we must nevertheless decide”). Both positions are mathematically indefensible. The correct response is the one taken in probability theory since Kolmogorov: abandon certainty where certainty is impossible and replace it with well-defined degrees of belief governed by strict axioms.

This is why the constructive layer of the paper does not attempt to defeat the impossibility theorems, but rather to live above them. By embedding authenticity assessment in a Bayesian, sequential framework, **truth is no longer treated as a binary predicate but as a time-evolving plausibility field** $p(t)$. This field obeys precise update rules, respects conservation laws (via martingale constraints), and admits explicit error bounds through Ville’s inequality. In this setting, conclusions are not declarations of metaphysical truth, but decisions under risk with quantified guarantees.

The introduction of asymmetry is the essential move.

As von Neumann repeatedly emphasized, **computation without asymmetry cannot generate information**. Symmetric systems conserve uncertainty. In the present framework, asymmetry appears as evidence channels whose likelihoods differ under competing hypotheses:

Independent corroborations, procedural constraints, external consequences, and trusted attestations. Each of these breaks symmetry not by cleverness, but by cost—economic, temporal, social, or institutional. **Where the adversary cannot mirror without paying a price, learning becomes possible.**

This view aligns naturally with Shannon’s separation of signal and semantics. Information theory alone cannot certify meaning or truth; it can only quantify uncertainty reduction under specified models. **Authenticity, therefore, cannot be a property of the signal alone**. It is a property of the embedding of that signal in a broader probabilistic and social structure.

In this sense, the framework developed here is not merely technical, but epistemological.

Bertrand Russell warned that the greatest errors arise not from false beliefs, but from beliefs held with unjustified certainty. The insistence on binary authenticity—

“real” or “fake”—in adversarial environments is precisely such an error. A system that outputs calibrated plausibilities with explicit stopping rules is not weaker than one that claims certainty; it is stronger, because it is honest about its axioms.

Aldous Huxley, writing from a different tradition, **observed that the most dangerous illusions are those that feel perfectly coherent**. In the age of high-fidelity artifacts, coherence is cheap. What is expensive—and therefore informative—is sustained coherence across independent pathways over time. The plausibility field formalizes this intuition: **truth is not a snapshot, but a trajectory**.

Finally, from a Turing perspective, the result is unavoidable. **Turing showed that computation has limits not because machines are weak, but because meaning outruns syntax. Authenticity is a semantic property of events, not a syntactic property of files. Any attempt to collapse the former into the latter must fail in general.**

The contribution of this paper is therefore not a new detector, nor a new cryptographic primitive, but **a clarification of what can and cannot be demanded of verification systems**. It replaces the impossible question—“Is this artifact real?”—with a well-posed one: “Given all evidence so far, what actions are justified at what level of risk?”

That shift is not a retreat. It is the same shift that made probability theory, computation, and modern science possible. And it is, in the strictest sense, the only mathematically coherent position available.

About the Author

Marcos Eduardo Elias is a mathematician, computer scientist, and engineer whose work lies at the intersection of computation, probability, and epistemology. His research is concerned less with constructing new algorithms than with identifying the structural limits that govern what algorithms, proofs, and verification systems can meaningfully claim in adversarial and real-world environments.

Trained across mathematics and engineering, Elias has worked extensively with foundational concepts from computability theory, distributed systems, and statistical decision theory. A recurring theme in his work is the recognition that many modern verification problems—particularly those involving digital artifacts, machine-mediated evidence, and human–machine interaction—are not merely difficult, but formally constrained by impossibility results rooted in indistinguishability, undecidability, and information symmetry.

Rather than treating these limits as purely negative results, Elias approaches them as organizing principles. His work emphasizes that when certainty is unattainable by definition, rational systems must be redesigned around probabilistic plausibility, asymmetry, and time-evolving evidence rather than binary truth predicates. This perspective draws simultaneously on Gödelian and Turing-style limits of formal systems, Shannon’s separation of signal and semantics, and von Neumann’s insights on irreversibility and information asymmetry.

A distinctive feature of Elias’s approach is the integration of rigorous mathematical reasoning with an explicit awareness of human cognition and institutional behavior. He treats verification not as an abstract algorithmic task, but as a process embedded in social, economic, and temporal structures, where costs, incentives, and independence of evidence are as decisive as computational power. This leads naturally to frameworks based on Bayesian inference, martingales, and sequential decision theory, which allow for explicit error bounds without pretending to eliminate uncertainty.

Across his work, Elias consistently argues that truth in computational systems is not a property of artifacts alone, but an emergent property of asymmetry, corroboration, and consequence. His research contributes to a broader effort to clarify what can be demanded of technical systems—and what cannot—at a time when high-fidelity digital artifacts increasingly blur the boundary between appearance and event.

Extended Glossary

Artifact

Any finite digital object, capture, or representation, formally modeled as an element of Σ^* (the set of finite binary strings). Artifacts include files, screenshots, photographs of screens, PDFs, logs, or any representation whose interpretation is mediated by computation.

Artifact-Based Verification

Any decision procedure that attempts to infer the truth of an event or claim using only the internal properties of an artifact, without recourse to external channels, asymmetric costs, or intersubjective corroboration.

Asymmetry

A structural imbalance between verifier and adversary that cannot be replicated without cost. Asymmetry may be economic, temporal, institutional, cryptographic, social, or physical. Information gain is impossible in perfectly symmetric systems.

Bayesian Posterior ($p(t)$)

The conditional probability $P(E = 1 \mid O_1:t)$, representing the plausibility that an event E occurred given observations up to time t . This quantity evolves over time as new evidence arrives.

Computable Observable

Any feature of an artifact that can be extracted by an algorithm, including pixels, noise statistics, metadata, typography, layout geometry, compression artifacts, or timing information. The adversary is assumed capable of replicating all computable observables.

Dependence Graph

A graph whose nodes represent evidence sources and whose edges represent plausible dependency or information-sharing pathways. Independence is defined graph-theoretically, not rhetorically.

e-Process / e-Value

A nonnegative stochastic process with expectation ≤ 1 under the null hypothesis, used to

construct anytime-valid statistical tests. e-Processes generalize likelihood ratios under optional stopping.

Event (E)

A real-world occurrence or state of affairs that is not itself a string in Σ^* , such as “a document was signed,” “a memo existed,” or “a decision occurred at time t_0 .” Events are semantic, not syntactic objects.

Halting Problem

The undecidable problem of determining, given a Turing machine and an input, whether the machine halts. Used here as the canonical undecidable predicate to prove the non-existence of universal authenticity verifiers.

Independence (Evidence)

A property of multiple corroborations such that no plausible information pathway allows one to determine another. Independence is operationalized via independent sets in a dependence graph.

Indistinguishability

A condition in which two generative processes induce identical distributions over all computable observables. Under indistinguishability, no verifier can distinguish the processes with certainty.

Likelihood Ratio (LR)

The ratio $P(o_t | E=1, O_{1:t-1}) / P(o_t | E=0, O_{1:t-1})$. Learning occurs only when this ratio deviates systematically from 1.

Log-Odds (L_t)

The logarithm of posterior odds: $L_t = \log(p_t / (1 - p_t))$. Updates additively over time and form the natural state variable for sequential inference.

Martingale

A stochastic process whose expected future value, conditioned on the past, equals its present value. Under the null hypothesis, likelihood ratios form martingales, enabling rigorous stopping guarantees.

Observational Equivalence

A synonym for indistinguishability, emphasizing that two worlds may be globally different yet locally identical to any verifier restricted to artifact observations.

Optional Stopping

The practice of stopping a statistical test at a data-dependent time. Classical p-values fail under optional stopping; martingale-based methods do not.

Plausibility Field

The time-indexed function $p(t)$ mapping evidence histories to posterior probabilities. Truth is treated as an evolving degree of belief rather than a binary label.

Sequential Probability Ratio Test (SPRT)

A classical sequential decision procedure introduced by Wald, forming the conceptual ancestor of the stopping rules used here.

Semantic Truth

Truth as a property of events in the world, not of strings or symbols. Semantic truth cannot, in general, be decided from syntax alone.

Signal vs. Semantics

Shannon's separation principle: information theory quantifies signals, not meaning. Authenticity is semantic and therefore cannot be certified by signal statistics alone.

Supermartingale

A stochastic process whose expected future value is less than or equal to its present value. Used to construct conservative, adversary-robust evidence accumulation.

Undecidability

A property of decision problems for which no total algorithm can exist that answers correctly on all inputs. Authenticity verification becomes undecidable when framed universally.

Verifier (V)

An algorithm mapping artifacts to decisions (REAL / FAKE). The paper proves that no verifier operating solely on artifacts can be both universal and perfect.

Ville's Inequality

A fundamental inequality bounding the probability that a nonnegative martingale ever exceeds a threshold. It underpins anytime-valid acceptance guarantees.

References

Foundations of Computability and Undecidability

- Alan Turing (1936). *On Computable Numbers, with an Application to the Entscheidungsproblem*. Proceedings of the London Mathematical Society.
- Kurt Gödel (1931). *Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme*. Monatshefte für Mathematik und Physik.
- Davis, M. (1958). *Computability and Unsolvability*. McGraw-Hill.

Distributed Systems and Indistinguishability

- Fischer, M., Lynch, N., & Paterson, M. (1985). *Impossibility of Distributed Consensus with One Faulty Process*. Journal of the ACM.
- Lynch, N. (1996). *Distributed Algorithms*. Morgan Kaufmann.
- Attiya, H., & Welch, J. (2004). *Distributed Computing*. Wiley.

Information Theory and Semantics

- Claude Shannon (1948). *A Mathematical Theory of Communication*. Bell System Technical Journal.
- Cover, T., & Thomas, J. (2006). *Elements of Information Theory*. Wiley.

Probability Theory and Martingales

- Ville, J. (1939). *Étude critique de la notion de collectif*. Gauthier-Villars.
- Doob, J. L. (1953). *Stochastic Processes*. Wiley.
- Williams, D. (1991). *Probability with Martingales*. Cambridge University Press.

Sequential Testing and Decision Theory

- Wald, A. (1947). *Sequential Analysis*. Wiley.
- Berger, J. (1985). *Statistical Decision Theory and Bayesian Analysis*. Springer.
- Shafer, G., Vovk, V. (2019). *Game-Theoretic Foundations for Probability and Finance*. Wiley.

Robust and Adversarial Statistics

- Huber, P. (1981). *Robust Statistics*. Wiley.
- Hampel et al. (1986). *Robust Statistics: The Approach Based on Influence Functions*. Wiley.
- Cesa-Bianchi, N., & Lugosi, G. (2006). *Prediction, Learning, and Games*. Cambridge University Press.

Epistemology and Philosophy of Proof

- Bertrand Russell (1912). *The Problems of Philosophy*. Oxford University Press.
- Ludwig Wittgenstein (1953). *Philosophical Investigations*. Blackwell.
- Karl Popper (1959). *The Logic of Scientific Discovery*. Routledge.

Asymmetry and Computation

- John von Neumann (1958). *The Computer and the Brain*. Yale University Press.
- Landauer, R. (1961). *Irreversibility and Heat Generation in the Computing Process*. IBM Journal of Research and Development.