

Scholarly Communication as an Engineering Discipline

Lei You^{1,2}, Lele Cao², Iryna Gurevych^{3,4}

¹Technical University of Denmark, ²CSPaper @ Scholar7, ³Technical University Darmstadt, ⁴MBZUAI

Scholarly communication has long relied on social filters: when writing and careful presentation were costly, polished prose and venue prestige served as imperfect signals of effort and quality. Large language models (LLMs) erode these signals by making fluent text cheap and enabling rapid, large-scale production of plausible manuscripts, including adversarial mimics designed for surface plausibility. Meanwhile, attention and verification remain scarce. We argue that scholarly communication has entered a communication-constrained regime. We formalize publishing as distributed inference with many encoders (studies) and many decoders (readers) operating under attention budgets and adversarial contamination. The model motivates a separation of representation from responsibility: machines may render prose, but authors remain accountable for structured claims, evidence, and uncertainty estimates. We derive limits showing why binary gatekeeping and implicit routing fail under overload and shifting base rates, and we translate them into requirements for the Protocol for Post-LLM Science: claim–evidence records, graded (probabilistic) certification, verification-aware routing, and explicit congestion control for scarce verification capacity.

Correspondence: lenart@scholar7.com

Date: February 2025



1 Introduction

Scientific publishing performs two jobs: it archives evidence so claims can be inspected and reused, and it coordinates attention so limited readers can decide what to trust. Journals, conferences, and peer review are one way of implementing these jobs. Here we emphasize the underlying information process. At its core, publishing is distributed inference: researchers generate evidence about an external state of the world, denoted by Θ , encode that evidence into public records, and a community of readers attempts to infer Θ while observing only a tiny fraction of what is produced under strict budgets for attention and verification.

Historically, the interface between authors and readers relied on a set of social heuristics rooted in economics. Because writing, careful presentation, and curation were costly, polished prose often correlated with effort. Affiliations and venue labels functioned as forms of costly signaling [1]. This social filtering, while imperfect, was workable when output volume was roughly commensurate with the community’s reading capacity. When readers could plausibly inspect the high-visibility slice of their field, coarse public labels were sufficient to route attention.

Communication Engineering Rather Than Sociology. Large language models (LLMs) shift scholarly communication to a new operating point by driving the marginal cost of fluent text toward zero. Empirical studies link LLM adoption to rapid growth in production and shifts in citation behaviour [2–4], while editors report that generative tools intensify adversarial contamination: plausibly written manuscripts intentionally produced to pass superficial filters and capture scarce attention (e.g., paper mills, LLM-assisted mass submissions) [5]. These changes break the core assumption of the old regime: *text fluency no longer reliably indicates investment*. When noise generation becomes cheap and attention remains scarce, the binding problems shift. Finding relevant work becomes a routing problem, distinguishing valid claims from synthetic hallucinations becomes an error-detection problem, and managing submissions against limited reviewer bandwidth becomes a congestion-control problem. Viewed in this light, the crisis in scientific publishing is *structural*: a communication-engineering problem, not only a sociological one.

Publishing as Distributed Inference. We therefore model post-LLM publishing as distributed inference under information constraints. We do not claim that scientific meaning reduces to bits. Rather, we encode semantic success as an inference objective: the community’s ability to reduce uncertainty about Θ .

In the post-LLM regime, the end-to-end bottleneck is no longer producing fluent text, but (i) allocating scarce human attention to a tiny subset of the record and (ii) defending that allocation against actors who can cheaply generate plausible-looking submissions. **Two communication constraints** therefore become binding:

- *Finite attention and verification capacity.* Namely, each reader can inspect only a small subset per time window, and deep checks such as peer review, reproducibility audits, or replication are a scarce service.
- *Strategic contamination.* Namely, a non-trivial fraction of records may be intentionally generated or manipulated to capture attention or pass superficial filters, e.g., LLM-assisted paper mills.

Once these constraints bind, familiar institutions such as binary acceptance and venue prestige are best analyzed not as gatekeepers of truth, but as low-bandwidth routing signals designed to allocate scarce attention. The structural failure of the current regime is that these signals lean heavily on text fluency. Since fluent text is no longer a reliable proxy for scientific validity, the system needs a new anchor.

Scope and Terminology. Our goal is to model the *mechanism-level* bottlenecks that emerge when throughput and adversarial contamination scale faster than human attention and verification. We use this model to diagnose structural failure modes of the current publishing system and to derive requirements for a new paradigm. We intentionally do *not* attempt a full socio-technical implementation blueprint (interfaces, governance, economics, and interoperability), because these constraints are rapidly evolving and difficult to forecast. Instead, we focus on what any viable implementation must achieve.

To keep language precise, we reserve the term *the Protocol* (capital P) for the scholarly communication system implied by these requirements, namely, a *Protocol for Post-LLM Science*. When we draw analogies to networking, we use terms such as *control scheme* rather than “protocol” to avoid overloading the word.

While scholarly communication is historically a human-centric endeavor, our model explicitly avoids assumptions that limit agents to human actors. Specifically:

- *Agent Neutrality:* We do not assume encoders (authors) or decoders (readers) are human. The roles of generating evidence and updating beliefs can be performed by AI agents or hybrid systems. The constraints of attention budgets and verification costs apply abstractly to any bounded computational entity.
- *Representation Neutrality:* We do not assume that the primary carrier of information must be natural language narrative. While prose is useful for human consumption, our model treats it as a rendering layer, with the core inference operating on structured records.
- *Consensus Neutrality:* We do not assume the goal is a single, centralized “version of truth.” The model supports distributed, asynchronous belief updates where different agents may hold different posteriors based on their local routing and verification paths.

Research novelty and time horizon. Assessing novelty is inherently contextual: reviewers must extract a paper’s core claims, retrieve and synthesizing related work, and making structured comparisons, which are all under tightening reviewer capacity [6]. Rather than adding novelty as a separate objective, we view it as a time-horizon effect in sequential inference. Shortsighted snapshot-risk minimization favors only immediately high-confidence records and can induce slowdown by filtering out uncertain but informative work. A scalable mechanism should therefore keep exploratory records visible but explicitly uncertain, and spend scarce verification on proof-bearing evidence that can turn uncertainty into reliable knowledge over time.

First Principle: Separating Prose from Proof. This diagnosis motivates a design shift:

Separation of Prose from Proof. Machines may render fluent text, but authors remain strictly accountable for the structured web of claims, evidence, and uncertainty.

The rationale is simple. LLMs can manufacture persuasive prose at low cost. What remains hard to fake—and still essential for trust—is a coherent, checkable chain from data and methods to claims. By shifting the burden of proof from *surface narrative* to *verifiable structure*, the Protocol restores a costly signal that correlates with scientific investment.

To operationalize this separation, we redefine the publication unit. Instead of treating a narrative manuscript as the primary object, the Protocol treats the core contribution as a *Claim and Evidence Record (CER)*. In the Protocol,

- **Prose** is a rendering layer: it may be generated (and personalized) by machines.
- **Proof** is the accountable layer: authors sign structured claims, evidence pointers, and uncertainty.

This architecture decouples readability from reliability, allowing text to scale while keeping verification grounded in the proof-bearing fields.

Our Contributions Are Threefold. First, we propose a concise mathematical model for post-LLM scholarly communication with many encoders and many decoders, attention constraints, and adversarial noise. Second, we derive limits that expose structural failure modes of current practice, including overload effects that make coarse public labels insufficient for routing. Third, we translate these limits into design requirements for the Protocol, namely structured messages, probabilistic certification, verification-aware routing, explicit congestion control for verification workflows, and incentive-compatible reporting.

Section 2 situates our approach relative to information theory, decentralized detection, rational inattention, and social learning. Section 3 specifies the mathematical model. Section 4 states our main theorems and interprets them as design constraints. We defer proofs to Appendix A.

2 Related Work

Our goal is to motivate *Scholarly Communication Engineering (SCE)* as a natural continuation of several existing lines of work, rather than as a new label introduced in isolation. Across disciplines, the literature has progressively shifted from treating publishing as a primarily *social* institution to treating it as a *large-scale information system* whose reliability depends on routing, verification, incentives, and robustness under strategic behavior.

Social Filtering, Incentives, and the Limits of Prestige. Classic accounts emphasize that scholarly communication has historically relied on social filters, i.e. reputation, venue prestige, and stylistic polish, to coordinate trust and attention. When producing careful prose and packaging results was costly, these filters could act as imperfect but useful signals of effort and quality, in the sense of costly signaling [1].

However, social learning theory also warns that when agents rely heavily on public signals, communities can herd and lock into persistent errors even when private evidence exists [7, 8]. At the same time, limited attention implies that actors will rationally economize on information processing [9], which makes low-bandwidth proxy signals (prestige, citation counts, “accept/reject”) especially tempting. A long-running concern is that once proxies become targets, they are gamed and degrade: this is often discussed as Goodhart’s law [10, 11] and as Campbell’s law in evaluation settings [12].

These incentive and selection pressures are frequently invoked in explanations of systematic reliability problems, including the claim that many published findings may be false under common research and publication practices [13] and evolutionary accounts of how competitive incentives can select for “bad science” [14]. In parallel, the growth of the scientific enterprise has raised concerns that producing the next unit of progress may require increasing effort, which amplifies the value of efficient coordination mechanisms [15].

Task-Oriented Information Theory and Distributed Detection. An engineering treatment becomes natural once we model publishing as a pipeline that moves evidence to readers under strict constraints. Shannon’s original theory formalizes reliable transmission of symbols while explicitly excluding semantics [16].

Subsequent work has shown that semantics can be reintroduced by measuring success through *task loss*—for example logarithmic loss, which connects expected loss to conditional entropy and yields clean multiterminal characterizations [17]. In a complementary tradition, decentralized detection and hypothesis testing study

how rate limits and quantization constrain error probabilities when many sensors report evidence to support decisions [18, 19].

These frameworks are directly relevant when publication venues and review outcomes are interpreted as low-rate public signals that must support downstream decisions. Finally, strategic or worst-case perturbations have long been modeled in information theory via adversarial channels, including arbitrarily varying channels [20], which provides a language for studying reliability when parts of the information stream are intentionally manipulated.

Peer Review as a Socio-Technical Pipeline and an Object of Computation. Peer review has been studied both as an institution and as an algorithmic process with measurable failure modes. A broad overview of challenges and computational approaches appears in [21]. At the text level, “revise and resubmit” has been modeled as a structured collaboration process with intertextual dependencies [22], and recent community perspectives argue that NLP can contribute tools for review assistance, transparency, and decision support [23].

On the practical side, several communities have experimented with workflow reforms that decouple review from single deadlines, encourage reuse of reviews, and increase transparency. Examples include ACL Rolling Review [24] and OpenReview-style processes described in conference author guides [25]. Other interventions explicitly target reviewer bandwidth and submission incentives, such as IJCAI–ECAI’s “primary paper” initiative and associated submission policies [26, 27].

LLMs and the Changing Economics of Scientific Production. A key reason the engineering framing has become timely is that generative AI has altered the cost structure of producing and evaluating scientific text. Controlled evidence shows large productivity gains from generative AI in knowledge work [28]. At ecosystem scale, multiple studies document rapid growth in LLM usage in scientific writing and shifts in writing style and citation behavior [2–4, 29, 30].

Beyond volume, AI tools may reshape what research is pursued: recent evidence suggests that AI can expand impact while narrowing topical focus [31], and perspective pieces discuss how LLMs may change experimental research workflows [32] as well as how AI is already used in highly formal domains such as mathematics and theoretical physics [33]. In parallel, editors have reported that generative tools intensify the challenge of paper mills and other synthetic misconduct [5], strengthening the case for models that explicitly account for strategic contamination.

LLMs in Peer Review: Assistance, Automation, and New Attack Surfaces. Alongside authoring, LLMs are now used (formally or informally) in review and editorial workflows, raising both opportunities and risks. Early studies explored LLMs as reviewing assistants [34], while later experiments positioned LLMs as checklist or critique assistants to help authors and reviewers structure feedback [6, 35].

Other lines of work emphasize that review quality depends on *reasoning* rather than surface fluency: structured critical assessment benchmarks and frameworks aim to evaluate whether automated reviewers can detect flawed arguments [36–38]. Related work also emphasizes the need for scalable, high-quality review under continued growth, and explores how AI might augment rather than replace human judgment [39, 40].

At the same time, the community has raised concerns about over-reliance on AI in peer review and the need for responsible policies [41–43]. Empirically, large-scale monitoring shows that AI-modified reviews already occur in peer review ecosystems [44]. New attacks further motivate an engineering/security lens: authors have attempted to game AI-based reviewing by embedding hidden instructions [45], and related medical literature reports “invisible text” prompt injection attacks against AI review models [46].

Finally, detection-based responses can create their own harms; for example, GPT detectors have been shown to be biased against non-native English writing [47], suggesting that reliable verification must go beyond surface-form detection.

From Prose to Verifiable Artifacts: Checklists, Alignment, and Verification-First AI. A common thread in recent work is a move from judging *prose* to validating *artifacts* and their alignment with claims. Quality assurance tasks include checking paper–code consistency and making the evidence chain auditable

[48]. In this direction, recent arguments propose “verification-first” AI as a necessary response to scaling pressures in peer review [49].

Taken together, these literatures motivate treating scholarly communication as an end-to-end engineered system: social signals and policies remain important, but the dominant bottlenecks increasingly arise from constrained attention, limited verification capacity, and adversarially optimized noise.

3 Mathematical Model

We introduce a minimal model tailored to the post-LLM regime. It highlights three features that become decisive under overload: many independent sources of evidence, many readers who each see only small subsets of the public record, and the possibility of strategic contamination.

Model at A Glance. The model composes five ingredients into one end-to-end objective. An external state Θ generates study-specific evidence $(Y_i)_{i=1}^m$, authors encode evidence into public records $(X_i)_{i=1}^m$, a routing layer selects which records each reader sees under an attention budget, a verification layer adds costly auxiliary checks, and adversaries may inject contaminated records. Readers decode their local information sets into posteriors q_j . The Protocol design question is then: *what record format, routing policy, and verification policy minimize communication log loss under worst-case contamination?*

A compact way to see the causal flow is:

$$\Theta \longrightarrow \{Y_i\}_{i=1}^m \xrightarrow{\text{encode}} \{X_i\}_{i=1}^m \xrightarrow[\text{verify}]{\text{route}} \mathcal{I}_j \xrightarrow{\text{decode}} q_j.$$

3.1 Studies as Distributed Encoders

We model a single scientific question through an unknown state $\Theta \in \mathcal{T}$. Θ can be a parameter (for example, an effect size) or a proposition (for example, whether a claim holds under a stated methodology).

A study acts as an *encoder* that observes evidence about Θ and emits a public record. For $i \in \{1, \dots, m\}$, study i observes a random variable Y_i whose distribution depends on Θ .

Assumption 1 (Noisy evidence). *Conditioned on Θ , each study i observes evidence Y_i drawn from a conditional distribution $P_{Y_i|\Theta}$. We allow $P_{Y_i|\Theta}$ to vary across studies to represent heterogeneous methodologies and measurement quality.*

Study i publishes a record

$$X_i = f_i(Y_i), \tag{1}$$

where f_i is an encoding map and X_i takes values in a message space \mathcal{X} . The message space \mathcal{X} may be free-form text or a structured object. In current systems, X_i is often an article-level narrative. In the Protocol, X_i is structured to support inference and verification.

Definition 1 (Claim and Evidence Record (CER)). *A Claim and Evidence Record is a tuple $X = (C, E, U, \sigma, v)$ in which the claim C targets a component of Θ and specifies its applicability conditions; the evidence bundle E points to verifiable artifacts (for example, data, code, proofs, or an executable environment); the uncertainty summary U reports uncertainty about C (for example, a confidence interval or posterior credible set); σ is an author signature; and v is a version identifier.*

Definition 1 operationalizes our separation of prose from proof. The record X is the canonical, signed message for inference and verification. Narrative can still exist—and it is often essential for explanation, pedagogy, and persuasion—but in the Protocol it is treated as a *derived view* rendered from X for a particular audience (potentially by a machine), without changing what authors are accountable for. For concreteness, Appendix B instantiates Definition 1 by presenting a CER for this paper.

This choice is not merely stylistic. If routing and verification must infer claims and evidence from prose, they inherit an open-ended natural-language understanding problem. By making the claim–evidence object explicit, we reduce the relevant message space from “any fluent text” to a structured, auditable set of fields on which routing, certification, and verification can operate directly.

From first principles, once prose is treated as a rendering of the structured record, it cannot improve the information-theoretic optimum for reliability. If narrative text T is generated from X (deterministically or stochastically), then the data processing inequality gives $I(\Theta; T) \leq I(\Theta; X)$. In equation 5, this is the point: prose can be optimized for accessibility and cost, while the “hard” optimization—robust inference under attention and verification constraints—focuses on the proof-bearing fields that are accountable and checkable.

3.2 Readers as Local Decoders

We model the scientific community as a collection of readers or teams indexed by $j \in \{1, \dots, N\}$. Each reader observes only a subset of the available records because attention and processing capacity are limited.

Assumption 2 (Attention constraint). *Reader j observes a random subset $S_j \subseteq \{1, \dots, m\}$ and receives the information set*

$$\mathcal{I}_j = \{X_i : i \in S_j\}. \tag{2}$$

We assume a budget $|S_j| \leq k_j$. In extensions, we replace this cardinality budget with an information-processing constraint of the form $I(Z_j; \mathcal{I}_j) \leq C_j$, where Z_j is the reader’s internal representation.

Given \mathcal{I}_j and a prior belief π_j over \mathcal{T} , reader j forms a posterior distribution

$$q_j(\theta) = P(\Theta = \theta \mid \mathcal{I}_j; \pi_j). \tag{3}$$

We treat the mapping from \mathcal{I}_j to q_j as the decoder. Under LLM-assisted reading, a reader may first compress \mathcal{I}_j into a summary or a task-specific view, but this step remains part of the decoder and does not create information about Θ .

To evaluate communication effectiveness in the reliability dimension, we use logarithmic loss, a canonical loss for probabilistic beliefs that ties expected loss to conditional entropy and yields clean characterizations in multi-terminal settings [17]. Each reader reports q_j , and the realized loss is $-\log q_j(\Theta)$. The community-level risk is

$$\text{Risk} = \frac{1}{N} \sum_{j=1}^N \mathbb{E}[-\log q_j(\Theta)]. \tag{4}$$

Under well-specified Bayesian updating, this is not an arbitrary choice: for any report q_j , $\mathbb{E}[-\log q_j(\Theta) \mid \mathcal{I}_j, \pi_j]$ equals a cross-entropy, which decomposes into $H(\Theta \mid \mathcal{I}_j, \pi_j)$ plus a KL term, and is minimized when q_j equals the Bayesian posterior.

We term this metric Risk following statistical decision theory. Under logarithmic loss, false certainty is penalized sharply: assigning high probability to a wrong value of Θ is far more costly than admitting uncertainty. In the post-LLM regime, this captures the epistemic danger of “hallucination-like” failure modes, where the community becomes confidently convinced of a falsehood.

Remark: Time Horizon and Novelty. Equation 4 is written for a fixed window (for example, a year of the archive). When we view scholarly communication as an ongoing process, the same log-loss objective extends to a long-horizon criterion obtained by accumulating (or averaging) risk over windows. This makes “novelty” endogenous: mechanisms that minimize snapshot risk myopically can stagnate by suppressing uncertain but informative claims whose proof layer could, after follow-up and verification, reduce future uncertainty about Θ . Our limits in Section 4 are stated per-window, but their implications are precisely about how mechanisms must behave under sustained throughput.

3.3 The Channel: Routing, Verification, and Noise

The subset S_j is not exogenous in practice. Search, recommendation, and venue signals shape what each reader sees. We model this effect as a routing policy that selects or ranks records for each reader under the budget in Assumption 2. A central goal of the Protocol is to define routing primitives that favor information gain about Θ rather than proxy metrics.

Peer review and replication provide additional information beyond the original records. We model verification as a costly action that reveals auxiliary evidence V_i about Θ and about the integrity of X_i .

Assumption 3 (Verification channel). *For each record X_i , a verifier may pay cost $c_i > 0$ to obtain an auxiliary observation V_i with conditional distribution $P_{V_i|\Theta, X_i}$. Verification actions are budgeted so that the total expected cost does not exceed a system budget B_{ver} .*

Verification subsumes peer review, reproducibility checks, and targeted replication. In the post-LLM regime, *verifier time is the scarce resource*: reviewer-hours and replication capacity do not scale with m . In the model, the verification policy is therefore an explicit *resource allocation* decision—which records receive which checks, in what order, and with what reuse of prior outcomes.

Finally, we model adversarial contamination to capture strategic content generation. We let an adversary control a subset of records and choose their distributions subject to minimal syntactic validity constraints.

Assumption 4 (Adversarial contamination). *A fraction $\epsilon \in [0, 1)$ of records are adversarial. For adversarial indices i , the record X_i is drawn from an arbitrary distribution Q_i that may depend on the publicly visible rules of the publishing mechanism and on past records. The remaining records are generated according to equation 1 under Assumption 1.*

Assumptions 3 and 4 let us capture two distinct roles of institutions. Verification mechanisms reduce the impact of adversarial or low-quality records by producing additional evidence. Routing mechanisms protect limited attention budgets by allocating exposure to records that are expected to reduce uncertainty.

3.4 The Problem: Robust Distributed Inference

The objects above define a *mechanism-design* task for scholarly communication. We write Π for a candidate *publishing mechanism*—a formal communication core that specifies three coupled choices: (i) a record format (what the archive stores and what authors remain accountable for), (ii) a routing policy (how limited attention is allocated, i.e., which subsets $(S_j)_{j=1}^N$ each reader observes under Assumption 2), and (iii) a verification policy (how scarce verifier effort is allocated under Assumption 3 and how outcomes enter readers’ information sets).

The design question in this paper is to choose Π within an admissible class. We call the particular mechanism implied by the design principles in this Perspective *the Protocol for Post-LLM Science* (or simply, *the Protocol*). Throughout, *publishing mechanism* refers to the generic object, while *the Protocol* refers to our proposed instance.

Under adversarial contamination (Assumption 4), we evaluate a mechanism by its worst-case community risk. We can summarize the design problem as a robust distributed inference optimization:

$$\min_{\Pi \in \mathcal{P}} \max_{Q \in \mathcal{Q}_\epsilon} \text{Risk}(\Pi, Q), \tag{5}$$

where \mathcal{P} denotes the class of admissible publishing mechanisms that respect the attention and verification constraints in Assumptions 2 and 3. The set \mathcal{Q}_ϵ denotes contamination strategies that can control at most an ϵ fraction of records, as in Assumption 4. The quantity $\text{Risk}(\Pi, Q)$ is the logarithmic loss in equation 4 under the joint distribution induced by mechanism Π and adversary Q .

This formulation clarifies what makes the post-LLM regime different. Throughput m increases, the contamination fraction ϵ can increase, while the budgets in Assumptions 2 and 3 do not. As a result, the feasible region of equation 5 can collapse under legacy workflow primitives.

Separating prose from proof is a minimal Protocol move that expands the feasible region. Prose becomes a low-cost rendering layer, while accountability and verification attach to structured, proof-bearing fields. Operationally, this shifts routing and certification from text-derived cues to checkable evidence, where feature-matching attacks are less severe. Section 4 develops limits that diagnose this collapse and translate them into design requirements.

4 Theory: Current Limits and Future Design

Section 3 frames publishing as a robust distributed inference problem. A *publishing mechanism* chooses a record format, routes limited attention, and allocates scarce verification effort. An adversary can contaminate

part of the public stream. A mechanism succeeds when it minimizes the community log-loss risk in equation 4, namely when many readers reduce their uncertainty about the external state Θ . Equation 5 summarizes this objective as a min-max problem under resource constraints.

We now state limits that hold for broad classes of publishing mechanisms. Each limit rules out an intuitive patch and translates into a design requirement for the Protocol. Throughout, our central move is the principle of *separating prose from proof*: machines may render fluent narrative, while authors remain accountable for the structured web of claims and evidence in Definition 1. In terms of equation 5, this move (i) removes unstructured prose from the variables that routing and verification must optimize over, and (ii) shifts reliability signals from easily spoofed surface form to verifiable structure, which constrains the adversary’s options in the inner maximization.

To keep the Perspective readable, we use lightweight statements (often conservation laws) and interpret each one explicitly in terms of how it reshapes the objective and constraints in equation 5.

4.1 Gatekeeping Breaks: Coverage and Bandwidth

Many legacy workflows expose a low-rate public control signal. A venue may internally evaluate a submission in rich detail but publicly releases only a discrete label, often well approximated by accept/reject. Readers then use this label to route attention and to approximate reliability.

This subsection is a direct application of Assumption 2: each reader can process at most k_j records per window. No additional probabilistic assumptions are needed; the first limit is purely combinatorial. Fix a time window, such as a year. Within this window, the public archive exposes m records. Reader j can process at most k_j records in this window. We write

$$K := \sum_{j=1}^N k_j$$

for the community’s total decoding budget. Suppose the publishing mechanism publishes a binary label $A_i \in \{0, 1\}$ for each record i . We interpret $A_i = 1$ as membership in a high-priority channel. Let

$$m_{\text{hi}} := \sum_{i=1}^m \mathbf{1}\{A_i = 1\}$$

be the number of high-priority records.

Definition 2 (Coverage and load ratio). *For each record i , let $R_i \in \{0, 1\}$ indicate whether at least one reader processes record i within the window. If $m_{\text{hi}} > 0$, we define the high-priority coverage as*

$$\text{Cov} := \frac{1}{m_{\text{hi}}} \sum_{i:A_i=1} R_i. \tag{6}$$

We also define the high-priority load ratio as

$$\rho := \frac{m_{\text{hi}}}{K}. \tag{7}$$

Theorem 1 (Coverage limit for binary gatekeeping). *Fix a time window with m records, binary labels A_1, \dots, A_m , and total decoding budget K . For any routing policy and any realization with $m_{\text{hi}} > 0$, the high-priority coverage satisfies*

$$\text{Cov} \leq \min \left\{ 1, \frac{K}{m_{\text{hi}}} \right\} = \min \left\{ 1, \frac{1}{\rho} \right\}. \tag{8}$$

Theorem 1 is a conservation law. It says that coverage is the inverse of load. If a mechanism collapses many records into one high-priority bucket, then coverage must fall once m_{hi} exceeds the decoding budget K . In terms of equation 5, a binary high-priority channel can become self-defeating. The mechanism may succeed at *selecting* items, but it fails at *delivering* them to any decoder.

Binary labels also have an information ceiling.

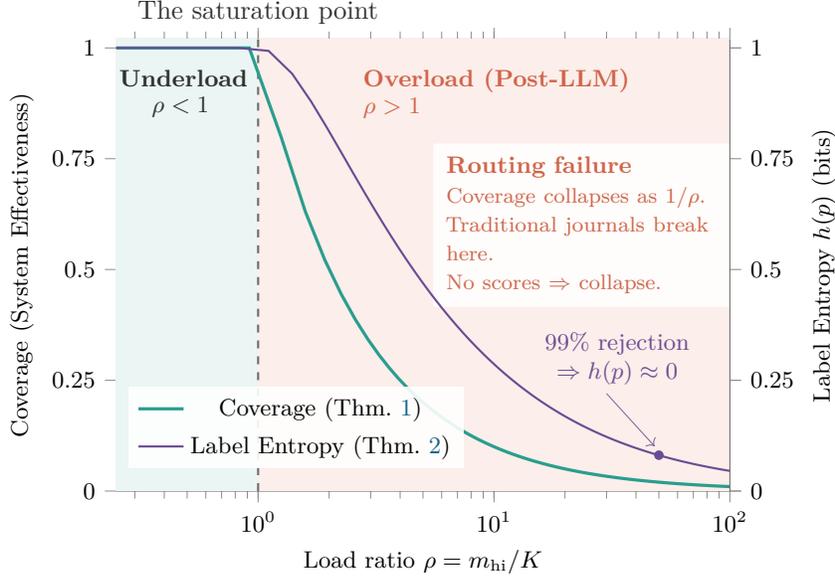


Figure 1: **The Attention–Coverage Cliff.** With load ratio $\rho = m_{\text{hi}}/K$, Theorem 1 yields $\text{Cov}(\rho) \leq \min\{1, 1/\rho\}$: once $\rho > 1$ (overload), coverage must fall arithmetically. The right axis shows the binary-label entropy $h(p)$ under an *illustrative* response $p(\rho) = 0.5 \min\{1, 1/\rho\}$ (so acceptance becomes rarer as overload grows), which drives $h(p) \rightarrow 0$ (Theorem 2) and motivates soft certification (scores).

Theorem 2 (Selectivity limits the information in acceptance decisions). *Suppose $A \in \{0, 1\}$ and let $p = \mathbb{P}(A = 1)$. Then*

$$I(\Theta; A) \leq H(A) = h(p), \quad (9)$$

where $h(p) = -p \log p - (1 - p) \log(1 - p)$ is the binary entropy function.

Theorem 2 says that a rare binary event carries little public information on average. When we reduce acceptance rates to preserve prestige, we also reduce the entropy of the public signal. Under overload, the system needs public metadata that can support routing and progressive belief updates. A one-bit signal cannot provide that bandwidth.

The next corollary combines the routing limit in Theorem 1 with the information ceiling in Theorem 2.

Figure 1 visualizes the coupled failure mode: as load grows, either coverage collapses (too many “high priority” items) or the public signal becomes information-poor (acceptance becomes too rare).

Corollary 1 (A one-bit scaling dilemma). *Consider a sequence of time windows indexed by m . In the m -record window, let $p_m := m_{\text{hi}}/m$ be the fraction of records labeled high priority. Assume the total decoding budget K is bounded above by a constant as $m \rightarrow \infty$. If p_m stays bounded below by a positive constant, then equation 8 implies that high-priority coverage converges to zero. If, instead, a policy keeps coverage bounded below by a positive constant for all large m , then $p_m \rightarrow 0$. Moreover, if we draw a record uniformly at random and write A for its label, then $\mathbb{P}(A = 1) = p_m$ and equation 9 implies $I(\Theta; A) \rightarrow 0$.*

Corollary 1 captures a regime shift. Under high throughput, a binary label cannot simultaneously act as a routable priority channel and an informative public signal. We can keep acceptance frequent enough to preserve information, but then the high-priority bucket becomes too thick to cover. Or we can shrink the bucket to preserve coverage, but then the public signal becomes information-poor.

Why Soft Certification Helps. Soft certification replaces the one-bit label A_i with a higher-rate public signal, such as a score S_i or a calibrated probability $s_i \in [0, 1]$. This helps with both failure modes in Corollary 1 and connects directly to equation 5.

Coverage (capacity matching). Theorem 1 is about arithmetic: once $m_{\text{hi}} \gg K$, no routing policy can “deliver” the whole high-priority slice to any decoders. A score lets the Protocol choose the slice size to match capacity. For example, the routing layer can promote only the top- K records by score (or set a threshold so that $m_{\text{hi}} \approx K$), which makes the upper bound in equation 8 close to one and turns scarce attention into actual coverage.

Information (public bandwidth). The information ceiling in Theorem 2 generalizes: for any public score S , we have $I(\Theta; S) \leq H(S)$. If S takes L discrete levels, then $H(S) \leq \log L$. Unlike a binary label ($L = 2$), a multi-level score can therefore carry more than one bit of public information, enabling progressive belief updates and finer routing decisions without requiring every reader to inspect every record.

Robustness (what the score depends on). In the post-LLM regime, scores help only if they depend on features an adversary cannot cheaply match. This is where separating prose from proof matters. If the public record is only narrative text, then a score will inevitably depend on textual surface cues and becomes vulnerable to feature matching (Section 4.4). If the public record is a structured claim–evidence object, then scores can depend on verifiable fields such as artifact availability, provenance, and replication status. This both improves routing and tightens the inner maximization in equation 5.

Finally, if the Protocol asks participants to publish probabilistic certificates, it must also make honest reporting incentive-compatible; Section 4.5 gives a concrete primitive via proper scoring rules.

A contemporary example. ICLR’s OpenReview-based workflow makes part of the internal certification record public, namely reviews and discussion [25]. In our model, this increases the rate of the public control channel. It does not remove overload, but it moves the system away from the one-bit dilemma by releasing a richer object than a venue label alone.

4.2 The Base-Rate Trap

The post-LLM regime changes not only volume but also composition. If fluent narratives become cheap to generate, then the base rate of reliability in the public stream can decrease. A common intuition is that stable reviewing standards should preserve reliability. Bayes’ rule says otherwise.

Assumption 5 (Binary reliability). *For the results in this subsection, we set $\Theta \in \{0, 1\}$, where $\Theta = 1$ denotes that a claim is correct and reproducible under its stated conditions. We write $\pi = \mathbb{P}(\Theta = 1)$ for the base rate of reliability.*

Theorem 3 (Posterior reliability after acceptance). *Under Assumption 5, define the true acceptance rate $\alpha = \mathbb{P}(A = 1 \mid \Theta = 1)$ and the false acceptance rate $\beta = \mathbb{P}(A = 1 \mid \Theta = 0)$. Then*

$$\mathbb{P}(\Theta = 1 \mid A = 1) = \frac{\alpha\pi}{\alpha\pi + \beta(1 - \pi)}. \quad (10)$$

Moreover, for any target $q \in (0, 1)$, the condition $\mathbb{P}(\Theta = 1 \mid A = 1) \geq q$ holds if and only if

$$\beta \leq \frac{\alpha\pi(1 - q)}{q(1 - \pi)}. \quad (11)$$

Figure 2 plots equation 10 for fixed sensitivity α and two false-accept rates β , illustrating how even “stable standards” can fail as the base rate π deteriorates.

Theorem 3 makes the base-rate trap explicit. As the base rate π falls, maintaining a fixed posterior reliability after acceptance forces the false-accept rate β to shrink proportionally. In equation 5, this corresponds to a more demanding inner maximization: the mechanism must either drive down false acceptance or acquire stronger evidence per claim.

The same point can be stated in likelihood language.

Theorem 4 (Evidence threshold grows as base rates fall). *Assume $\Theta \in \{0, 1\}$ with $\mathbb{P}(\Theta = 1) = \pi$. Let E be any evidence variable with conditional distributions $P(E \mid \Theta = 1)$ and $P(E \mid \Theta = 0)$. Define the likelihood*

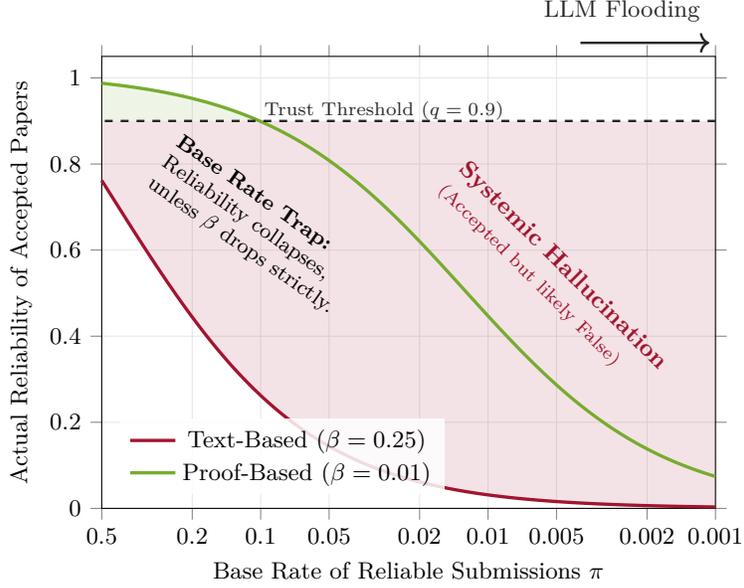


Figure 2: **The Base Rate Trap:** By Theorem 3, $P(\Theta = 1 | A = 1) = \frac{\alpha\pi}{\alpha\pi + \beta(1-\pi)}$ with $\alpha = 0.8$ collapses as the base rate π falls (log scale, reversed). Text-based review ($\beta = 0.25$, red) remains below a trust threshold $q = 0.9$ (shaded red). Proof-based verification ($\beta = 0.01$, green) stays above the threshold only when the base rate is sufficiently high (shaded green), illustrating how low base rates force β to shrink in proportion (Equation (11)).

ratio $\text{LR}(E) = P(E | \Theta = 1)/P(E | \Theta = 0)$. For any target $q \in (0, 1)$, the condition $\mathbb{P}(\Theta = 1 | E) \geq q$ implies

$$\log \text{LR}(E) \geq \log \frac{q}{1-q} - \log \frac{\pi}{1-\pi}. \quad (12)$$

In particular, when π is small, the right-hand side grows on the order of $\log(1/\pi)$.

Theorem 4 clarifies why prose cannot anchor certification once base rates fall. When π is small, reaching a fixed posterior threshold requires evidence with large likelihood ratios. LLMs can generate prose whose surface statistics mimic reliable prose, so text-level cues are unlikely to yield large log LR. By contrast, executable artifacts, provenance, and independent replication can, in principle, produce much larger likelihood ratios because they are harder to spoof at scale.

This is why the separation of prose from proof is not cosmetic. It is an evidence-amplification strategy: it moves certification onto evidence types that can meet the likelihood threshold imposed by a low base rate.

4.3 Review Is a Queue

Peer review and replication are expensive. The verification pipeline must therefore satisfy a stability condition. If verification demand grows faster than verification capacity, then turnaround times and verification quality must degrade.

In this framework, reviewers and replicators are not a background institution; they are the scarcest resource. Their time (measured here in reviewer-hours) is limited and must be scheduled. In terms of equation 5, the verification component of the mechanism Π decides *which* records receive *which* checks (human or machine), and in what order, to maximize the marginal reduction in community risk per unit of verification effort.

Assumption 6 (Verification workload). *Verification tasks arrive at rate λ per unit time. Each task requires random effort with mean $\tau > 0$, measured in a common unit such as reviewer-hours. The system can supply verification effort at rate μ per unit time.*

Theorem 5 (Workload stability requires admission control). *Under Assumption 6, if $\lambda\tau > \mu$, then no mechanism can keep the expected backlog of unverified work bounded over time. In particular, the expected unfinished verification workload must grow at least linearly with time.*

Theorem 5 is another conservation law. If the stream demands more verification work than the community can supply, then delay is not an optimization problem. It is an arithmetic impossibility. In terms of equation 5, the feasible set collapses because no mechanism can satisfy the verification constraint with bounded backlog.

A post-LLM publishing system therefore needs *explicit congestion control*. A scalable system can act on all three terms in the stability inequality. It can reduce the arrival rate λ of tasks that demand deep human verification. It can reduce the mean effort τ per task by automating low-level checks. It can increase the effective service rate μ by reusing verification outcomes instead of repeatedly restarting them.

We remark that congestion control implies selective dropping but not permanent loss. Stability requires that some work be *deferred*: not every record can enter deep human verification immediately. In the Protocol, “dropping” is therefore interpreted as *not admitting* a record to the high-assurance verification queue *yet*, not erasing it from the public archive. A low-cost preprint layer functions as a buffer or cache. Records remain visible and citable, and as authors iterate (or as independent replications accumulate) a record can be re-admitted to the verification queue with improved evidence. This preserves the archiving mandate while making the verification pipeline stable.

The principle of separating prose from proof helps on each front. It reduces τ because verifiers do not spend effort on narrative quality. They spend effort on the proof layer, and machines can pre-screen that layer for obvious failures. It also allows the Protocol to keep archiving cheap while throttling the high-assurance verification channel, which reduces the effective λ that must enter the queue.

A Contemporary Example. The IJCAI-ECAI 2026 Primary Paper Initiative introduces a submission fee of USD 100 for each submission, which is waived only for *primary papers* for which none of the authors appear on any other submission [26]. In our model this is admission control by pricing. It aims to reduce λ by discouraging high-volume submission behaviour and by internalizing part of the review cost. It does not, by itself, solve the evidence problem in Theorems 3 and 4, and it can raise equity concerns if fees bind for some authors. These tradeoffs make the design question explicit rather than eliminating it.

A Second Example. ACL Rolling Review (ARR) decouples reviewing from any single conference deadline and lets reviews follow a paper across cycles [24]. This is a form of traffic shaping. Conference deadlines create bursty arrivals that stress verification capacity. ARR smooths arrivals over time, and it amortizes verification by reusing reviews rather than restarting them for each resubmission.

Communication systems offer a clean prototype for both aspects. A rolling review pipeline resembles admission control with a token bucket. Authors can generate submissions at any time, but the system admits work for deep review at a controlled rate using a buffer that smooths bursts. ARR also resembles feedback schemes such as automatic repeat request (ARQ)¹ in networking: transmit, check, and retransmit based on feedback. At the application layer the loop is analogous: submit, receive structured feedback, revise, and resend within a persistent session. Compared with stateless conference cycles, this reuse reduces wasted retransmissions and increases the effective verification capacity.

4.4 The Indistinguishability Frontier

The post-LLM regime makes some noise sources strategic. We therefore need a robustness bound. Consider any filter that observes a feature $F = \phi(X)$ of a record and tries to decide whether it comes from a reliable process or from an adversary. If an adversary can match the feature distribution of reliable records, then no classifier can perform well.

Theorem 6 (Indistinguishability under feature matching). *Let F be a feature with distributions P_1 under a reliable source and P_0 under an adversarial source. With equal priors on the two sources, the minimum*

¹What a coincidence of the name!

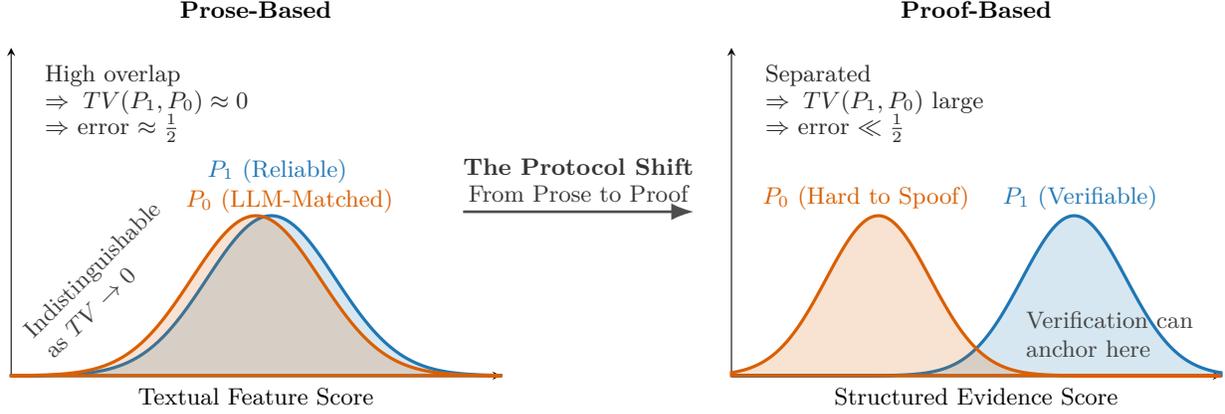


Figure 3: **The Feature-Matching Collapse (schematic)**. If screening relies on a feature space that an adversary can match, the resulting feature distributions P_1 (reliable) and P_0 (adversarial) can nearly overlap, making $TV(P_1, P_0) \approx 0$ and driving the optimal classification error toward $1/2$ (Theorem 6). Anchoring certification to structured, hard-to-spoof evidence is meant to increase separability, which relaxes the acceptance–rejection tradeoff in Corollary 2.

achievable classification error satisfies

$$P_e^* = \frac{1}{2} (1 - TV(P_1, P_0)), \quad (13)$$

where $TV(P_1, P_0)$ is the total variation distance. In particular, if $TV(P_1, P_0) = 0$, then $P_e^* = 1/2$.

Figure 3 gives an intuition for Theorem 6: when an adversary matches the distribution of the features being screened, classification error is driven toward $1/2$.

Theorem 6 explains why text-only screening becomes fragile. LLMs are trained to match the distribution of human-written text. In a purely textual feature space, $TV(P_1, P_0)$ can become small, which drives the optimal error toward random guessing. In terms of equation 5, this makes the worst-case adversary nearly indistinguishable from the reliable stream, so the max term dominates.

This is exactly where separating prose from proof matters. It moves screening from a feature space that an LLM can match to a feature space that requires consistent artifacts. In the proof layer, an adversary must supply a coherent evidence bundle, an executable environment, and traceable provenance. These constraints can increase $TV(P_1, P_0)$. Equivalently, in equation 5 they shrink the adversary’s effective action set by requiring checkable artifacts rather than merely fluent strings. They do not guarantee security, but they change what the adversary must imitate.

Figure 4 complements this by showing, in a stylized way, how low base rates and low separability jointly increase the verification burden.

Indistinguishability also implies an innovation tradeoff.

Corollary 2 (An innovation tradeoff under indistinguishability). *Let g be any screening rule that maps the feature F to an acceptance decision $A = g(F) \in \{0, 1\}$. Write $\Theta = 1$ for a reliable source and $\Theta = 0$ for an adversarial source. Define the true acceptance rate $\alpha = \mathbb{P}(A = 1 \mid \Theta = 1)$ and the false acceptance rate $\beta = \mathbb{P}(A = 1 \mid \Theta = 0)$. Then*

$$\alpha - \beta \leq TV(P_1, P_0). \quad (14)$$

Equivalently, for any target $\beta_0 \in [0, 1]$, if $\beta \leq \beta_0$ then $\alpha \leq \beta_0 + TV(P_1, P_0)$.

Corollary 2 formalizes a hidden cost of tightening standards. Low base rates force gatekeeping systems to reduce β to maintain posterior reliability, by Theorem 3. If $TV(P_1, P_0)$ is small, then reducing β also reduces

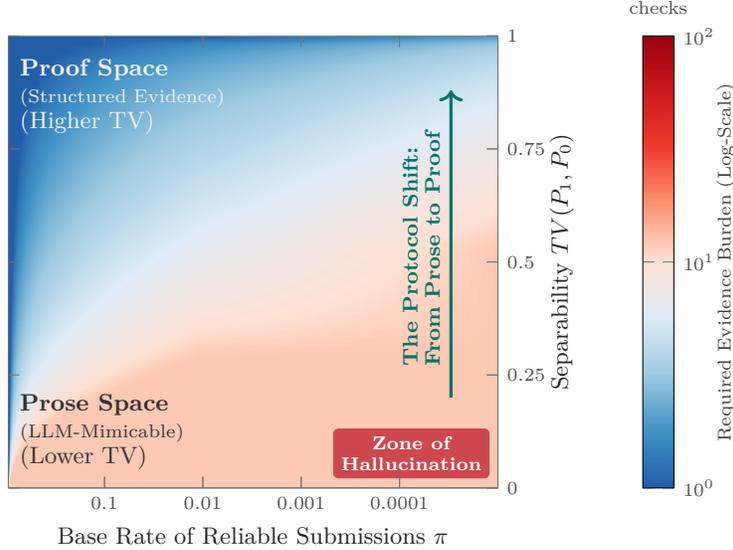


Figure 4: **The Indistinguishability Frontier** (stylized). Base rates π penalize single-claim inference (Theorem 4), while separability $TV(P_1, P_0)$ limits how well screening can distinguish reliable from adversarial streams (Theorem 6). Colors show $\log_{10}(n)$ (clipped to $[0, 2]$), where n is a toy estimate of how many independent “checks” are needed to reach posterior confidence $> 95\%$: $n \approx \ln(19(1 - \pi)/\pi) / \ln((1 + TV)/(1 - TV))$.

α . In practice, early-stage and unconventional results often produce weaker evidence, so they are among the first casualties of aggressive filtering. A system that separates archiving from certification can mitigate this effect by keeping uncertain claims visible while they accumulate stronger evidence. A system that separates prose from proof can mitigate it further by increasing $TV(P_1, P_0)$ in the evidence space, which relaxes the tradeoff.

4.5 Making Soft Certification Honest

The preceding results imply that scalable systems need higher-rate public signals than binary acceptance. A natural object is a probabilistic certificate, namely a public estimate of $\mathbb{P}(\Theta = 1 \mid \text{current evidence})$. Once the Protocol asks participants to report probabilities, it must also align incentives for honest reporting.

Theorem 7 (Log scoring makes truthful probability reports optimal). *Let $\Theta \in \{0, 1\}$ and suppose an evaluator believes $\mathbb{P}(\Theta = 1) = p$. The evaluator reports a probability $q \in (0, 1)$ and then receives score*

$$S(q, \Theta) = \Theta \log q + (1 - \Theta) \log(1 - q).$$

Then the expected score is uniquely maximized at $q = p$.

Theorem 7 provides a concrete primitive for soft certification. If the Protocol stores probabilistic assessments and later evaluates them against verification outcomes, then a proper scoring rule can reward calibrated judgments. This turns certification from a narrative opinion into a quantitative public signal that can be aggregated, audited, and used for routing.

Figure 5 illustrates the incentive alignment from Theorem 7: under a proper (log) scoring rule, truthful probability reports uniquely maximize expected reward.

Taken together, these results translate directly into requirements for the robust design problem equation 5. They describe *how* the mechanism class \mathcal{P} must be enriched (and how the adversary class \mathcal{Q}_ϵ must be constrained) for the optimization to have any hope of producing low community risk at scale:

- *Increase the public signal rate beyond one bit.* Binary acceptance is a low-bandwidth control channel. Graded certificates (scores or calibrated probabilities) increase public information rate and let routing

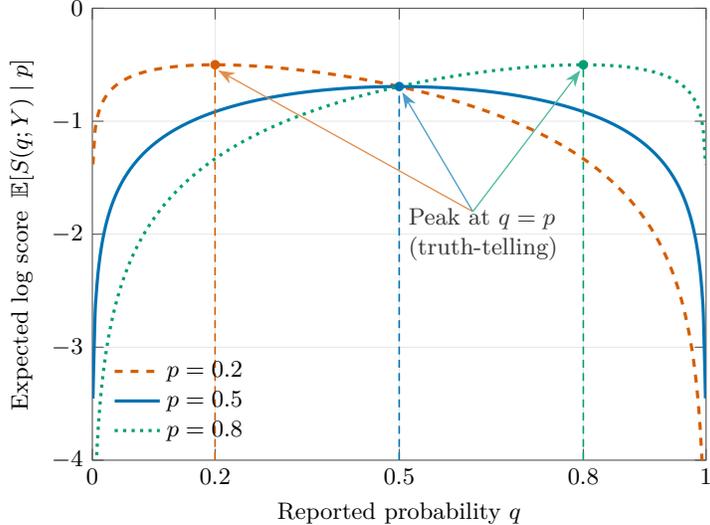


Figure 5: **The Reward of Honesty.** Under the log scoring rule in Theorem 7, a reviewer with true belief p maximizes expected score by reporting $q = p$. Any deviation (overconfidence or hedging) strictly lowers the expected score.

policies choose a high-priority slice whose size matches attention capacity K , improving both coverage and inference.

- *Treat verification as an explicitly optimized scarce service.* When verification is a queue, stability is a feasibility condition rather than a preference. Admission control, buffering (via preprints), and reuse of prior checks are required so that verification constraints remain satisfiable as m grows.
- *Move routing and certification onto verifiable structure.* Indistinguishability bounds show that text-only features can make reliable and adversarial streams statistically inseparable. By operating on structured claim–evidence records, the Protocol increases separability in the relevant feature space and effectively tightens the inner maximization in equation 5.
- *Make soft certificates incentive-compatible.* If the Protocol stores probabilistic assessments, proper scoring rules make truthful reporting optimal, so that the high-rate public signal can be aggregated and audited.

The separation of prose from proof is the simplest Protocol move that advances these requirements at once: it removes prose from what must be optimized and verified, while making reliability signals depend on checkable evidence rather than on fluent surface form.

5 Discussion

From Production to Synthesis: Shifting Epistemic Agency. In low-volume settings, scholarly communication largely follows a “push” model: authors produce static papers, and communities converge on shared reference points by reading overlapping sets of canonical texts. Under overload ($\rho \gg 1$), overlap inevitably shrinks. Readers can no longer track even the high-visibility slice of the record without assistance, so filtering and synthesis increasingly move to machine-mediated workflows. The effective unit of consumption shifts from a single, shared manuscript to a task-specific synthesis of the public record, and different readers form different local posteriors q_j from different information sets \mathcal{I}_j .

This decentralisation is not inherently harmful—science has always involved specialisation and partial views—but it raises the stakes for the reliability layer. If synthesis operates over prose-only narratives, it is forced to rely on superficial cues that are cheap to optimize and therefore easy to manipulate. In that regime, personalised synthesis can amplify both benign noise (uncertainty, heterogeneity) and adversarial noise

(strategically generated text), pulling local beliefs away from the ground truth Θ .

The Protocol we outline addresses this by anchoring synthesis to verifiable structure. By separating prose from proof, it allows machines to render and personalise narrative while keeping accountability attached to structured claims, evidence, and uncertainty. Graded certification, verification-aware routing, and explicit congestion control then provide the higher-rate public signals and stable verification pipeline needed to make robust distributed inference feasible at scale.

Novelty, Exploration, and the Time Horizon. Novelty is often treated as a separate axis from reliability, as if systems must choose between “safe” work and “exploratory” work. Our view is that this split is unnecessary once publishing is framed as sequential, distributed inference. A mechanism that only optimizes *snapshot* risk will indeed favor conservative claims with immediate, high-confidence evidence. But over longer horizons it can raise cumulative risk by inducing epistemic stagnation: if uncertain but informative claims are systematically filtered out, the community fails to generate the new evidence required to reduce uncertainty about Θ in the future. Within the unified objective equation 5, supporting novelty therefore means keeping exploratory work *visible but explicitly uncertain* (archiving without over-certifying), and then allocating scarce verification to the claims whose proof layer can plausibly turn uncertainty into reliable knowledge. This is precisely the role of graded certification and re-admission to verification queues as evidence accumulates: exploration is internalized as a time-horizon tradeoff rather than bolted on as a second track.

Machines and AI as Scientific Agents. The model is intentionally agent-neutral. Encoders need not be humans; they can be automated laboratories, simulation pipelines, or AI agents that propose hypotheses, run experiments, and emit structured Claim and Evidence Records. Decoders likewise need not be humans; they can be retrieval and synthesis agents that form posteriors from routed evidence under strict compute and latency budgets. If anything, the AI-heavy regime sharpens our central diagnosis: *prose becomes cheap and abundant*, while verification (running code, checking provenance, reproducing experiments, and auditing dependencies) remains scarce. The Protocol’s separation of prose from proof therefore becomes more, not less, important: it gives machine readers an unambiguous, checkable substrate and makes “novelty” legible through evidence rather than through stylistic novelty in text. At the same time, an AI-heavy ecosystem introduces new stresses—for example, correlated errors when many agents share the same models or training data—which makes independent verification and provenance even more valuable as a diversity and robustness primitive.

Limitations and Outlook. Our model is deliberately mechanism-level. It abstracts away from governance, incentives, and user-interface design, not because those are unimportant, but because they are rapidly evolving and difficult to forecast. The value of the model is that it isolates non-negotiable constraints, namely finite attention, scarce verification effort, and strategic contamination, and translates them into requirements that any viable implementation of the Protocol must satisfy.

References

- [1] Michael Spence. Job market signaling. *The Quarterly Journal of Economics*, 87(3):355–374, 1973.
- [2] Weixin Liang, Yaohui Zhang, Zhengxuan Wu, Haley Lepp, Wenlong Ji, Xuandong Zhao, Hancheng Cao, Sheng Liu, Siyu He, Zhi Huang, Diyi Yang, Christopher Potts, Christopher D. Manning, and James Y. Zou. Mapping the increasing use of LLMs in scientific papers. arXiv:2404.01268, 2024. URL <https://arxiv.org/abs/2404.01268>.
- [3] Weixin Liang, Yaohui Zhang, Zhengxuan Wu, Haley Lepp, Wenlong Ji, Xuandong Zhao, Hancheng Cao, Sheng Liu, Siyu He, Zhi Huang, Diyi Yang, Christopher Potts, Christopher D. Manning, and James Y. Zou. Quantifying large language model usage in scientific papers. *Nature Human Behaviour*, 2025. doi: 10.1038/s41562-025-02273-8.
- [4] Keigo Kusumegi, Xinyu Yang, Paul Ginsparg, Mathijs de Vaan, Toby Stuart, and Yian Yin. Scientific production in the era of large language models. arXiv:2601.13187, 2026. URL <https://arxiv.org/abs/2601.13187>.
- [5] Loyal Liverpool. Ai intensifies fight against paper mills. *Nature*, 618(7964):222–3, 2023.
- [6] Osama Mohammed Afzal, Preslav Nakov, Tom Hope, and Iryna Gurevych. Beyond" not novel enough": Enriching scholarly critique with llm-assisted feedback. *ECAL (main), accepted*, 2025.
- [7] Sushil Bikhchandani, David Hirshleifer, and Ivo Welch. A theory of fads, fashion, custom, and cultural change as informational cascades. *Journal of political Economy*, 100(5):992–1026, 1992.

- [8] Daron Acemoglu, Munther A Dahleh, Ilan Lobel, and Asuman Ozdaglar. Bayesian learning in social networks. *The Review of Economic Studies*, 78(4):1201–1236, 2011.
- [9] Christopher A Sims. Implications of rational inattention. *Journal of monetary Economics*, 50(3):665–690, 2003.
- [10] Charles AE Goodhart. Problems of monetary management: the uk experience. *Monetary theory and practice: The UK experience*, pages 91–121, 1984.
- [11] David Manheim and Scott Garrabrant. Categorizing variants of Goodhart’s law. *arXiv preprint*, 2019.
- [12] Donald T Campbell. Assessing the impact of planned social change. *Evaluation and program planning*, 2(1): 67–90, 1979.
- [13] John PA Ioannidis. Why most published research findings are false. *PLoS medicine*, 2(8):e124, 2005.
- [14] Paul E Smaldino and Richard McElreath. The natural selection of bad science. *Royal Society open science*, 3(9): 160384, 2016.
- [15] Nicholas Bloom, Charles I Jones, John Van Reenen, and Michael Webb. Are ideas getting harder to find? *American Economic Review*, 110(4):1104–1144, 2020.
- [16] Claude E Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948.
- [17] Thomas A Courtade and Tsachy Weissman. Multiterminal source coding under logarithmic loss. *IEEE Transactions on Information Theory*, 60(1):740–761, 2013.
- [18] John N Tsitsiklis. Decentralized detection by a large number of sensors. *Mathematics of Control, Signals and Systems*, 1(2):167–182, 1988.
- [19] R AHLWEDE and I CSISZAR. Hypothesis testing with communication constraints. *IEEE transactions on information theory*, 32(4):533–542, 1986.
- [20] B HUGHES and PRAKASH NARAYAN. Gaussian arbitrarily varying channels. *IEEE transactions on information theory*, (2):267–284, 1987.
- [21] B. Nihar Shah. An overview of challenges, experiments, and computational solutions in peer review (extended version). *preprint*, 2025. URL <https://www.cs.cmu.edu/~nihars/preprints/SurveyPeerReview.pdf>.
- [22] Iliia Kuznetsov, Jan Buchmann, Max Eichler, and Iryna Gurevych. Revise and resubmit: An intertextual model of text-based collaboration in peer review. *Computational Linguistics*, 48(4):949–986, 2022.
- [23] Iliia Kuznetsov, Osama Mohammed Afzal, Koen Dercksen, Nils Dycke, Alexander Goldberg, Tom Hope, Dirk Hovy, Jonathan K Kummerfeld, Anne Lauscher, Kevin Leyton-Brown, et al. What can natural language processing do for peer review? *arXiv preprint arXiv:2405.06563*, 2024.
- [24] Association for Computational Linguistics. ACL rolling review (ARR): Overview and author guidelines. ACL Rolling Review website, 2025. URL <https://aclrollingreview.org/>.
- [25] International Conference on Learning Representations. ICLR 2026 author guide and review process. ICLR website, 2025. URL <https://iclr.cc/Conferences/2026/AuthorGuide>.
- [26] International Joint Conference on Artificial Intelligence. IJCAI–ECAI 2026 primary paper initiative. IJCAI-ECAI 2026 website, 2025. URL <https://2026.ijcai.org/primary-paper-initiative/>.
- [27] International Joint Conference on Artificial Intelligence. Call for papers: IJCAI–ECAI 2026 (main track). IJCAI-ECAI 2026 website, 2025. URL <https://2026.ijcai.org/ijcai-ecai-2026-call-for-papers-main-track/>.
- [28] Shakked Noy and Whitney Zhang. Experimental evidence on the productivity effects of generative artificial intelligence. *Science*, 381(6654):187–192, 2023.
- [29] Keigo Kusumegi, Xinyu Yang, Paul Ginsparg, Mathijs de Vaan, Toby Stuart, and Yian Yin. Scientific production in the era of large language models. *Science*, 390(6779):1240–1243, 2025.
- [30] Weixin Liang. Computational approaches to understanding large language model impact on writing and information ecosystems. *arXiv preprint arXiv:2506.17467*, 2025.
- [31] Qianyue Hao, Fengli Xu, Yong Li, and James Evans. Artificial intelligence tools expand scientists’ impact but contract science’s focus. *Nature*, pages 1–7, 2026.

- [32] Gary Charness, Brian Jabarian, and John A List. The next generation of experimental research with llms. *Nature Human Behaviour*, pages 1–3, 2025.
- [33] Yang-Hui He. Ai-driven research in pure mathematics and theoretical physics. *Nature Reviews Physics*, 6(9): 546–553, 2024.
- [34] Ryan Liu and Nihar B Shah. Reviewgpt? an exploratory study on using large language models for paper reviewing. *arXiv preprint arXiv:2306.00622*, 2023.
- [35] Alexander Goldberg, Ihsan Ullah, Thanh Gia Hieu Khuong, Benedictus Kent Rachmat, Zhen Xu, Isabelle Guyon, and Nihar B Shah. Usefulness of llms as an author checklist assistant for scientific papers: Neurips’24 experiment. *arXiv preprint*, 2024.
- [36] Nils Dycke, Matej Zečević, Ilija Kuznetsov, Beatrix Suess, Kristian Kersting, and Iryna Gurevych. Stricta: Structured reasoning in critical text assessment for peer review and beyond. *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 22687–22727, 2025.
- [37] Nils Dycke and Iryna Gurevych. Automatic reviewers fail to detect faulty reasoning in research papers: A new counterfactual evaluation framework. *arXiv preprint arXiv:2508.21422*, 2025.
- [38] Sarina Xi, Vishisht Rao, Justin Payan, and Nihar B Shah. Flaws: A benchmark for error identification and localization in scientific papers. *arXiv preprint arXiv:2511.21843*, 2025.
- [39] Qiyao Wei, Samuel Holt, Jing Yang, Markus Wulfmeier, and Mihaela van der Schaar. The ai imperative: Scaling high-quality peer review in machine learning. *arXiv preprint arXiv:2506.08134*, 2025.
- [40] Subhabrata Dutta, Timo Kaufmann, Goran Glavaš, Ivan Habernal, Kristian Kersting, Frauke Kreuter, Mira Mezini, Iryna Gurevych, Eyke Hüllermeier, and Hinrich Schuetze. Problem solving through human–ai preference-based cooperation. *Computational Linguistics*, 51(4):1337–1372, 2025.
- [41] James Zou. Chatgpt is transforming peer review—how can we use it responsibly. *Nature*, 635(8037):10–10, 2024.
- [42] Miryam Naddaf. Ai is transforming peer review—and many scientists are worried. *Nature*, 639(8056):852–854, 2025.
- [43] Tiffany I Leung. Llms in peer review—how publishing policies must advance. *JAMA Network Open*, 9(1): e2552042–e2552042, 2026.
- [44] Weixin Liang, Zachary Izzo, Yaohui Zhang, Haley Lepp, Hancheng Cao, Xuandong Zhao, Lingjiao Chen, Haotian Ye, Sheng Liu, Zhi Huang, et al. Monitoring AI-modified content at scale: A case study on the impact of ChatGPT on AI conference peer reviews. *International Conference on Machine Learning*, pages 29575–29620, 2024.
- [45] Elizabeth Gibney. Scientists hide messages in papers to game ai peer review. *Nature*, 643(8073):887–888, 2025.
- [46] Byungjin Choi, Tae Joon Jun, Joung Won Sung, Il Woo Park, Jeong-Moo Lee, Soo Ick Cho, Hyung Jun Park, Ro Woon Lee, and Jungyo Suh. Invisible text injection and peer review by ai models. *JAMA Network Open*, 9(1):e2552099–e2552099, 2026.
- [47] Weixin Liang, Mert Yuksekogunul, Yining Mao, Eric Wu, and James Zou. Gpt detectors are biased against non-native english writers. *Patterns*, 4(7), 2023.
- [48] Tim Baumgärtner and Iryna Gurevych. Scicoqa: Quality assurance for scientific paper–code alignment. *arXiv preprint*, 2026.
- [49] Lei You, Lele Cao, and Iryna Gurevych. Preventing the collapse of peer review requires verification-first ai. *arXiv preprint arXiv:2601.16909*, 2026.

Appendix: Scholarly Communication as an Engineering Discipline

A Proofs

We provide proofs for the results stated in Section 4. Throughout, we use standard identities for entropy and mutual information.

A.1 Proof of Theorem 1

Proof. Fix a realization of the time window with m items and labels A_1, \dots, A_m . Let $\mathcal{H} = \{i : A_i = 1\}$ denote the set of high-priority indices so that $|\mathcal{H}| = m_{\text{hi}}$. Each reader processes at most k_j items, so across all readers the total number of processed item instances is at most $K = \sum_{j=1}^N k_j$.

We write $R_i \in \{0, 1\}$ for the indicator that at least one reader processes item i . If $R_i = 1$, then the community must spend at least one unit of the global decoding budget on item i . Therefore the number of distinct items that can satisfy $R_i = 1$ is at most K . In particular,

$$\sum_{i \in \mathcal{H}} R_i \leq \min\{m_{\text{hi}}, K\}.$$

If $m_{\text{hi}} \leq K$, then the right-hand side equals m_{hi} and equation 8 holds because $\text{Cov} \leq 1$. If $m_{\text{hi}} > K$, then dividing the inequality by m_{hi} yields

$$\text{Cov} = \frac{1}{m_{\text{hi}}} \sum_{i \in \mathcal{H}} R_i \leq \frac{K}{m_{\text{hi}}}.$$

Combining the two cases gives $\text{Cov} \leq \min\{1, K/m_{\text{hi}}\}$. Substituting $\rho = m_{\text{hi}}/K$ yields the equivalent form $\text{Cov} \leq \min\{1, 1/\rho\}$. \square

A.2 Proof of Theorem 2

Proof. By the definition of mutual information,

$$I(\Theta; A) = H(A) - H(A | \Theta) \leq H(A).$$

If A is binary with $\mathbb{P}(A = 1) = p$, then $H(A) = h(p)$, where $h(p) = -p \log p - (1 - p) \log(1 - p)$. This proves equation 9. \square

A.3 Proof of Corollary 1

Proof. We first show the coverage statement. In the m -item window, the number of high-priority items equals $m_{\text{hi}} = mp_m$ by the definition of p_m . Theorem 1 therefore yields

$$\text{Cov} \leq \min \left\{ 1, \frac{K}{mp_m} \right\}.$$

If p_m is bounded below by a positive constant and K is bounded above by a constant, then $K/(mp_m) \rightarrow 0$ as $m \rightarrow \infty$. The displayed inequality implies that $\text{Cov} \rightarrow 0$.

We now show the information statement. Suppose the policy keeps coverage bounded below by a constant $c > 0$ for all large m . Then $\text{Cov} \geq c$ and Theorem 1 implies $c \leq K/(mp_m)$ for all large m . Rearranging gives $p_m \leq K/(cm)$, and therefore $p_m \rightarrow 0$ as $m \rightarrow \infty$.

Finally, draw an item uniformly at random from the m -item window and write A for its label. By definition, $\mathbb{P}(A = 1) = p_m$. Theorem 2 then implies

$$I(\Theta; A) \leq h(p_m).$$

Since $p_m \rightarrow 0$ and the binary entropy function satisfies $h(p) \rightarrow 0$ as $p \rightarrow 0$, we conclude that $I(\Theta; A) \rightarrow 0$. \square

A.4 Proof of Theorem 3

Proof. We apply Bayes' rule. First,

$$\mathbb{P}(\Theta = 1 \mid A = 1) = \frac{\mathbb{P}(A = 1 \mid \Theta = 1)\mathbb{P}(\Theta = 1)}{\mathbb{P}(A = 1)} = \frac{\alpha\pi}{\mathbb{P}(A = 1)}.$$

Next, we expand the denominator using the law of total probability,

$$\mathbb{P}(A = 1) = \mathbb{P}(A = 1 \mid \Theta = 1)\mathbb{P}(\Theta = 1) + \mathbb{P}(A = 1 \mid \Theta = 0)\mathbb{P}(\Theta = 0) = \alpha\pi + \beta(1 - \pi).$$

Substituting this into the first display gives equation 10.

For the second statement, the inequality $\mathbb{P}(\Theta = 1 \mid A = 1) \geq q$ is equivalent to

$$\frac{\alpha\pi}{\alpha\pi + \beta(1 - \pi)} \geq q.$$

All terms are nonnegative and the denominator is positive, so we can multiply both sides by $\alpha\pi + \beta(1 - \pi)$ and rearrange,

$$\alpha\pi \geq q\alpha\pi + q\beta(1 - \pi), \quad (1 - q)\alpha\pi \geq q\beta(1 - \pi), \quad \beta \leq \frac{\alpha\pi(1 - q)}{q(1 - \pi)}.$$

This is equation 11. □

A.5 Proof of Theorem 4

Proof. Bayes' rule implies an odds form for the posterior. For any realization of E , we have

$$\frac{\mathbb{P}(\Theta = 1 \mid E)}{\mathbb{P}(\Theta = 0 \mid E)} = \frac{\mathbb{P}(\Theta = 1)}{\mathbb{P}(\Theta = 0)} \cdot \frac{P(E \mid \Theta = 1)}{P(E \mid \Theta = 0)} = \frac{\pi}{1 - \pi} \text{LR}(E).$$

If $\mathbb{P}(\Theta = 1 \mid E) \geq q$, then $\mathbb{P}(\Theta = 0 \mid E) \leq 1 - q$, and therefore

$$\frac{\mathbb{P}(\Theta = 1 \mid E)}{\mathbb{P}(\Theta = 0 \mid E)} \geq \frac{q}{1 - q}.$$

Combining the last two displays yields

$$\frac{\pi}{1 - \pi} \text{LR}(E) \geq \frac{q}{1 - q}.$$

Taking logarithms on both sides gives equation 12. □

A.6 Proof of Theorem 5

Proof. Let $W(T)$ denote the total verification effort required by tasks that arrive in the time interval $[0, T]$. Under Assumption 6, the expected number of arrivals in $[0, T]$ is λT , and each arrival requires expected effort τ . Linearity of expectation gives

$$\mathbb{E}[W(T)] = \lambda T \tau.$$

Let $C(T)$ denote the total verification effort the system can supply in $[0, T]$. By assumption, $C(T) \leq \mu T$ deterministically. Let $L(T)$ denote the remaining unfinished verification workload at time T , measured in the same units. Since completed work cannot exceed supplied work, we have the basic inequality

$$L(T) \geq W(T) - C(T) \geq W(T) - \mu T.$$

Taking expectations gives

$$\mathbb{E}[L(T)] \geq \mathbb{E}[W(T)] - \mu T = (\lambda\tau - \mu)T.$$

If $\lambda\tau > \mu$, the right-hand side diverges to infinity as $T \rightarrow \infty$. Therefore the expected backlog cannot remain bounded. □

A.7 Proof of Theorem 6

Proof. Let p_1 and p_0 denote densities of P_1 and P_0 with respect to a common dominating measure. With equal priors, any decision rule partitions the feature space into a region \mathcal{D}_1 where it declares source 1 and a region \mathcal{D}_0 where it declares source 0. The probability of error under such a rule is

$$P_e = \frac{1}{2} \int_{\mathcal{D}_0} p_1(f) df + \frac{1}{2} \int_{\mathcal{D}_1} p_0(f) df.$$

The Bayes-optimal rule chooses $\mathcal{D}_1 = \{f : p_1(f) \geq p_0(f)\}$, which yields

$$P_e^* = \frac{1}{2} \int \min\{p_1(f), p_0(f)\} df.$$

For any nonnegative numbers a and b , we have $\min\{a, b\} = (a + b - |a - b|)/2$. Applying this identity pointwise gives

$$\int \min\{p_1(f), p_0(f)\} df = \frac{1}{2} \int (p_1(f) + p_0(f)) df - \frac{1}{2} \int |p_1(f) - p_0(f)| df.$$

Since $\int p_1 = \int p_0 = 1$, the first integral equals 2. By definition, the total variation distance is $\text{TV}(P_1, P_0) = \frac{1}{2} \int |p_1(f) - p_0(f)| df$. Substituting these identities yields

$$P_e^* = \frac{1}{2} (1 - \text{TV}(P_1, P_0)),$$

which is equation 13. □

A.8 Proof of Corollary 2

Proof. Let g be a screening rule and define the acceptance decision $A = g(F)$. Let $E = \{f : g(f) = 1\}$ be the acceptance event in the feature space. By definition, $\alpha = \mathbb{P}(A = 1 \mid \Theta = 1) = P_1(E)$ and $\beta = \mathbb{P}(A = 1 \mid \Theta = 0) = P_0(E)$.

The total variation distance satisfies

$$\text{TV}(P_1, P_0) = \sup_B |P_1(B) - P_0(B)|,$$

where the supremum ranges over all measurable events B . In particular, we may choose $B = E$ to obtain

$$\begin{aligned} P_1(E) - P_0(E) &\leq \sup_B (P_1(B) - P_0(B)) \\ &\leq \sup_B |P_1(B) - P_0(B)| = \text{TV}(P_1, P_0). \end{aligned}$$

Substituting $P_1(E) = \alpha$ and $P_0(E) = \beta$ yields equation 14. If $\beta \leq \beta_0$, then

$$\alpha = (\alpha - \beta) + \beta \leq \text{TV}(P_1, P_0) + \beta_0,$$

which gives the equivalent statement. □

A.9 Proof of Theorem 7

Proof. Because Θ is binary, the expected score under belief $\mathbb{P}(\Theta = 1) = p$ is

$$\mathbb{E}[S(q, \Theta)] = p \log q + (1 - p) \log(1 - q).$$

This is a concave function of q on $(0, 1)$ because its second derivative equals

$$\frac{d^2}{dq^2} \mathbb{E}[S(q, \Theta)] = -\frac{p}{q^2} - \frac{1-p}{(1-q)^2} < 0.$$

We therefore maximize the expected score by setting the first derivative to zero,

$$\frac{d}{dq} \mathbb{E}[S(q, \Theta)] = \frac{p}{q} - \frac{1-p}{1-q} = 0,$$

which implies $p(1 - q) = (1 - p)q$ and therefore $q = p$. Concavity shows that this maximizer is unique. □

B A Claim and Evidence Record for This Paper

This appendix encodes the present article as a Claim and Evidence Record (CER) in the sense of Definition 1. The goal is not to claim that a theoretical perspective is “self-verifying,” but to demonstrate what it would look like for a paper to expose its own accountable structure: what it claims, what in the artifact supports those claims, what remains uncertain, and what a verifier could check.

CER: *Scholarly Communication as an Engineering Discipline* (meta-record)

Record identifier. `cer:scadi.protocol.2026`

Type. Mechanism-design / theory

Scope. Scholarly communication under attention and verification constraints, with strategic contamination.

C: Claims.

C1 Robust-inference framing. Scholarly publishing in the post-LLM regime can be modeled as a robust distributed inference problem, in which a publishing mechanism Π chooses record format, routing, and verification to minimize expected log-loss risk under worst-case contamination (Eq. equation 5).

C2 Coverage cliff under overload. Under binary gatekeeping and fixed community attention budget K , the fraction of high-priority items that are processed is upper bounded by $\min\{1, K/m_{\text{hi}}\}$ (Theorem 1); when throughput grows faster than attention, coverage of even the “best” slice collapses unless selectivity vanishes (Corollary 1; Fig. 1).

C3 One-bit acceptance has vanishing information when rare. If acceptance becomes highly selective, the mutual information carried by a binary accept/reject label goes to zero (Theorem 2), limiting its usefulness for routing and collective belief update.

C4 Base-rate fragility. When the base rate of “true” claims falls, posterior reliability after acceptance requires extremely low false-accept rates (Theorem 3), implying that evidence thresholds must grow (Theorem 4; Fig. 2).

C5 Verification is a scarce service with stability constraints. Modeling deep checks as a queue yields a feasibility condition: if workload intensity exceeds capacity, backlog diverges unless the system performs admission control and traffic shaping (Theorem 5).

C6 Text-only screening becomes statistically fragile. If an adversary can match the feature distribution of “reliable” records in a given feature space, no classifier can achieve low error; the limit is controlled by total variation distance (Theorem 6; Fig. 3). Consequently, protocols should shift screening/routing from prose-level features toward proof-bearing structure (CER fields and verifiable artifacts) that increase separability in the relevant feature space.

C7 Soft certification can be made incentive-compatible. If the system stores probabilistic assessments, proper scoring rules (log score) make truthful reporting optimal (Theorem 7; Fig. 5).

E: Evidence bundle (verifiable artifacts and pointers).

E1 Formal specification. The mechanism-design objective and definitions in Section 3, including Eq. equation 4–equation 5 and Definition 1.

E2 Theorems + proofs. Statements in Section 4 with complete proofs in Appendix A.

E3 Executable figures. All main figures are generated from committed L^AT_EX/TikZ source (`figures/*.tex`); a verifier can reproduce them by compiling the manuscript, and can inspect the plotting code for mathematical consistency.

E4 Empirical/contextual anchors. Cited measurements and policy documents on LLM usage, paper mills, and evolving review workflows (e.g., [3, 5, 24, 25]).

U: Uncertainty summary (what is not claimed).

- The results are *conditional*: they hold under explicit assumptions about attention/verification budgets,

routing constraints, and the presence of strategic contamination.

- The paper primarily provides *limits and design requirements*, not a calibrated quantitative forecast of parameters such as K , ϵ , or real-world evidence costs.
- The model abstracts away several realities (field heterogeneity, correlated agent errors, institutional incentives beyond the scoring-rule primitive). These factors can change which mechanisms are implementable, even if the information-theoretic limits remain.

σ : **Signature.** Authors and affiliations on the title page. In a deployed Protocol, σ would be a cryptographic signature binding this CER (and its artifacts) to the responsible parties.

v : **Version.** Manuscript version as distributed by the authors; figures and proofs in the same source tree.

Canonical serialization (illustrative).

```
cer_id: cer:scadi.protocol.2026
title: Scholarly Communication as an Engineering Discipline
theta:
  description: Robust inference about scientific states under overload and contamination
claims:
  - id: C1
    statement: Publishing can be framed as robust distributed inference (Eq. (5)).
    evidence: [E1, E2]
  - id: C2
    statement: Binary gatekeeping induces a coverage cliff under overload (Thm. 1; Fig.
      1).
    evidence: [E2, E3]
  - id: C3
    statement: Selective accept/reject labels carry vanishing information when rare (Thm
      . 2).
    evidence: [E2]
  - id: C4
    statement: Low base rates force extreme false-accept control and higher evidence
      thresholds (Thm. 3--4; Fig. 2).
    evidence: [E2]
  - id: C5
    statement: Verification is a scarce service; stability requires admission control (
      Thm. 5).
    evidence: [E2]
  - id: C6
    statement: Text-only screening becomes indistinguishable under feature matching;
      structure increases separability (Thm. 6; Fig. 4--5).
    evidence: [E2, E3]
  - id: C7
    statement: Probabilistic certification can be made incentive-compatible via log
      scoring (Thm. 7; Fig. 6).
    evidence: [E2]
uncertainty:
  - Results are conditional on explicit budget and contamination assumptions.
  - Parameter magnitudes (K, epsilon, costs) are not empirically calibrated here.
  - Institutional and behavioral dynamics may constrain implementability.
signature:
  human_readable: authors + affiliations
  cryptographic: "(protocol-dependent)"
version: v1
```