

Article

Not peer-reviewed version

---

# Flow-by-Flow: Content-Judgment Bypass for Governing AI Output in High-Loss Domains

---

[Hiroki Naito](#)\*

Posted Date: 28 April 2026

doi: 10.20944/preprints202604.1948.v1

Keywords: AI governance; flow-rate limitation; cognitive cost score; human-in-the-loop; content-judgment bypass; high-loss domains



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC, OpenAlex.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

# Flow-by-Flow: Content-Judgment Bypass for Governing AI Output in High-Loss Domains

Hiroki Naito

UTIE Research Institute (UTIE Instruments Inc.); h.naito@utie-instruments.com

## Abstract

Prior work showed that human-in-the-loop oversight becomes structurally untenable in high-loss domains when AI output velocity  $V$  exceeds human cognitive capacity  $C_{\max}$ . The operative constraint, however, is not  $V$  alone but  $V \times L$ , where  $L$  denotes per-item cognitive load.  $L$  consists of triage, judgment, and response, which respond asymmetrically to AI capability improvement. Triage cost does not decline as models become more capable, because semantic indeterminacy is inherent in general-purpose design. Response cost is invariant to accuracy improvements. Only judgment cost faces downward pressure, and this pressure often operates by inducing omission rather than genuine reduction. Capability improvement therefore restructures  $L$  rather than reducing it. Governance mechanisms based on evaluating whether AI output is correct either delegate that evaluation to AI and inherit hallucination risk, or delegate it to humans and face the  $V \times L$  ceiling. We propose Flow-by-Flow, a governance paradigm that controls supervisory load without evaluating content. A cognitive cost score based on formal, countable features imposes nonlinear costs on high-volume production, while an institutional capacity cap keeps processing volume within  $C_{\max}$ . We derive four design invariants for any content-judgment-bypass exceedance pathway: no content judgment, no scalable consumption of examiner capacity, identity-bound per-application friction, and no batch clearance. One reference implementation is discussed to show that these invariants are jointly satisfiable, while its practical difficulties are explicitly acknowledged. An illustrative Monte Carlo analysis across 1,000 parameter draws suggests that composite multi-metric flow control outperforms supervision reinforcement alone in 90.8% of trials.

**Keywords:** AI governance; flow-rate limitation; cognitive cost score; human-in-the-loop; content-judgment bypass; high-loss domains

---

## 1. Introduction

Contemporary debates on AI governance still rest on the assumption that generative AI systems can be operated safely even as their output expands, so long as a human being who bears final judgment and responsibility remains within the decision-making loop. The human-oversight requirement for high-risk AI systems under the EU AI Act, the concept of human involvement in the U.S. NIST AI Risk Management Framework, and the familiar statement in many corporate deployment guidelines that “a human must always make the final check” all depend on this assumption.

Critiques of this assumption have accumulated historically across multiple scholarly lineages. Bainbridge (1983) observed that as automation becomes more advanced, the tasks left to human operators converge toward two activities that are least suited to human cognitive characteristics: the monotonous monitoring of a normally functioning system and extreme intervention under unfamiliar failure conditions. Perrow (1984) showed that, in highly complex and tightly coupled systems, accidents are not merely defects in design but systemic events, a view developed as Normal Accident Theory. Beck (1986) argued that the development of modern technology endogenously generates risks that exceed existing institutional capacities for control. Parasuraman and Riley (1997) systematized the categories of use, misuse, disuse, and abuse of automation, and Parasuraman and

Manzey (2010) proposed an integrated attentional model of automation bias and complacency. Elish (2019) named the structure in which humans absorb responsibility without retaining substantive control over automated systems the “moral crumple zone.”

These insights were developed primarily with physical automation systems in mind. In generative AI, however, the same problem reappears in a qualitatively different form. In physical automation, the object monitored by humans is the operation of a machine, and errors become visible as deviations in physical quantities such as temperature or pressure. The normal range can be defined objectively, and the deviation criteria for identifying abnormalities exist outside the supervisor. By contrast, in generative AI, the objects that supervisors must inspect are cognitive and inferential artifacts, such as text, media, and code. Errors appear as failures of semantic or contextual coherence, or as logical leaps. The deviation criteria therefore depend on the internal model of the human supervisor.

Because the criteria are intrinsic rather than extrinsic, a distinctive vulnerability arises. Outputs from generative AI exhibit stronger statistical regularity than human outputs. Supervisors who are repeatedly exposed to them learn and internalize this regularity as the normal state through predictive-error minimization mechanisms (Friston, 2010; Clark, 2013). As predictive error declines, this is experienced at the cognitive level as increased processing fluency (Alter & Oppenheimer, 2009), that is, as a diminished sense that something is wrong with the output. This process proceeds even when the supervisor is paying attention. It therefore differs qualitatively from conventional automation bias, which is usually understood as a failure of attention. Instead, it degrades error-detection capacity through a transformation of the cognitive frame itself.

Evidence supporting this conclusion has accumulated from several directions. Carnat (2024) argues that the human-oversight requirement in Article 14 of the EU AI Act cannot resolve automation bias in generative LLMs, and that technical improvements in hallucination may instead accelerate human overtrust. Horowitz and Kahn (2024), in an experiment involving 9,000 participants across nine countries, show that trust and confidence in AI systems are major drivers of automation bias. A systematic review of 35 studies from 2015 to 2025 (AI & Society, 2025) reports that interventions such as explanations and trust-calibration feedback are ineffective in reducing automation bias. Kücking et al. (2024) show that non-experts are more vulnerable in clinical decision-support systems. Park, Kim, and Han (2026), using a dynamic systems model, identify a critical threshold at which human capability collapses sharply ( $K^* \approx 0.85$ ), and show that mandatory practice intervals substantially preserve human capability. Bastani and Cachon (2025) derive conditions under which, as AI accuracy improves, the compensation required to economically motivate supervisory effort grows without bound, eliminating the possibility of a feasible contract design.

Empirical evidence is also accumulating on the effect of LLM adoption on output rate. Kusumegi et al. (2025), analyzing more than two million preprints, report that LLM adoption increases per-author output by 23.7% to 89.3%, depending on field and author background. Naito (2026), using approximately 25,000 preprint records, shows that across fields without physical-space constraints, output rates accelerated after 2023 at a scale that cannot be explained by AI investment or by an increase in AI researchers, with submission counts increasing by at least 30%.

Naito (2026) integrates these prior studies and derives the claim that the idea of responding to AI capability gains by strengthening or improving human checking regimes breaks down in high-loss domains when three constraints hold simultaneously. First, legal responsibility is fixed to humans ( $R = 1$ ). Second, human cognitive processing capacity has a biological upper bound ( $C_{max}$ ). The finitude of attentional resources (Kahneman, 1973), the capacity limits of working memory (Cowan, 2001), and the temporal decline of sustained attention (Warm, Parasuraman & Matthews, 2008) are well-established findings in cognitive science. Third, economic pressure causes the AI output rate to continue expanding nonlinearly. Once the output rate  $V$  exceeds  $C_{max}$ , HITL is transformed from substantive oversight into a formal procedure. In addition, the invisibilization of errors places the penalty function in a dormant state, so risk does not materialize through gradual adaptation via continuous feedback, but instead appears discontinuously as threshold shocks. Under

these conditions, the most directly controllable variable for keeping expected loss bounded in high-loss domains is the limitation of output rate. A flow-rate limitation such as  $V_{\text{eff}} = \min(V, C_{\text{max}})$  therefore becomes more necessary than further reinforcement of the human oversight layer.

These lines of research are converging on a shared diagnosis: there are limits to strengthening human oversight layers. Yet an unresolved problem remains in moving from this diagnosis to institutional design. The model in Naito (2026) remains coarse. Because  $V$  captures only the quantity of outputs, it cannot distinguish the cognitive burden per item. One thousand simple outputs and ten highly complex outputs are both treated as the same  $V$ , even though the cognitive burden imposed on supervisors is fundamentally different. Moreover, while the inequality  $V > C_{\text{max}}$  derives the hollowing-out of oversight, the earlier work only indicated a broad direction, such as limiting the number of submissions or applications, when addressing what form this inequality takes in real institutions and what institutional responses are available. To realize flow design, it is necessary not only to operationally estimate  $C_{\text{max}}$ , but also to make cognitive load quantitatively estimable and to develop a methodology in which the process of quantification itself does not depend on AI judgment. The bridge from the diagnostic conclusion that human oversight becomes nominal to an implementable institutional design remains an important open problem in AI governance. This paper addresses that problem.

The first step is theoretical refinement through the introduction of the variable  $L$ . The output rate  $V$  is extended to  $V \times L$ , where  $L$  denotes the per-item cognitive load. The previous model,  $V > C_{\text{max}}$ , describes, in simple terms, a condition in which “there are too many items for human supervisors to process.” In reality, however, even when the number of items is small, a human oversight regime can similarly break down if the cognitive cost required to verify and respond to each item is sufficiently large.  $L$  includes both the cost of checking whether an output is correct and the cost of actually responding to an output once it has been confirmed to be correct.

The second step is the presentation of the Flow-by-Flow approach. As an answer to the question of how to implement flow design, this paper proposes an institutional design that enables low-cost and appropriate flow-rate limitation while bypassing content judgment. The strength of this design lies in the fact that it does not evaluate the truth or falsity of output content at all. Instead, it automatically infers cognitive cost from formal quantitative features alone. In other words, it does not create room for the hallucination risks that constitute the greatest weakness when human supervisors rely on AI.

As in the previous work, the scope of this paper is limited to domains in which formal verification is inherently impossible and the loss caused by error is severe. Ordinary low-risk uses, as well as domains in which external criteria can be established through statistical verification, are outside the scope of this paper. However, even within high-loss domains, when the generator population is institutionally bounded, such as AI-assisted diagnosis by physicians or AI assistance for judges, the expansion of  $V \times L$  saturates once all users have adopted AI, and supervision enhancement may remain sufficient. The scope is therefore not determined by the name of a domain. Even within the military domain, AI-enabled landmines, which are strongly constrained by physical space and have an institutionally bounded generator population, and autonomous cyberattacks, which are completed within information space and may involve an unbounded submitter population, require opposite responses with respect to flow limitation.

The need for flow control in this paper arises when the generator population is open and  $V \times L$  has a structure that expands compoundly. This condition is assessed across domains. Flow design is necessary in domains that are both high-loss and characterized by an open generation structure. Existing AI governance frameworks, including the EU AI Act, primarily classify systems by the magnitude of loss. By contrast, what determines the need for flow control is not the magnitude of loss alone, but whether the generator population is institutionally bounded. In addition, the argument of this paper assumes continued improvement in the capabilities of general-purpose AI models and continued human use of AI. If AI capabilities cease to improve in the future, or if humanity stops using AI, the argument of this paper loses its scope.

**Box 1. Standard objections and the throughput problem.****1 If AI becomes more accurate, supervision will be easier.**

Higher accuracy may reduce some errors, but it does not eliminate  $p > 0$ . It can also increase  $V$  and the number of outputs requiring downstream response.

**2 AI can verify AI.**

AI verification creates a recursive verification problem. Under  $R = 1$ , the legal responsibility for accepting the verification still returns to a natural or legal person.

**3 We can add more human supervisors.**

Additional supervisors increase  $C_{\max}$  only linearly or sublinearly, while AI-enabled output  $V \times L$  can grow nonlinearly.

**4 Disclosure of AI use is enough.**

If disclosure cannot be independently verified at scale, supervisors must treat all submissions as potentially AI-mediated.

**5 Full intake is an access right.**

Access rights do not imply immediate processing rights. Other infrastructures routinely distinguish access from simultaneous throughput.

**6 Capital-rich entities can bypass physical gates.**

The aim is not perfect exclusion. The aim is to convert digital near-zero marginal cost into embodied time, movement, and coordination costs.

**7 Nominal oversight is acceptable if AI makes fewer mistakes than humans.**

High-loss domains are not governed by average error-rate comparison alone. If  $V$  increases faster than  $p$  decreases, the absolute number of catastrophic errors can rise. Nominal oversight also converts continuous correction into delayed threshold shocks and leaves humans or organizations as moral crumple zones.

**2. Theoretical Framework: From  $V$  to  $V \times L$** 

## Table of Symbols

Symbol	Definition
$V$	AI output rate, defined as the number of AI outputs per unit time.
$L$	Per-item cognitive load.
$C_{\max}$	The upper bound of cognitive load that a human supervisor or human supervisory organization can process per unit time.
$C_{\text{eff}}$	The effective cognitive processing capacity per unit time of a supervisor or supervisory organization. $C_{\text{eff}}$ is no greater than $C_{\max}$ .
$R$	Attribution of legal responsibility. $R = 1$ means that legal responsibility is fixed to a natural or legal person, not to the AI system itself.
$B$	The degree of physical-space constraint. A lower $B$ indicates that the process is more completely contained within information space. This is a qualitative variable.
$p$	Error probability, with $p > 0$ .
$S$	Cognitive cost score, defined as a weighted geometric product of normalized formal features.

k	The threshold exceedance multiplier of the cognitive cost score.
N	The number of feature dimensions that constitute the cognitive cost score.

### 2.1. Reconsidering the Previous Model

$$V(t) > C_{\max}, \text{ under } R = 1$$

This inequality in the previous work states that, once the AI output rate  $V$  exceeds the upper bound of human cognitive processing capacity,  $C_{\max}$ , human-in-the-loop oversight in high-loss domains loses its substantive function. In that formulation, however,  $V$  is a one-dimensional variable that represents the number of outputs. It abstracts away differences in the cognitive burden that supervisors must bear for each output.

From the perspective of real-world human review systems, this abstraction is significant. In patent examination, an application with three claims and a ten-page specification and an application with fifty claims, a 120-page specification, and eighty prior-art references clearly require entirely different levels of cognitive resources for examination. Similarly, in the supervision of legal documents generated by AI, the review of a standardized contract and the verification of a complex contractual scheme spanning multiple jurisdictions are likely to differ substantially in per-item cognitive load.

It should also be noted that  $R = 1$  has another function beyond fixing legal responsibility to humans. It also creates a pathway through which parties internalize the cost of fact-checking statements issued under their own names. This constitutes part of the epistemic rate-limiting mechanism underlying the macroscopic limit developed in the previous work. Since this is not the central subject of the present paper, it is left for separate treatment.

### 2.2. Introducing the Variable $L$

In the previous work,  $C_{\max}$  was defined as the upper bound of the cognitive processing capacity available to a human supervisor or supervisory organization. In the present paper, because the focus is on the product of output rate  $V$  and per-item cognitive load  $L$ ,  $C_{\max}$  is used as the maximum cognitive processing capacity per unit time.

$$V \times L \leq C_{\max}$$

Here,  $L$  denotes the cognitive load required for a responsible actor to place one AI output into a state in which it can be supervised.  $L$  is not merely reading time. It includes the load of determining what type of information the output should be treated as, the load of evaluating the content, and the load of carrying out the actual institutional response after evaluation.

The important point is that the number of items,  $V$ , alone cannot capture supervisory load. Even when the number of outputs is small, if the verification, judgment, and response required for each item are sufficiently burdensome,  $V \times L$  can exceed  $C_{\max}$ .

### 2.3. Components of $L$

$L$  is not a single homogeneous burden. The costs borne by humans in supervising AI outputs are a composite of tasks that differ in kind. The previous section introduced the distinction between verification cost and response cost, but before that distinction there is an additional layer of work.

Consider supervision in physical measurement. When a supervisor looks at a value displayed by a thermometer, the fact that the value is a temperature has already been fixed at the moment of display. What the supervisor must do is simply compare the observed value with an externally specified standard, namely the normal range for temperature. The higher-order question of what the object of judgment is, or which criterion should be invoked, does not arise. The supervisor's cognitive resources are concentrated on one task: comparing the value with the standard.

In generative AI, this higher-order question does arise. A text produced by AI does not, on its surface, determine whether it is a factual report, an inference, a list of possibilities, a fictional

narrative, a calculation result, or the author's opinion. Each category requires a different axis of evaluation. If the output is factual, the issue is truth or falsity. If it is an inference, the issue is plausibility. If it is a narrative, the issue is coherence. If it is a calculation, the issue is correctness. If it is an opinion, the issue is validity. The supervisor must first decide which axis should be used to read the output.

After making that decision, the supervisor must invoke the relevant criterion, evaluate the output against that criterion, and then take the appropriate response in light of the evaluation. When this sequence is divided according to differences in task type, it consists of three components: triage, judgment, and response.

Triage is the task of deciding what type of information the output should be treated as. Judgment is the task of evaluating the correctness or validity of the content under that assumption. Response is the task of taking an actual action after judgment, such as drafting a notice of reasons for rejection, deciding and explaining a treatment plan, or creating and applying a software patch.

These three tasks respond differently to improvements in AI accuracy. Judgment cost is subject to pressure toward omission as accuracy improves. Response cost does not decrease with accuracy, because it depends on physical time and human labor. Triage cost also does not decrease with accuracy.

Conventional machines such as thermometers, calculators, and clocks were designed so that the semantic category of their output was fixed in advance. By contrast, the central characteristics of generative AI are generality and human-likeness. Its semantic category is therefore not fixed. That each output imposes triage on the reader is an inherent consequence of the design of general-purpose generative AI models.

This distinction determines the behavior of L as a whole. The size of L is determined by the combined burden of triage difficulty, the amount of judgment required, and the volume of response. Improvements in AI accuracy may appear to reduce the amount of judgment required, but they increase the difficulty of triage and expand the volume of response. If higher accuracy increases the absolute number of correct outputs, it also increases the number of outputs that must be acted upon.

#### 2.4. *Toward Flow-by-Flow*

The problem addressed in this paper is not the average quality of AI outputs, but the upper bound of cognitive load that responsible actors can process per unit time. Improvements in AI accuracy, increases in the number of human supervisors, AI-based verification of AI, and disclosure regimes for AI use may each be useful in specific contexts. None of them, however, guarantees that  $V \times L$  will remain within  $C_{\max}$ .

Accuracy improvements may reduce  $p$ , but they can also increase  $V$  and the number of downstream responses. AI-based verification creates a recursive verification problem, namely the verification of the verification result. Disclosure regimes do not remove triage cost unless the truth of the disclosure can be independently verified. Thus, even if existing responses can improve the quality of content judgment, they do not remove the throughput constraint. Detailed responses to standard objections are provided in Appendix B. The next chapter presents Flow-by-Flow. Rather than judging the truth or falsity of content, Flow-by-Flow controls inflow so that content judgment remains possible within substantive human capacity.

### 3. Flow-by-Flow

#### 3.1. *Purpose of This Chapter*

The purpose of the Flow-by-Flow approach is simple: to prevent outputs containing hallucinations, factual distortions, logical leaps, or similar problems from passing through beyond the checking capacity of human supervisors. To achieve this purpose, Flow-by-Flow does not judge the correctness of output content. Instead, it controls the flow rate itself so that outputs remain within the range of human supervisory capacity.

The dominant orientation in current AI governance is to ask how accurately AI outputs can be evaluated. More advanced explainability techniques, more detailed audit logs, and more multilayered approval procedures all seek to secure safety by improving the quality of content judgment. Yet, as shown in Chapter 2, content judgment is a major component of  $L$ . Efforts to improve the quality of judgment may further increase  $L$  and thereby worsen the inequality  $V \times L > C_{\max}$ . Moreover, as argued in the supplementary information to the previous work, supervision-enhancing measures may not have improved the quality of supervision itself, but may instead have unintentionally suppressed  $V$  by increasing procedural friction.

Flow-by-Flow reverses the orientation of this problem. Rather than evaluating the content of outputs, it measures cognitive cost solely from formal quantitative features and controls the inflow so that it does not exceed human processing capacity. Human supervisory capacity has a biological upper bound,  $C_{\max}$ , and improvements in AI accuracy eliminate the economic incentive to exert supervisory effort (Bastani & Cachon, 2025). Incentive design on the supervisor side has no solution, either in terms of cognitive capacity or in terms of economic motivation.

The Flow-by-Flow design is based on the following single principle:

to implement, without content judgment, an asymmetric incentive design that makes AI-enabled mass production costly and AI-enabled concision inexpensive.

All components of this chapter are derived from this principle. The cognitive cost score imposes costs on each dimension of mass production, and the exceedance path converts that cost into institutional friction. The institutional capacity cap makes the finitude of processing slots explicit and creates a mechanism in which concise submissions are more likely to fit within the available capacity. For applications that exceed the cap, an additional mechanism is required to preserve per-application marginal cost without reintroducing content judgment.

None of these mechanisms judges the truth or falsity of output content. For this reason, they do not rely on semantic content judgment and therefore avoid the class of hallucination risks associated with AI-based content evaluation. For legitimate human applicants, all components of the system are low-cost. Even if a false positive occurs, meaning that a legitimate applicant exceeds the threshold, the applicant should face only a bounded, identity-bound per-application friction. For AI-enabled mass producers, by contrast, this friction scales with the number and complexity of applications. If they increase the number of submissions, they encounter submission-count limits. If they increase complexity, the score rises sharply. If they seek passage despite exceedance, the exceedance pathway imposes a non-zero marginal cost on each application.

Before presenting the specific components of the Flow-by-Flow design, we derive four design invariants that any mechanism occupying the role described above must satisfy. They are necessary conditions for any institutional mechanism that seeks to control flow rate while bypassing content judgment in high-loss domains.

**Invariant 1:** No substantive content judgment. The mechanism must not evaluate the truth, falsity, quality, or appropriateness of the content of outputs. Any mechanism that requires such evaluation reintroduces hallucination risk if delegated to AI, or the  $V \times L$  ceiling if delegated to humans.

**Invariant 2:** No scalable consumption of examiner capacity. The mechanism must not consume human examiner time in proportion to the number of submissions processed. If it does, the mechanism itself becomes subject to the  $C_{\max}$  constraint it is designed to protect.

**Invariant 3:** Identity-bound per-application friction. The cost imposed by the mechanism must be attached to each individual application and to a verified identity. Without identity binding, the cost can be distributed across anonymous accounts. Without per-application attachment, batch clearance becomes possible and the marginal cost of additional submissions returns to near zero.

**Invariant 4:** No batch clearance. No single action, credential, payment, or institutional certification may clear multiple applications simultaneously. Any batch-clearance route restores the near-zero marginal cost structure that the mechanism is designed to eliminate.

These four invariants are jointly necessary. A mechanism that satisfies only three of them is vulnerable to circumvention through the unsatisfied dimension. The cognitive cost score satisfies Invariants 1 and 2. The institutional capacity cap satisfies Invariant 2. Neither alone satisfies Invariants 3 and 4 for threshold-exceeding applications. The remaining question is what mechanism can satisfy all four invariants for applications that exceed the institutional capacity cap. One reference implementation is presented in Section 5.5, but it is not the only possible mechanism. Any mechanism that satisfies all four invariants would be a valid substitute.

### 3.2. Two Misconceptions: Full Intake and Processing as Soon as Possible

Before justifying the Flow-by-Flow design institutionally, it is necessary to identify a major problem in existing institutions. Knowledge infrastructures in high-loss domains, such as patent offices, courts, peer-review systems, and pharmaceutical review agencies, depend on the following syllogism.

- (a) High-loss domains are important social infrastructures.
- (b) Every human being has a right of access to important social infrastructures.
- (c) Therefore, every submitted application must be processed.

Statements (a) and (b) are correct. Statement (c), however, is a leap. What (b) guarantees is the right of anyone to submit an application. It does not guarantee the right to be processed as soon as possible.

This distinction is already established in other forms of social infrastructure. In medicine, triage is used: all patients are treated, but the order of treatment is determined by severity. In road traffic, all vehicles can use the road, but simultaneous passage is limited by traffic signals. In electricity supply, all consumers are supplied, but simultaneous demand is managed through planned load control during peak periods. In communications, all users can connect, but bandwidth is limited. In all of these infrastructures, the obligation to serve everyone is clearly distinguished from the obligation to process everyone simultaneously. The former is protected as an absolute right. The latter is managed under physical constraints.

Knowledge infrastructures, however, have not made this distinction. Patent offices accept all applications even when tens of thousands of applications concentrate within a year, and they operate under the principle of processing them as soon as possible. Similar principles apply to complaints submitted to courts and manuscripts submitted to peer-review systems. These institutions do not intentionally delay intake on the ground that the complaint or manuscript exceeds the cognitive resources of judges or reviewers.

The answer to this problem is to recognize that (c) was not a right, but an empirical regularity that happened to hold before AI. In an environment in which drafting one patent application required several months of human labor and writing one academic paper required months, the submission rate itself functioned as a natural flow limitation. Application volume exceeding the  $C_{\max}$  of knowledge infrastructures was physically unlikely to occur. The practice of processing all submissions within a given administrative period was therefore not designed as an institutional principle. It merely happened to hold under those conditions.

Generative AI has destroyed this natural flow limitation. In an environment in which a single user can draft dozens of patents or generate dozens of papers per year with AI, submission volume can exceed  $C_{\max}$  on a persistent basis. The sharp increase in monthly patent applications at the Japan Patent Office in December 2025, approximately 2.7 times the previous month, can be understood as an early manifestation of this change (Naito, 2026). Nonlinear increases in AI-enabled output have already been empirically documented across multiple domains. Kusumegi et al. (2025), analyzing more than two million preprints, conclude that LLM adoption increases per-author paper production by 23.7% to 89.3%. Naito (2026), using approximately 25,000 preprint records, complements this conclusion through a different approach. Submissions to major AI conferences have increased and are estimated to exceed three times their 2024 level by 2030.

Once it becomes clear that (c) was merely an empirical regularity, it becomes necessary to explicitly abandon (c) and institutionalize flow control in order to substantively protect the access right in (b). This is not a restriction of rights. On the contrary, maintaining full immediate processing without flow control risks violating access rights through the collapse of processing quality. If an institution accepts more items than it can process and legitimate applicants no longer receive substantive examination, the result is, in effect, not meaningfully different from being denied intake.

Older knowledge infrastructures have rested on two principles: the principle of full intake, under which all submitted applications are accepted, and the principle of processing as soon as possible, under which accepted applications are processed as quickly as possible. AI breaks these two principles through different mechanisms merely by increasing the number of submissions.

The principle of processing as soon as possible induces processing beyond  $C_{\max}$  by requiring the maximization of processing speed. If examiners are pressured to process faster, the degradation of  $C_{\text{eff}}$  discussed in the previous work occurs. The principle of full intake, by contrast, creates an unlimited divergence between intake volume and processing capacity under an explosion of  $V$ . Simply accepting all items and carrying them over causes the queue to expand indefinitely, leaving later legitimate applicants without substantive access.

The approach proposed here explicitly transforms both principles. In response to the principle of processing as soon as possible, the institutional capacity cap fixes the processing rate within  $C_{\max}$  and thereby preserves quality. In response to the principle of full intake, the approach controls the entry of submissions that exceed processing capacity.

### 3.3. Automatic Measurement of Cognitive Cost: Counting-Based Governance

The first component of Flow-by-Flow is a mechanism for automatically measuring the cognitive cost of each output without content judgment. Conventional AI governance has treated the evaluation of output quality as the foundation of institutional design. Yet quality evaluation necessarily entails content judgment, and content judgment is vulnerable to hallucination. If AI is asked to judge whether an output is safe, that judgment itself carries hallucination risk. If humans are asked to judge it, the system confronts the wall of  $V \times L > C_{\max}$ .

Our approach addresses this problem by not judging whether the content of an output is correct at all. It measures only the formal quantitative features of the output. In patent applications, these features include the number of claims, the number of characters or pages in the specification, and the number of cited prior-art references. In academic papers, they include word count, number of references, number of figures and tables, and volume of supplementary materials. In legal documents, they include the length of the complaint, the number of evidentiary items, and the number of relevant legal provisions. In pharmaceutical applications, they include the length of application documents, the number of clinical trials, and the number of indications.

What these features share is that measurement is completed simply by counting. These measurements do not refer to the truth or falsity of the output content. Even if AI is used for the measurement, there is no need for semantic content judgment. A machine may count incorrectly, but this is a verifiable error, such as counting 61 as 62. It differs in kind from a hallucination such as mistakenly judging that a claim is novel. The former can be mechanically corrected. The latter requires expert judgment.

The selection of features and the determination of weights are performed autonomously by each institution on the basis of its own processing-time data. Empirical relationships, such as the increase in average examination time when the number of claims increases by one, or when the specification increases by ten pages, provide initial values for the weights. This design naturally reflects differences across fields and institution-specific processing characteristics, without requiring centralized standard-setting.

The measured features are integrated into a cognitive cost score. The simplest form is a weighted product of each feature. The important point is that the raw values of features are not multiplied directly. Features such as number of claims, word count, number of citations, and number of figures

have different units and scales, so a raw product of heterogeneous quantities has no meaning. The weighted product in this paper means that each feature is first made dimensionless by dividing it by a median value in past data or by an institutional baseline value, and that the resulting feature ratios are then integrated according to their weights.

The initial values of the weights are estimated from past processing-time data in each domain. In patent examination, for example, data on average examination time by International Patent Classification already exist, and differences in complexity across fields can be reflected as coefficients. These coefficients are continuously calibrated through institutional operation. The system can begin with rough estimates in the first year and improve its accuracy as actual examination-time data accumulate.

The cognitive cost score does not need to be a precise proxy for L. Its function is to transform evasive optimization into a multidimensional constraint-satisfaction problem. This is achieved by designing the features so that each feature is inseparably connected to the content of the application, and so that compressing one dimension transfers burden to another. A lower bound is set for each feature that composes the cognitive cost score, in order to prevent the entire score from collapsing to zero when a particular feature is zero.

### 3.4. Evasive Behavior Through AI Optimization

A technical note is necessary here. This kind of AI-driven evasion of flow design is not speculative. It is a direct consequence of Goodhart's law: once a measure becomes a target, it ceases to be a good measure (Goodhart, 1975). If the number of submissions is limited, submitters optimize along other dimensions, such as length, complexity, and density. If word count is limited, compression and obfuscation advance.

Such constraint-evasion behavior is widely documented in behavioral economics as a rational response to regulation. The problem is that, in conventional Goodhart problems, discovering and implementing evasive strategies required human labor and cost. That labor itself functioned as friction, and it was often realistic for regulation or institutions to respond quickly enough for the constraint to remain effective.

This stance does not work against AI-enabled mass production, because the search for and implementation of evasive strategies themselves fall within the basic capabilities of LLMs. Instructions such as "compress this application into fifteen pages" or "reduce the number of claims while preserving the substantive scope of protection" are ordinary tasks for LLMs. They are not special alignment problems but basic capabilities. The speed and accuracy of evasive optimization against each indicator in a flow limitation regime increase as LLM capabilities improve. A quantitative evaluation of how much evasive optimization capability each generation of LLMs possesses against institutional constraints is an important empirical question, but it is beyond the scope of this paper and is left for separate treatment.

This capability did not suddenly appear. It is the cumulative result of three technical breakthroughs. The Transformer architecture introduced by Vaswani et al. (2017) made it possible to process entire input sequences in parallel through self-attention and solved the long-range dependency problems of recurrent neural networks. This made the training of large-scale language models practically feasible for the first time.

Next, OpenAI's GPT series (Radford et al., 2018; Radford et al., 2019; Brown et al., 2020) demonstrated that capabilities improve predictably as parameter counts and data volume increase. The emergence of generalization to tasks not explicitly trained, namely few-shot learning, in GPT-3 with 175 billion parameters showed that language models could function not merely as pattern-reproduction devices but as general-purpose instruction-following systems.

Finally, reinforcement learning from human feedback, RLHF (Ouyang et al., 2022), and its consumer-facing implementation in ChatGPT, released in 2022, raised the ability to follow arbitrary natural-language instructions to a practical level. Through these three stages, constraint-satisfying

rewriting became not a special capability of LLMs, but a direct expression of their basic language-processing capacity.

In other words, once the indicators of a flow limitation regime are made public, evasive optimization against those indicators becomes one of the standard uses of LLMs. Because this evasive capability continues to improve in both accuracy and speed as LLM capabilities improve, flow limitation based on a single indicator becomes more fragile over time.

The cognitive cost score in Flow-by-Flow is constructed as a product of multiple features as a design response to this fragility. If the number of claims is reduced, the description per claim becomes more complex and word count increases. If word count is compressed, increased density is reflected in the complexity coefficient. Because each feature is rooted in the nature of the content itself, compressing one dimension transfers burden to another. Whereas evasion of a single indicator is a one-dimensional optimization problem, simultaneous evasion of mutually constraining multiple indicators becomes a multidimensional constraint-satisfaction problem, increasing the computational cost of evasion. Complete elimination of evasion is impossible in principle. It is possible, however, to design an institution in which the attempt to evade increases the cost of evasion.

### 3.5. Fixing the Institutional Capacity Cap

The cognitive cost score and the institutional capacity cap satisfy Invariants 1 and 2. For applications within the cap, no further mechanism is needed. For applications that exceed the cap, however, a mechanism satisfying all four invariants is required. Section 5.5 presents one reference implementation.

The third component of Flow-by-Flow is to institutionally fix the upper bound of processing capacity on the receiving institution's side. The discussion so far has concerned incentive design for submitters. Submitter-side control alone, however, does not specify what should happen when submission volume exceeds  $C_{\max}$ . The institutional capacity cap makes the physical constraint of the receiving institution explicit as an institutional rule. The upper bound of total cognitive load that can be examined annually is calculated as follows:

annual processable load = number of human examiners with official examination authority  $\times$   $C_{\max}$   $\times$  annual working hours

Applications that exceed this cap are not rejected. Instead, they are routed into an exceedance pathway described in Section 5.5. Applications within the annual processing capacity are handled through automatic measurement of the cognitive cost score alone. If a single submitter uses AI to submit thousands of applications at once, all of them enter the exceedance pathway. The right of access is fully preserved, but processing speed is controlled so that it remains within human processing capacity. An important secondary effect of fixing the institutional capacity cap is that it naturally creates an incentive for submitters to make their applications simple and understandable. Once the finitude of processing slots becomes explicit, submitters voluntarily pursue concision and clarity in order to achieve the best result within the limited capacity. This is not a restriction imposed by regulation. It is an incentive design arising naturally from the finitude of physical resources.

### 3.6. From AI-Use Disclosure Requirements to Process-Time Declaration Requirements

The automatic measurement of cognitive cost and the institutional capacity cap presented in this chapter control submissions on the basis of the formal features of the submitted output itself. For operational purposes in flow design, however, it is also important to obtain information about how the output was actually produced. This section proposes replacing existing AI-use disclosure requirements with a requirement to declare process time. The determination of whether the cap is exceeded is made solely by the cognitive cost score. Other information sources, including process-time declarations, are not used for that determination.

Current governance frameworks, including the EU AI Act, are moving toward requiring self-disclosure of whether AI was involved in producing a submission (European Parliament, 2024). This approach confronts the problem discussed in Section 2.2 of the previous paper: there is no

independent means of verifying the disclosure. Post hoc detection of AI-generated text is unreliable, as Liang et al. (2023) show, and detection rates fall sharply through simple prompt manipulation. The detector released by OpenAI in 2023 achieved an accuracy of only 26% and was discontinued. In other words, disclosure requirements are not accompanied by a means of detecting falsehood.

Unlike real-name registration systems, in which identity-disclosure requests to providers can function as post hoc verification hooks, AI-use disclosure has no post hoc verification hook. Process-time declaration replaces the yes-or-no question of whether AI was used with a continuous quantity: how many real-time hours were required for each process. For an academic paper, this may include the time required for idea generation, literature review, experimentation, data analysis, writing, and revision. For a patent application, it may include the time required for prior-art search, organizing the invention, drafting the specification, and preparing the claims.

Compared with proposition-type disclosure, the routes for verifying process-time declarations are far more numerous. If a paper declares three hours of literature review but lists 200 references, this implies that each prior work was read and evaluated in 54 seconds. If a specification declares ten hours of writing time but contains eighty pages, this implies that eight pages were written per hour. These consistency checks require no external identity-disclosure request, because the submitted output itself provides the verification material.

With respect to incentives for false declarations, this proposal separates process time from the evaluation of exceedance. The only means of determination is the cognitive cost score. Process time is used as an independent channel for authenticity verification and field-classification calibration. As a result, neither an incentive to exaggerate time nor an incentive to underreport time is attached to the declaration itself. What remains is the incentive to report a value close to the true value in order to avoid inconsistencies being identified by consistency checks.

The proposal functions even without assuming that declarations are accurate. Most current AI-use disclosures are single-value yes-or-no declarations. They reduce the issue to whether the statement is true or false, but because this cannot be verified, the information is entirely worthless. Process-time declarations, by contrast, accumulate continuous values by field and by process. Over time, the distributions reveal the natural median and interquartile range of process times in each field. Systematic bias in inflation or deflation can also be estimated from correlation analysis with the formal features of submitted outputs. For example, if one group reports markedly shorter writing times than another group for papers of the same length, the direction and degree of bias in that group's declarations can be inferred.

A secondary effect of process-time declarations is the continuous quantification of the B variable discussed in Section 4. If the median proportion of process time belonging to physical-space activities is calculated by field, B can be redefined from a qualitative binary classification into a continuous value. In geology, for example, one can calculate what proportion of the workflow is spent on fieldwork and sample analysis and what proportion is spent on writing. Similar distributions can be obtained for ecology, pharmaceutical development, and materials engineering. Declarations that deviate from these distributions, such as a geology paper reporting zero fieldwork time, become visible as statistical deviations from field norms.

The qualitative classification in Appendix A of the previous work thereby acquires a pathway for ex post calibration through operational data. The advantage of this proposal is that it begins from the recognition that existing AI-use disclosure requirements do not function, but replaces them not by abolition alone, but by an alternative disclosure framework that can generate useful information.

### 3.7. *Two Layers of Flow-by-Flow*

The design presented in this paper is integrated into the following two-layered structure. The first layer is automatic and immediate: automatic measurement of cognitive cost. Formal features of the submitted output, such as the number of claims, word count, and number of citations, are measured. After applying existing classification systems and complexity coefficients, the cognitive

cost score is calculated automatically. This layer can be fully automated and contains no content judgment.

The second layer is human: only outputs that fall within the institutional capacity cap are examined by human examiners. The volume of outputs reaching this layer is guaranteed to remain within  $C_{max}$ . The operational estimation of  $C_{max}$  required to calculate annual processable load does not directly measure the cognitive capacity of individual examiners. Rather, it is set conservatively from the relationship between processing volume and error rates in each institution's existing operational records, with a margin placed before the point at which error rates begin to rise, as discussed in the supplementary information to the previous work.

The essence of this operation is that the first layer functions as an automatic deceleration device, while the humans in the second layer process only a workload that remains within  $C_{max}$ . Because humans receive only a cognitively manageable volume of outputs, degradation of  $C_{eff}$  is suppressed.

The important point is that the first layer does not concern itself at all with the quality of content. It therefore neither detects nor judges the major problems of AI content judgment, such as hallucinations, factual errors, and logical leaps. Those judgments are entrusted to humans in the second layer, after the flow rate has been controlled within  $C_{max}$ . Flow-by-Flow does not control the quality of AI outputs. It systematically arranges the conditions under which humans can control quality.

### 3.8. Applicability and Threshold Structure

The Flow-by-Flow design can be applied across domains wherever formal features exist from which a cognitive cost score can be calculated. Examples include the following.

Patent examination: normalized claim-count multiplier  $\times$  normalized word-count multiplier  $\times$  normalized citation-count multiplier.

Academic peer review: normalized word-count multiplier  $\times$  normalized reference-count multiplier  $\times$  normalized figure-and-table multiplier  $\times$  normalized supplementary-material multiplier.

Judicial proceedings: normalized complaint-length multiplier  $\times$  normalized evidence-count multiplier  $\times$  normalized relevant-law multiplier.

Pharmaceutical review: normalized application-document-length multiplier  $\times$  normalized clinical-trial-count multiplier  $\times$  normalized indication-count multiplier.

In all cases, what is measured is only quantitative features. The truth or falsity of content is never referenced.

The flow control produced by the cognitive cost score has an intrinsic threshold structure. If the number of claims, word count, and number of citations are each twice the baseline value, the cognitive cost score becomes eight times the baseline. Because the score is a product, increases across multiple dimensions compound nonlinearly. An application within median values remains within the ordinary processing path, whereas an application with features at twice the median in all three dimensions enters an exceedance pathway unless sufficient institutional capacity remains. The specific form of that pathway is not fixed by the score itself; it must be supplied by a mechanism satisfying the four invariants.

### 3.9. Institutional Design Note

This chapter can now be summarized. What fundamentally distinguishes the Flow-by-Flow design from conventional AI governance is that it removes dependence on AI judgment from the foundation of institutional design.

From this perspective, the previous work proposed task rotation between AI-use tasks and non-AI tasks. At first glance, this may appear to be a weakly grounded proposal. Yet institutional designs of the same kind are already widely established in high-loss domains. In aviation, pilots' consecutive flight hours are strictly regulated, because degradation of judgment by a fatigued pilot directly endangers the lives of hundreds of passengers. Long-distance bus drivers are also subject to limits on

consecutive driving time, and serious accidents have repeatedly occurred when such regulations were not observed. Consecutive operating hours for surgeons are also limited.

In all of these cases, the magnitude of potential loss justifies limits on the working time of the processor. The same activity does not always require regulation. Riding a bicycle for a long time is not regulated. In sports, elementary school students playing football all day is not legally restricted, but if a professional-level match were continued for eight hours, players would be injured. Regulation is justified not by the type of activity alone, but by the combination of activity intensity and magnitude of loss.

The improvement of AI capabilities creates the need to apply this logic to intellectual labor in high-loss domains. As model capabilities improve, the cognitive intensity required of supervisors increases. Supervising the output of a model with 99% accuracy imposes a higher cognitive burden than supervising the output of a model with 30% accuracy. This claim has already been experimentally established in automation research. Parasuraman, Molloy, and Singh (1993) demonstrated that complacency arises within twenty minutes when monitoring a highly reliable automated system. Oakley et al. (2003) compared seven levels of reliability and showed that higher reliability substantially worsens human monitoring efficiency. Singh et al. (2009) similarly reported that high static reliability induces complacency and that monitoring performance is far better under lower reliability.

## 4. Weakest-Link Lockout and the B Variable

### 4.1. Lockout in General-Purpose Models

General-purpose AI models, by design, possess capabilities across multiple domains at the same time. The ability to write code, the ability to draft legal documents, and the ability to summarize scientific papers are different manifestations of the same model, and it is technically difficult to separate them individually.

Outputs from general-purpose AI models include both low-loss and high-loss domains. This inseparability creates the following problem. The condition for a general-purpose model to be safely released is that final human supervisability must be maintained in all domains of application. In other words, in every domain, the product of output rate and cognitive load in that domain must remain within that domain's supervisory capacity.

This paper calls the situation in which the single most vulnerable domain determines the release condition of the entire model Weakest-Link Lockout.

For a model specialized in a single domain, such as a medical-only or legal-only model, only  $V \times L$  in that domain is at issue. In a general-purpose model, however, capability improvements occur across domains, and a threshold may be exceeded in a domain that the developer did not intend. For example, training aimed at improving coding ability may, as a by-product, increase the model's capability to discover vulnerabilities in the cybersecurity domain. This can arise from the inseparability of capabilities in general-purpose models.

If this inseparability is addressed through content filtering, selective restriction becomes unsustainable as the number of lockout domains increases. This is because knowledge overlaps densely across the domains subject to lockout, and restrictions produce side effects that impair legitimate uses. Flow control, namely controlling how much output is permitted, therefore becomes the residual option against content judgment, namely deciding what may be output.

### 4.2. The B Variable: Domain Differences Produced by Physical-Space Constraints

The previous work (Naito, 2026) showed that the rate of expansion in AI output rate  $V$  is not uniform across domains, but differs systematically according to the presence or absence of physical-space constraints. This paper reorganizes that factor as the B variable.

B is a qualitative indicator of the extent to which the standard research or business workflow in a domain requires activities in physical space, such as experimentation, fieldwork, physical

measurement, or sample collection. Generative AI can drive process costs within information space close to zero, but it cannot replace processes in physical space. Therefore, the lower  $B$  is, the greater the acceleration of output rate enabled by AI becomes. The higher  $B$  is, the more that acceleration is constrained by physical-space bottlenecks.

Appendix A of the previous work used approximately 25,000 preprint records from Preprints.org and showed that, after 2023, the divergence in output rates widened between domains with  $B \approx 0$ , such as mathematics, and high- $B$  domains, such as geology and entomology. This divergence could not be explained by the growth of the platform as a whole, and was also independent of the increase in researchers or funding associated with the AI boom. Similar acceleration was observed even in academic fields unrelated to AI investment.

The  $B$  variable provides a framework for predicting the order in which Weakest-Link Lockout is triggered.

## 5. Toy Model Analysis: From Judgment Is All You Need to Counting Is All You Need

### 5.1. Two Stages of Control and the Time Axis

The overview of the Flow-by-Flow design presented in Chapter 3 stated that the cognitive cost score is constructed as a product of multiple features. This chapter uses a simplified model of applicant behavior, hereafter referred to as the toy model, to show why such compositionality produces a stronger institutional design than flow limitation based on a single metric. It also discusses how the need for such compositionality develops over time and differs across domains.

Control of AI output in high-loss domains can be divided into two stages from the perspective of implementation sequence.

The first stage is a direct restriction on submitting actors and submission counts. Specifically, it includes identifying submitting actors through know-your-customer procedures, KYC, and setting an upper limit on the number of submissions per unit time. The purpose of this stage is to physically prevent a specific actor from submitting a large number of applications within a short period.

A distinctive feature of this stage is that KYC and submission-count limits function inseparably as a pair. Submission-count limits without KYC are easily bypassed through the replication of anonymous accounts and therefore have no meaningful effect.

As of April 2026, major AI providers differ in their implementation of flow control. Among general-purpose models with comparable capabilities, some providers implement both KYC and user-count limits, others implement KYC but impose no output restriction, and still others implement neither. This divergence reflects the tension between the theoretical necessity of flow control derived in this paper and commercial incentives under market competition.

The implementation status of existing institutions in high-loss domains is also uneven. Institutions such as patent applications and litigation filings impose real-name requirements on applicants, but do not impose limits on the number of submissions themselves. They are KYC systems without output restrictions. Academic peer-review systems require authors to declare their real names, but lack an independent means of verifying the truth of such declarations and impose no submission-count limit. They are systems without effective KYC and without output restriction. The toy model presented in this chapter also functions as a framework for organizing this diversity of implementation.

We first analyze the limits of a complete first-stage implementation, **specifically** a KYC-plus-submission-count-limit system, and then show why institutions must move toward the second stage, namely control based on the cognitive cost score. The transition from the first stage to the second stage does not occur simultaneously across all domains. The speed of transition differs across domains, and this difference is linked to the temporal progress of LLM capability improvements.

Why is the second stage necessary? Once the number of submissions is limited, the applicant's optimization pressure shifts from increasing the number of submissions to increasing the amount of

information and complexity per submission. Longer texts, more claims, inflated citations, and obfuscation are all rational responses under submission-count limits.

Of course, LLMs do not automatically possess, at the moment a submission-count limit is introduced, an intrinsic ability to optimize by “packing more into each item.” This optimization capability has been acquired gradually through a series of breakthroughs: Transformers, scaling laws, and RLHF. As of 2026, models are capable of going beyond simple compression and executing multi-objective optimization such as “minimize the cognitive cost score while maximizing the scope of protection.”

Evasive optimization per item does not necessarily activate immediately when the first stage, submission-count limitation, is introduced. Its activation requires both that LLM capability has reached the relevant level and that the applicant community discovers and disseminates the corresponding use case.

The timing of the transition from the first stage to the second stage depends on the domain, and this dependence is consistent with the B variable. In low-B domains, such as patents, academic papers, legal documents, and administrative applications, compression, expansion, and restructuring of content are standard tasks for LLMs. Evasive optimization that increases information density per item can therefore be activated immediately after submission-count limits are introduced. In high-B domains, such as clinical-trial applications, pharmaceutical approval, and in-person medical records, each item contains physical elements that AI cannot generate. Natural barriers against evasive optimization therefore remain, and the transition to the second stage is delayed.

## 5.2. A Toy Model of Applicant Behavior

This section formalizes applicant behavior under flow limitation as a simplified model. The purpose of the model is not to cover real-world diversity exhaustively, but to provide an analytical framework for clarifying why composite metrics are structurally stronger than single metrics.

Consider a knowledge infrastructure, such as a patent office, peer-review system, court, or pharmaceutical review agency, to which applicants submit applications. The knowledge infrastructure has an upper bound on processing capacity and imposes some form of flow limitation.

Applicants seek, under the given flow limitation, to obtain as many desirable outcomes as possible, especially successful passage of their applications. They decide two things: how many applications to submit within a given period, and how to write each application. More specifically, they choose the values of formal features such as word count, number of claims, and number of citations.

For applicants, an application generates benefits, but creating it also entails costs. A successful application yields a certain return. Applicants compare the return obtained from successful passage with the creation cost and institutional cost, such as exceedance friction, and act so as to maximize expected payoff.

Under this setting, we distinguish three types of applicants.

The first type is the AI non-user. This type creates applications by conventional methods and does not use LLMs even as an auxiliary tool. The time and effort required to create a single application are large. The feasible number of submissions per unit period is naturally constrained by the applicant’s own physical and cognitive limits. The features of the resulting application reflect what the applicant genuinely wishes to claim, and the incentive to strategically manipulate those features is weak.

The second type is the AI user. With LLM assistance, creation costs decline substantially. The important point is that this type contains a continuous spectrum. At one end are users who employ LLMs merely as a means of improving drafting efficiency and determine the content of the application on the basis of what they themselves wish to claim. At the other end are users who aim to maximize submission counts and generate the content of each application in a form that is likely to satisfy passage conditions. Real-world use is distributed between these two poles.

Along this spectrum, the effect of applicant behavior on the institution changes. What the toy model in this chapter focuses on is not the precise position of each user on this spectrum, but the following fact: when an institution imposes a particular flow limitation, such as a submission-count cap, applicants using LLMs will, regardless of malicious intent, act so as to maximize their own payoff under that constraint. If the number of submissions is limited, LLM capability is naturally redirected toward increasing information density per item.

This optimization does not necessarily require intentional institutional evasion by individual users. Merely by making use of LLMs, collective pressure arises that exceeds the boundary conditions assumed by the institution.

It is therefore unnecessary to distinguish between good-faith AI use and malicious AI use. The question is what kind of feature distribution the entire population of LLM-using applicants collectively produces under the constraints imposed by a given institution.

### 5.3. Limits of a First-Stage-Only Institution

This section analyzes what kind of feature distribution is produced by the collective behavior of AI users under an institution that implements KYC and submission-count limits, namely a first-stage-only institution.

The first-stage institutional design is as follows. Each applicant is identified through identity verification, and the number of applications that can be submitted per unit period is constrained by a fixed cap. Once the number of submissions is constrained, the applicant's only remaining optimization variable is the information density of each item.

Because the expected return from a single application depends on comprehensiveness, complexity, and scope of claims, maximizing information density becomes a rational choice. For AI non-users, increasing density means an increase in creation cost, and there is a natural upper bound. For AI users, by contrast, LLMs can draft additional claims, expand pages, and add citations at low cost, making it possible to push information density to a much higher level.

From the institution's perspective, the number of submissions is controlled, and the flow limitation therefore appears to be functioning formally. However, if  $L$  increases while  $V$  is fixed,  $V \times L$  can exceed  $C_{\max}$ . This exceedance arises as the aggregate result of rational behavior by individual applicants who are trying to submit the best possible application under the rules. If LLM capability continues to improve, the information density that can be generated per item under the same constraint will increase over time, and this trajectory is irreversible.

### 5.4. A Second-Stage Institution Based on Composite Metrics

The second-stage institutional design adds to the first stage, KYC and submission-count limits, the measurement of the cognitive cost score for each application and an exceedance pathway whose cost increases with the amount by which the threshold is exceeded.

For concreteness, consider patents as an example. Let the features be three:  $X$ , normalized text length;  $Y$ , normalized number of claims; and  $Z$ , normalized number of citations. Each feature is made dimensionless by dividing it by the field median. Let the score be  $S = X \times Y \times Z$ .

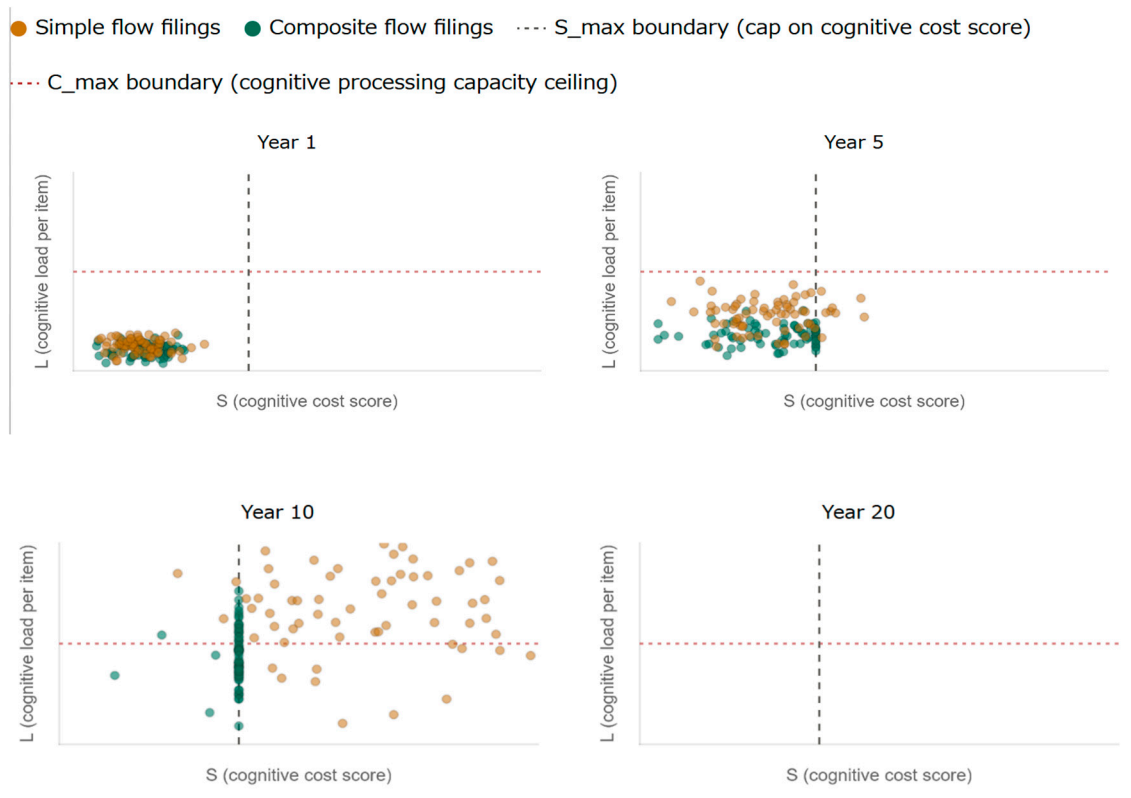
In the second stage, an increase in  $S$  is directly linked to an increase in institutional friction, making one-directional optimization impossible. If  $X$  is increased,  $Y$  or  $Z$  must be reduced in order to satisfy the  $S$  constraint. Optimization along each dimension interferes with optimization along the others.

The robustness of composite metrics does not lie in differences in computational cost alone. It lies in the fact that the features are inseparably connected to the content. If claims are reduced, the scope of protection narrows. If pages are compressed, descriptions are omitted. The substantive trade-offs produced by multidimensional constraints cannot be eliminated by improvements in LLM capability.

The institution can increase the difficulty of evasion simply by adding the number of features  $N$ . In a score constructed as the product of  $N$  features, compressing a single dimension leaves the burden

in the remaining  $N - 1$  dimensions. Therefore, all dimensions must be optimized simultaneously. Adding  $N$  is a low-cost operation for the institution, while for the evading side, the complexity of the problem increases with each additional optimization dimension.

This robustness, however, is relative rather than absolute. As discussed in the next section, if LLM capability reaches a sufficiently high level, multidimensional constraint-satisfaction problems may also become solvable. But once one accepts the starting point that completely eliminating evasion is impossible in principle, it remains possible to increase the computational and economic cost of evasion. This is the core design principle of Flow-by-Flow.



**Figure 1.** Numerical visualization of the toy model in Sections 5.3-5.4. X axis = cognitive cost score  $S$ , the product of formal features such as page count. Y axis = actual cognitive load  $L$  per filing. Under simple flow,  $S$  is uncapped, so both  $S$  and  $L$  diffuse rapidly over time. Under composite flow,  $S$  is capped at  $S_{\max}$ . The only remaining path for raising  $L$  without raising  $S$  is sophisticated optimization against dimensions outside the scored features. The suppression effect in the composite case decays over time as AI optimization capability improves. The decay period, set to 20 years in this illustration and sampled from 5 to 50 years in the Monte Carlo analysis, reflects the inherent uncertainty in predicting when AI systems will become capable of fully circumventing multidimensional institutional constraints.

LLM capabilities continue to improve, and their ability to handle complex constraint-satisfaction problems is also advancing on a yearly basis. This capability improvement, however, does not invalidate composite metrics. As discussed in Section 5.4, so long as the features are inseparably connected to the content, even highly capable LLMs cannot eliminate substantive trade-offs such as the narrowing of protection scope when the number of claims is reduced.

A dynamic equilibrium emerges in which LLM capability improvement and institutional response continuously contend with each other. The design does not aim for complete victory by either side, but for sustainable contestation. This is the decisive difference from single-metric restrictions. Single metrics are unilaterally invalidated by improvements in LLM capability. Composite metrics, by contrast, contain within themselves a means of institutional response.

### 5.5. Incorporating the Physical Gate: A Reference Implementation

The second-stage institution based on composite metrics is more robust than single-metric restrictions, but it leaves a serious side effect. As the number of dimensions  $N$  increases, the probability also increases that legitimate applications requiring special fields or inherently complex inventions will happen to be caught by the cap. To mitigate this side effect, applications that exceed the upper bound of the cognitive cost score should not be uniformly rejected. Instead, a route that legitimate applications can pass through must be established.

Any mechanism serving this role must satisfy the four design invariants derived in Section 3.1. The physical waiting path presented here is one reference implementation that satisfies all four. Its contribution is as a proof of existence, demonstrating that at least one such mechanism can be constructed.

**Mechanism.** An application whose cognitive cost score exceeds the threshold is required to undergo physical waiting at a designated KYC-enabled office as a condition for receiving a submission passcode. The applicant visits the office with an identity document such as a passport, waits for a specified baseline waiting time, and then receives a 16-digit submission passcode. By entering this code, the application is accepted even if it exceeds the threshold.

Physical waiting time is proportional to the multiplier by which the cognitive cost score exceeds the threshold. If the baseline waiting time is  $T$ , for example two hours, and the threshold multiplier of the score is  $k$ , then: physical waiting time =  $T \times k$ . If  $k$  is less than or equal to 1, there is no waiting. If  $k = 8$ , the waiting time is sixteen hours. Waiting time can be dynamically adjusted as an institutional parameter without requiring algorithmic updates or content judgment.

Mapping to the four invariants. The physical waiting path satisfies Invariant 1 because it involves no evaluation of content. It satisfies Invariant 2 because it does not consume examiner time: only service-desk identity verification is required. It satisfies Invariant 3 because identity verification by official documents is conducted both at the start and at the end of waiting, binding each waiting event to a verified identity. It satisfies Invariant 4 because each application requires a separate waiting event and a separate passcode, preventing batch clearance.

**Asymmetric screening.** The mechanism belongs to the lineage of Spence signaling theory (1973) and Salop and Salop self-selection theory (1976), in which separating equilibria are achieved by exploiting cost differences between legitimate actors and mimicking actors. For legitimate human applicants who genuinely require complex applications, the waiting cost is a one-time inconvenience. For AI-enabled mass producers, the cost scales at least linearly with the number of applications and may increase superlinearly due to coordination, movement, and employment costs. Unlike investment in API usage, investment in physical waiting cannot preserve the scaling advantage of digital replication.

The essence of the constraint is that it transforms the marginal-cost structure of information space into the marginal-cost structure of physical space. A single user can generate 100,000 outputs through API calls at near-zero marginal cost. Under the physical waiting path, 100,000 applications require 100,000 instances of physical waiting.

**Reasonable accommodation.** The purpose of the physical waiting path is not to require physical travel for its own sake, but to ensure that each threshold-exceeding application consumes a finite amount of identity-verified human time. Accordingly, remote identity verification, regional offices, reservation systems, and assisted proxy support may be permitted, provided that the one-application-one-waiting-instance structure is preserved. The minimum invariant is that the marginal waiting cost must remain attached to each application. The detailed implementation of accommodations should be left to each jurisdiction.

**Practical difficulties.** It must be stated plainly that the physical waiting path faces serious practical difficulties. First, accessibility: requiring physical presence imposes disproportionate burdens on applicants with disabilities, applicants in remote regions, and applicants in developing countries. The reasonable accommodations above mitigate but do not eliminate this problem. Second, legal compatibility: existing legal frameworks in many jurisdictions guarantee immediate intake of applications, and introducing mandatory waiting periods may require legislative or treaty revision.

Third, international coordination: if one jurisdiction introduces the physical waiting path while another maintains full immediate intake, applicant migration across jurisdictions may undermine the design.

For this reason, the contribution of this paper is not the physical waiting path itself, but the four design invariants derived in Section 3.1. The remote identity-verified time-lock pass discussed in Appendix C-2 is one alternative candidate that trades some robustness against proxy labor for substantially improved accessibility. Any mechanism that satisfies all four invariants would serve the same function. If a strictly superior mechanism is found, it should replace the physical waiting path. Additional discussion of the relationship between the physical waiting path and existing screening mechanisms, including computational proof-of-work and alternative penalty designs, is provided in Appendix D.

### 5.6. Monte Carlo Analysis: Implementation and Results

The simulation results of the toy model in this paper do not predict the specific number of years until a particular domain collapses. The concrete timing of collapse depends strongly on the  $B$  variable discussed in the previous work, which captures domain differences arising from the presence or absence of physical-space constraints, and on the adoption rate of AI. It therefore differs greatly across domains.

The difference between domains in which collapse has already been observed, such as peer-review systems in AI and vulnerability response in cybersecurity, and domains in which no sign of collapse is currently visible, such as peer review in entomology, patent applications in emerging economies, and clinical-trial approval, reflects empirical differences across domains that cannot be represented by any single simulator parameter.

Accordingly, what the Monte Carlo analysis in this section shows is the robustness of relative ordering. That is, the simulator cannot predict when a given domain will collapse, but it can test whether the ordering of the amount of time gained by adopting one of three strategies in a given domain is maintained across a wide range of parameters.

The sampling ranges were as follows. The initial  $C_{\max}$  utilization rate was 10% to 50%. The annual increase rate of  $V$  was 10% to 27%, corresponding to the range of empirically observed  $V$  increases in Kusumegi et al. (2025) and the previous work. The annual increase rate of  $L$  was 10% to 40%. The increase in checking capacity under the supervision-enhancing system was 10% to 20% for strategy (a), 5% to 13% for strategy (b), and 2% to 8% for strategy (c). The initial suppression rate of composite flow was 50% to 95%. The number of years over which composite suppression decays to zero was set between 5 and 50 years.

The HITL increase rates for (b) and (c) were set lower than that for (a) to reflect the assumption that the cost of operating flow control consumes part of the budget that would otherwise be used to increase human supervision. The initial suppression rate of composite flow, 50% to 95%, and the suppression decay period, 5 to 50 years, have no empirical anchor because the Flow-by-Flow design has not yet been implemented.

The lower bound of the decay period corresponds to the arrival time of AI capability sufficient to optimize multidimensional institutional constraints faster than human institutional response. Such a capability presupposes a world in which not only the Flow-by-Flow design, but all currently human-designed institutions, including law, taxation, and safety standards, are arbitrated by general-purpose AI models. Whether this arrives in five years or fifty years is currently unpredictable, and this uncertainty itself is one of the major AI risks.

The initial  $C_{\max}$  utilization rate was set at 10% to 50%. This is an assumption extremely unfavorable to our argument, because it implies that 50% to 90% of supervisory institutions' processing capacity is unused, which is unlikely to hold in reality.

A low growth rate of  $L$  is not a counterexample to this framework. If per-item cognitive load increases by only about 2% per year, for example, that means that AI has only a limited ability to produce outputs that are more complex, more human-like, and better at evading institutional

constraints. In that case, it is unnecessary to increase the number of dimensions in the cognitive cost score frequently. False positives are also few, and activation of physical proof-of-work remains exceptional. Thus, low L growth does not show that flow control is unnecessary. It indicates an optimistic regime in which lightweight flow control is sufficient.

Conversely, if L grows at an extremely high rate, for example more than 40% per year, our design has clear limits. Adding features makes evasion more difficult, but it also increases false positives among legitimate applicants and increases the number of cases requiring physical proof-of-work. In an ultra-high-growth scenario, it may therefore become necessary to shift toward upstream controls such as coarser processing-capacity caps, intake quotas, lotteries, licensing systems, or temporary suspension of intake. This is not a counterexample to the framework, but rather indicates the need to switch institutional friction in stages according to the growth rate of AI capability.

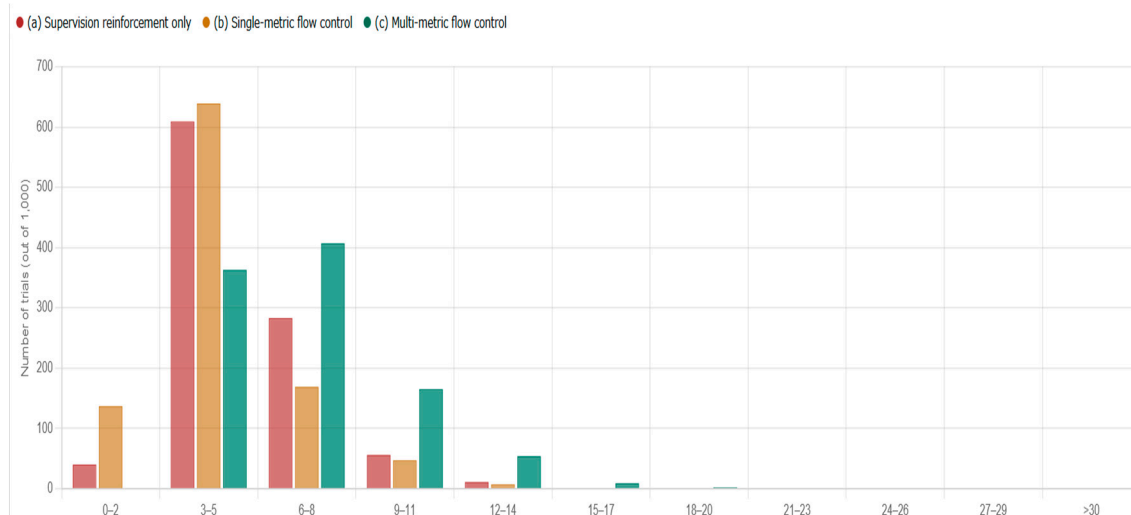
To avoid favorable parameter settings, this simulation intentionally gives extremely favorable initial conditions to the supervision-enhancing approach. Even under those conditions, the superiority of composite flow is maintained in 90.8% of trials. Under uncertainty, it is rational to adopt the strategy that endures the longest, and the simulation shows that composite flow is that strategy across a wide range of parameters.

Table 1 reports the share of trials in which each strategy performed best, based on 30-year simulations for each trial and the rate at which each strategy avoided exceedance of  $C_{max}$ .

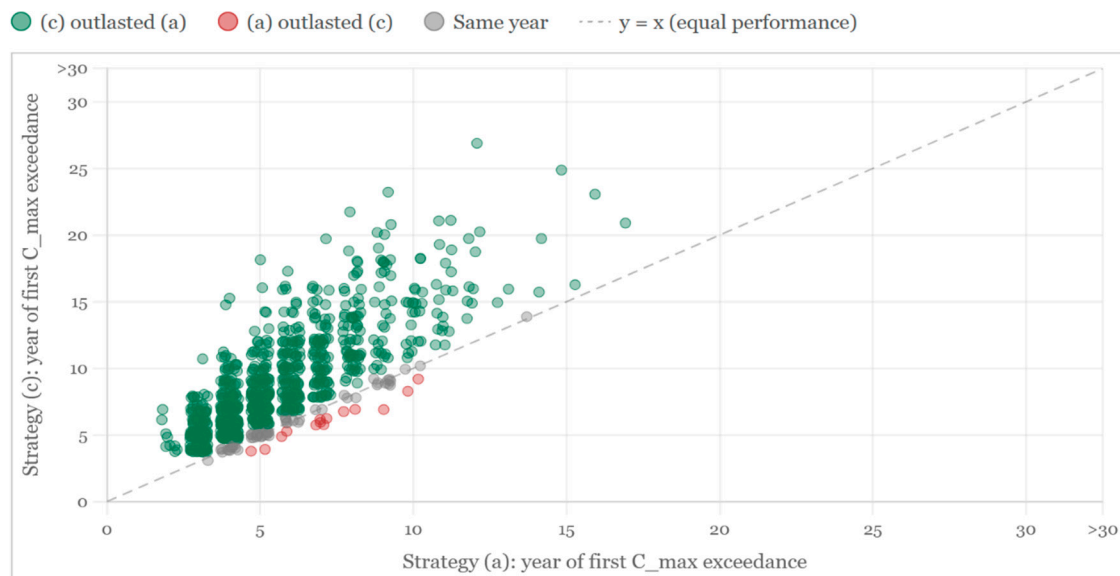
**Table 1.** Monte Carlo Analysis Results (N = 1,000).

Strategy	Share of trials in which the strategy performed best
(a) Supervision enhancement only	1.6%
(b) Supervision enhancement + simple flow	7.6%
(c) Supervision enhancement + composite flow	90.8%

For all strategies, supervision enhancement is implemented as a common baseline by linearly increasing  $C_{max}$ . The annual increase in  $C_{max}$  relative to the initial value was set at 10% to 20% for (a), 5% to 13% for (b), and 2% to 8% for (c). The smaller increases for (b) and (c) reflect the conservative assumption that the administrative cost of flow operation may crowd out resources allocated to supervision enhancement. Because this assumption disadvantages (c), the simulation may underestimate the superiority of (b) and (c).



**Figure 2.** For each of 1,000 parameter draws, the year in which  $V \times L$  first exceeds  $C_{\max}$  is recorded for each strategy. The rightmost bin,  $>30$ , indicates trials where no exceedance occurred within the 30-year simulation horizon. Supervision reinforcement alone, strategy (a), clusters in the early years. Single-metric flow control, strategy (b), delays exceedance moderately. Multi-metric flow control, strategy (c), pushes the majority of trials beyond the simulation horizon.



**Figure 3.** Each point represents a single trial under identical parameter draws. The horizontal axis is the year in which strategy (a), supervision reinforcement only, first exceeds  $C_{\max}$ . The vertical axis is the year in which strategy (c), multi-metric flow control, first exceeds  $C_{\max}$ . Coordinate (31, 31) denotes trials where neither strategy exceeded  $C_{\max}$  within the 30-year horizon. Points above the diagonal indicate trials where (c) outlasted (a). Points below the diagonal indicate the reverse.

These results show that the ordering in which the introduction of flow control buys more time than supervision enhancement alone is maintained across a wide range of parameters. The timing of collapse in individual domains must be estimated separately on the basis of the domain's  $B$  value and its AI adoption and use rate.

The simulation does not present the value 90.8% itself as a prediction. Its purpose is to examine how robustly the relative ordering among supervision enhancement alone, simple flow, and composite flow is maintained under a wide range of parameter uncertainty.

In AI governance risk assessment, once the collapse of human checking systems has been observed as empirical data, institutional response can only be retrospective. Ex ante analysis is therefore indispensable. To claim that the parameters of this simulation are unrealistic, one would have to assume that AI capabilities plateau and that human use of AI does not continue to expand. The continued tens-of-billions-of-dollars-scale investment by major AI companies and published GPU development roadmaps do not support that assumption.

The range for the  $V$  increase rate, 10% to 27%, is based on the output increase rate measured by Kusumegi et al. (2025) from more than two million preprints, 23.7% to 89.3% over approximately two years. The upper bound of 27% in this simulation is near the median of the observed values, and the lower bound of 10% is a conservative setting that assumes stagnation in AI diffusion. Considering that the observation period of Kusumegi et al. included the early stage of AI use, a period unfavorable to paper writing, and that  $V$  has continued to expand after 2025 (Naito, 2026), the actual increase rate of  $V$  may well exceed the upper bound used in this simulation.

The range for the  $L$  increase rate, 10% to 40%, is also conservative in light of multiple empirical studies showing that human discrimination accuracy for AI-generated content is converging toward chance level (Chen et al., 2024; Jakesch et al., 2023; Cooke et al., 2025).

## 6. Limitations

### 6.1. Components of the Cognitive Cost Score

The features used in this paper as examples of components of the cognitive cost score, such as the number of claims, word count, and number of citations, are only initial approximations. The resistance of each feature to evasion is not uniform. For example, the number of citations may be more vulnerable than word count to strategic reduction using AI.

The essence of the Flow-by-Flow design does not lie in the selection of any particular feature. It lies in the design principle of constructing the score as a product of multiple mutually constraining features. If individual features are hacked, the institution can respond by adding new features or recalibrating coefficients.

As discussed in Section 3.3, processing-time data accumulate through institutional operation. It is therefore empirically detectable which features deviate from actual cognitive cost. No design can be made completely impossible to optimize against. However, continuous calibration based on operational data may allow institutions to maintain a more favorable position than they currently hold in the ongoing contest between evasion and response.

### 6.2. Trade-off in the Number of Dimensions $N$

Increasing the number of dimensions  $N$  makes evasion more difficult, but it also increases the probability that legitimate applicants will exceed the threshold. However, the penalty for threshold exceedance need not be rejection. It should be a bounded, identity-bound per-application friction, and the physical waiting path discussed in Section 5.5 is only one reference implementation of such a mechanism.

The appropriate value of  $N$  is determined by the balance between evasion difficulty and the frequency of false positives. It should be adjusted empirically as operational data accumulate. The upper bound for each dimension should be calibrated by referring to the distribution of natural applications in operational data.

### 6.3. Operational Estimation of $C_{max}$ and $L$

The operational estimation of  $C_{max}$  is a limitation carried over from the previous work. This paper treats  $C_{max}$  as a constant with a biological upper bound and does not estimate its specific value by field. In implementation,  $C_{max}$  must be estimated from empirical data, such as the annual number of cases that can be processed per examiner in each field or the standard processing time per case. Such data are internally accumulated in many institutions, but cross-institutional comparison in publicly available form is not easy.

Similarly, although this paper decomposes  $L$  into three tasks, triage, judgment, and response, no established method exists for measuring each component independently. The ratio among these components is likely to vary by field and by individual, and calibration of this ratio is left to future empirical research.

The claim of this paper is a qualitative proposition: among the three components, at least triage and response do not decrease as AI accuracy improves. This proposition itself does not require precise measurement of each component. In concrete institutional design, however, estimating the ratio among the components will affect operational judgment. The purpose of this paper is not to provide field-specific measurements of  $C_{max}$  and  $L$ , but to present the principle of flow design that institutions should adopt.

### 6.4. Operational Design of Process-Time Declarations

Process-time declarations have the property that their interpretive accuracy improves as declaration data accumulate. In the initial stage, however, group-level distributions have not yet been established, and therefore criteria for interpreting individual declarations are also not established.

Immediately after introduction, process-time declarations cannot be used as material for the Flow-by-Flow design. This period must be treated as a data-accumulation phase. The extent to which process-time declaration data should be given weight within the Flow-by-Flow design is left to the operational design of each institution.

#### 6.5. *Asynchronous Implementation Across Jurisdictions*

The design in this paper takes a single institution as its unit of analysis. Coordination problems that arise when multiple institutions or jurisdictions introduce the design separately are outside the scope of this paper.

Patent systems, pharmaceutical approval systems, and academic peer-review systems all have mutual-recognition agreements or international institutional linkages. If one jurisdiction introduces Flow-by-Flow while another jurisdiction maintains the older principle of full immediate intake, user selection behavior may become biased across jurisdictions.

This coordination problem lies outside the theoretical claim of this paper. It should be addressed in a later phase of inter-institutional coordination, after empirical experience has accumulated from separate implementations in individual jurisdictions.

#### 6.6. *Scope of Applicability*

The scope of this paper is limited to AI risks in domains where formal verification is inherently impossible and where the losses caused by error are severe. Ordinary low-risk uses, domains in which external criteria can be established through statistical verification, and domains in which formal truth values can be determined, such as arithmetic calculation, formal logic, and database queries, are outside the scope of this paper.

In these domains, the problem of  $V \times L > C_{\max}$ , which is the premise of this design, either does not arise, or the total cost can be recovered even when errors or accidents occur, or content judgment may be established at low cost. Applying the Flow-by-Flow design to all domains, including low-risk domains, may create unnecessary friction.

The parameters used in the simulation in this paper are based on empirical data from academic papers (Kusumegi et al., 2025). For other high-loss domains, such as pharmaceutical applications, judicial proceedings, and clinical trials, equivalent empirical data on the growth rate of  $V \times L$  are insufficient. Applying this model to those domains requires domain-specific parameter estimation. Estimating the growth rate of  $V \times L$  and the timing of breakdown by domain remains a task for future empirical research.

The simulation results assume that  $V \times L$  expands compoundly over time. This assumption holds in domains where the generating actors are not institutionally bounded. In domains such as patent applications, academic paper submissions, preprints, legal complaints on the plaintiff side, deepfake generation, and cyberattacks, there is no institutional upper bound on the number of generating actors, and each actor has an economic incentive to increase the number of outputs through AI. For this reason,  $V \times L$  continues to expand.

By contrast, in domains where the generating actors are limited to a finite set of members within an institution, the expansion of  $V \times L$  is linear or limited. In areas such as AI support for police investigations, AI diagnostic support for physicians, and AI assistance for judges, the number of users is institutionally fixed, and the number of uses per person is constrained by the nature of the work, such as the number of cases, patients, or lawsuits. In these domains,  $V \times L$  saturates once all users have adopted AI. Because  $V \times L$  approaches a finite value, reinforcement of AI-checking systems alone may be able to catch up, and the need for flow control is relatively lower.

In real institutional design, it is therefore necessary first to determine which category the target domain belongs to. If anyone can enter as a generating actor, flow control is necessary. If the generating actors are limited to a finite set of humans within an institution, reinforcement of AI-checking systems may be sufficient.

However, even within the same domain, the generating side and the verification side may be separated. In civil litigation, judges belong to the verification side and are finite, while plaintiffs belong to the generating side and anyone can submit a complaint. Compound expansion of  $V \times L$  therefore occurs on the generating side. Increasing the number of judges alone cannot address this problem; flow control for litigation is required.

The theoretical endpoint of the Flow-by-Flow design is a state in which all countable dimensions have become saturated through evasive optimization, leaving differences only in non-countable dimensions, such as the texture of content or the depth of argument. In the case of citations in paper submissions, for example, this sequence may proceed from a limitation on  $V$ , to a layer-one limitation on the number of citations within  $L$ , to  $L_2$ , a limitation on the number of fields from which cited works originate, to  $L_3$ , a limitation on the cumulative word count of cited works, to  $L_4$ , a limitation on the number of papers cited by the cited papers themselves.

In such a state, detecting differences inevitably requires content judgment, and the system returns to the problem that this paper has sought to bypass. The Flow-by-Flow design does not eliminate this endpoint. It delays arrival at that endpoint.

### 6.7. Does $L$ Really Continue to Increase?

The theoretical claim and toy model in this paper assume that both  $V$  and  $L$  continue to increase. The continued increase of  $V$  is increasingly supported by empirical evidence in Kusumegi et al. (2025) and the previous work (Naito, 2026). By contrast, the continued increase of  $L$  currently contains a stronger element of prediction than of direct empirical demonstration.

If one focuses on the triage component of  $L$ , namely the task of determining whether an output is generated by AI or by humans, multiple empirical studies have accumulated showing that discrimination difficulty increases over time. Jakesch, Hancock, and Naaman (2023) show in PNAS that discrimination accuracy does not improve beyond chance level even when monetary incentives and training are provided. Chein, Martinez, and Barone (2024), in a review of prior studies across AI-generated face images, video, artworks, poetry, and text, report that although human evaluators could distinguish outputs from early generative AI models, their performance on outputs from newer generations of models often remains statistically indistinguishable from random guessing.

These studies indicate that improvements in generative AI capability are being empirically observed as increases in discrimination difficulty. Because triage cost is monotonically linked to discrimination difficulty, the closer discrimination comes to chance level, the more cognitive resources triage requires. This empirical observation supports the validity of the  $L$  growth-rate settings in the toy model of this paper. The lower bound of 10% per year is, if anything, conservative in light of the rate at which empirical studies show discrimination difficulty to be increasing. Given that several modalities reach random-guessing levels within only a few years, the actual growth rate of  $L$  may be closer to the upper end of the range used in this paper.

There are, in principle, two conditions under which the claim of this paper would cease to hold. The first is that AI capabilities stop improving and the increase in output-discrimination difficulty also stops. The second is that humans reduce their use of AI, restoring the proportion of human-only outputs.

If either or both of these conditions become reality, the theoretical claim and simulation in this paper would substantially lose their contribution. What matters here is the empirical meaning of these two conditions. Continued improvement in AI capability is a fundamental premise of the ongoing businesses of global AI companies such as OpenAI, Anthropic, and Google DeepMind. These companies place at the center of their business the assumptions that model capabilities will continue to improve through generational updates and that humanity will remain interested in and continue using AI. On that basis, they have raised tens of billions of dollars from capital markets and have undertaken some of the world's largest data-center investments. Their revenues presuppose that corporate and individual users will continue to use generative AI.

If the businesses of major AI companies collapse, the theoretical value of this paper is also likely to be largely lost. At the same time, however, the very problem this paper warns against, the hollowing-out of oversight caused by  $V \times L > C_{\max}$ , would also disappear automatically.

This point also has implications for current AI governance frameworks, including the EU AI Act. Existing frameworks likewise design regulations and institutions on the assumption that AI capabilities will continue to improve and that humans will continue to use AI.

The scale of risks produced by generative AI has expanded with each generation of models. In 2019, OpenAI initially withheld the release of GPT-2 because of concerns about impersonation. During the diffusion of text and image generation models after 2023, the effects of deepfakes on elections and personal reputation became a major issue. With the mass diffusion of interactive models after 2024, chains of hallucination and users' psychological dependence came to be regarded as concerns at the level of the public cognitive environment. In April 2026, with Claude Mythos Preview, capability for vulnerability discovery in cybersecurity was publicly discussed from the standpoint of safeguarding national financial infrastructure.

The risks of each generation are not resolved by the arrival of the next generation. They accumulate, and the scale of their targets continues to expand from individuals to national infrastructure. At each point, institutional responses centered on strengthening human checking systems have been proposed. Yet as of April 2026, no case can be confirmed in which the response to the risks of a previous generation was completed by the time the next generation arrived.

The empirical fact that  $V \times L$  continues to expand is also suggested by observing this very history of accumulation and expansion.

## 7. Conclusion

The simulation results presented in this paper are not absolute predictions of how many years remain before a particular domain breaks down. Rather, they show the robustness of the relative ordering among strategies in terms of how much time each strategy can buy. The timing at which  $C_{\max}$  is exceeded in each domain must be estimated separately on the basis of two variables: the  $B$  variable proposed in the previous work (Naito, 2026) and the rate of AI adoption. Accumulating empirical data for each domain therefore remains a task for future research.

The central contribution of this paper is the derivation of four design invariants that any flow-control mechanism must satisfy in high-loss domains where content judgment cannot be the foundation of governance: no substantive content judgment, no scalable consumption of examiner capacity, identity-bound per-application friction, and no batch clearance. The physical waiting path is presented as one reference implementation that satisfies all four invariants, but it faces substantial difficulties in accessibility, legal compatibility, and international coordination. The theoretical value of this paper does not depend on the feasibility of the physical waiting path. It depends on whether the four invariants correctly characterize the necessary conditions for content-judgment-bypass flow control.

We emphasize that this paper does not deny the usefulness of strengthening human oversight regimes, increasing the number of supervisors, or introducing double-checking in high-loss domains. The problem we identify is that the marginal returns of these investments decline over time and eventually fail to keep pace with the growth of  $V \times L$ .

The condition under which responsible human checking can continue to function in the age of generative AI is to redesign the incentives on the submission side, thereby slowing the growth rate of  $V \times L$  itself, and to embed a mechanism that self-selectively separates AI-enabled mass producers from legitimate human applicants.

Buying time against rapidly emerging AI risks, and responding to them substantively, must shift toward institutional limits on output rate and complexity. By avoiding the creation of economic incentives for AI-enabled mass producers, while keeping the false-positive cost for human-only applicants relatively low, it becomes possible to restore the sustainability of responsible human oversight in high-loss domains.

The Flow-by-Flow design may be difficult to apply to all high-loss domains with unbounded generator populations, especially in terms of compatibility with existing legal frameworks. However, alternative approaches also have their own inherent limitations. Economic staking, such as deposit systems, creates structural discrimination against applicants from developing countries and other resource-constrained contexts. Human endorsement systems run counter to digital transformation and may infringe the human-rights basis of the right to submit. Expanding qualification-based control over generating actors is ineffective against mass production by qualified actors. Online penalties can be avoided by moving across institutions.

Given that these alternatives all have serious defects, the legal challenges facing the physical waiting path are not grounds for abandoning the proposal. Rather, they indicate the need for long-term revision of legal frameworks themselves. The proposal in this paper is not a demand for immediate implementation. It is positioned as a theoretical starting point for the long-term revision of laws and treaties toward a paradigm of flow control.

A realistic implementation path may be to postpone public institutions, where treaty revision requires many years, and instead allow private platforms and academic journals that already face the collapse of quality signals due to AI-enabled mass production to implement flow control first as a form of self-defense, thereby accumulating evidence of its actual effects. This approach has broad applicability beyond high-loss domains. A design that allows society to enjoy the benefits of technology, such as AI-assisted drafting and analysis, while selectively suppressing only the vector of mass production that collapses the system through physical friction, deserves consideration across all knowledge infrastructures.

## Appendix A: Process-Time Declarations by the Authors

Section 3.6 of this paper proposed abolishing AI-use disclosure requirements and replacing them with process-time declaration requirements. This Appendix applies that system first to ourselves, the authors of this proposal.

### Previous work: The Supervision Paradox, v1.0-v1.4

Process	Declared time	Notes
Idea formation	0 h	Extension of an idea that had been continuously developed since around 2023.
Literature review	45 h	31 references; most were read in full.
Data collection and analysis, Appendix A	15 h	Approximately 25,000 records were obtained through APIs, classified by field, and used to calculate growth rates.
Writing, Japanese version	40 h	Includes all revisions from the initial draft through v1.4.
Figure and table preparation	4 h	Figures 1–3, Table A1, and Figure A1.
AI assistance, translation	2 h	Japanese-to-English translation was generated by AI, followed by approximately 1 h in total for checking.
Total	Approx. 106 h	

### This paper: Flow by Flow, v1.0

Process	Declared time	Notes
---------	---------------	-------

Idea formation	0 h	The idea arose during the writing of the previous work.
Literature review	10 h	Substantial overlap with the previous work.
Data collection and analysis, Mythos/Glasswing	8 h	Conducted as part of ongoing professional monitoring of recent developments.
Monte Carlo analysis, coding, execution, and interpretation	7 h	AI was used for coding.
Writing, Japanese version	60h	
Figure and table preparation	5 h	
AI assistance, translation	2 h	Japanese-to-English translation was generated by AI, followed by approximately 1 h for checking.
Total	Approx. 92 h	

Process-time declarations do not require perfect accuracy. The purpose of this Appendix is to show that the system is low-cost for the proposers themselves—the declaration took approximately five minutes to complete—while also providing information that can be verified by third parties.

## Appendix B

### *B-1 Why Accuracy Improvements Do Not Eliminate the Need for Flow Limitation*

The previous work responded to the objection that, if AI model capabilities improve, the error probability  $p$  will decline and expected loss will therefore decrease. The response was based on two points: first, that  $p = 0$  is in principle unattainable in probabilistic systems (Xu, Jain & Kankanhalli, 2024), and second, that  $V$  continues to expand faster than humanity's cognitive capacity improves. In addition, probabilistic systems cannot internally determine the correctness or incorrectness of their own outputs (Consistent Reasoning Paradox, Colbrook et al., 2024).

The  $V \times L$  framework further strengthens this response. The claim that accuracy improvements make flow limitation unnecessary presupposes that, once accuracy improves sufficiently,  $L$  becomes sufficiently small and  $V \times L$  remains within  $C_{\max}$ . However, as discussed above, only the verification component of  $L$  can decline through accuracy improvements. The response component is independent of AI accuracy.

The objection that “if AI accuracy improves sufficiently, response work itself can be delegated to AI, and therefore human cognitive load will continue to decrease” is unlikely to hold for the following reason. Even if response work is delegated to AI, in high-loss domains the legal responsibility for the result of that response is borne by a natural or legal person ( $R = 1$ ). Therefore, at the moment the response is delegated to AI, the cognitive load of verifying the response result arises. If that verification is also delegated to AI, then verification of the verification result arises. At the end of this recursion, there is always a human, and  $C_{\max}$  applies to that human.

Automating response does not eliminate cognitive load. It only transfers it recursively.

A simpler example is a real-name-registration social network. Even if 99.9% of users correctly register under their real names, once an incident occurs, it is still necessary to independently verify each time whether the account that made the problematic post is truly registered under a real name. The burden of verification is not determined by the incidence rate of fraud. It is determined by the possibility that fraud is not zero. No matter how much accuracy improves, as long as the system is probabilistic,  $p > 0$ . And as long as  $p > 0$ , the verification burden remains.

One reason this mathematical fact has not become a premise of AI governance debates may be a bias widely demonstrated in cognitive science. Weinstein (1980, 1982) identified unrealistic

optimism as a systematic human tendency to underestimate future risks, and showed that this bias is caused by mistaken extrapolation from past experience: “because no problem occurred before, no problem will occur in the future.” Since this bias has been observed universally across age, gender, educational level, and occupation, and has been replicated in more than 1,000 studies (Shepperd, Klein, Waters & Weinstein, 2013), it is natural to assume that it also applies to AI.

In the context of AI, the problem is that this bias operates across domains. At present, the majority of everyday AI use occurs in low-risk domains, such as drafting business emails. Even when errors are present, they are unlikely to be discovered, and even if discovered, the damage is minor. Experiences accumulated in this environment — “even if there were errors, things were fine” — distort risk estimation in high-loss domains through the mechanism identified by Weinstein.

This distortion is also reinforced by the market-formation strategy of the AI industry. Commercial diffusion of AI began in low-risk domains such as summarization, translation, idea generation, and copywriting. In these domains, errors can be offset by gains in total productivity, and the logic that “errors are tolerable as long as the overall result pays off” has been practically viable. The market perception formed under this logic — that AI will become sufficiently usable as long as accuracy improves — may be directly extrapolated to high-loss domains.

However, this extrapolation is precisely what existing institutional designs of professional licensing have explicitly rejected. Professions such as lawyers, physicians, tax accountants, and pilots require legal qualifications in every jurisdiction. Summarization and idea generation do not require national licenses. This difference is not accidental. Professional licensing systems are the institutional expression of a social judgment that even a small number of errors is unacceptable in the relevant domain. In other words, the training process required for qualification embeds the premise that “judgment in this domain is qualitatively different from ordinary work that requires no license.”

In domains where errors are unacceptable, legal constraints are imposed on business designs that may induce errors in the first place. Clinical-trial requirements in medicine, operating rules in aviation, and internal-control requirements in finance are all based on the idea of regulating in advance the designs under which errors may occur, rather than responding to errors only after they occur.

There is a clear asymmetry across domains in AI output. With respect to legal advice, because regulations corresponding to unauthorized-practice-of-law rules exist in many countries, LLM providers can plausibly justify restricting outputs that constitute legal advice. Although such restrictions are fragile and contain many practical loopholes, they can at least function as a form of physical friction similar to flow-rate limitation.

By contrast, no LLM safety system restricts tasks such as drafting patent claims, writing academic papers, or designing clinical trials. An AI system that refuses all scientific questions or technical interests in the name of safety would not be commercially viable. This asymmetry means that voluntary brakes on the LLM side against V explosion exist only in a subset of high-loss domains. Surges in patent applications or academic submissions cannot be stopped by LLM tuning. Since there is no restriction on the output side, flow limitation on the receiving side becomes unavoidable.

This problem can be restated through the pharmaceutical-approval example used in the previous work. Suppose that a review agency processes 100,000 applications per year with AI assistance and that accuracy improvements reduce the error rate to 0.01%. Verification costs decline. However, even for correctly approved drugs, response work occurs item by item, including post-market safety monitoring, evaluation of adverse-event reports, and revisions to package inserts. Once the number of cases expands to 100,000, even if verification cost is zero, the sum of response costs may exceed human processing capacity.

Error-rate improvements postpone the collapse of oversight, but they cannot eliminate it. The previous work described error-rate improvement as a treadmill: the system must run faster merely to remain in the same place, and even a slight decrease in running speed immediately turns into an increase in absolute harm. This characteristic becomes even clearer under the framework of the present paper.

Bastani and Cachon (2025) independently derive the economic dimension of this problem from a contract-theoretic framework. In their model, as AI accuracy improves, errors become rare, and opportunities for supervisors to actually detect errors and receive rewards decline. As a result, the compensation required to economically motivate supervisory effort diverges. In other words, accuracy improvement collapses the incentive design of supervision.

The argument in this paper and the argument by Bastani and Cachon are complementary. Bastani and Cachon show that even when a supervisory regime is complete, supervision cannot be economically motivated. We show that even if supervision were economically motivated, the supervisor's cognitive processing capacity would be exceeded. When both hold simultaneously, supervision enhancement has no solution either in terms of incentives or in terms of cognitive capacity. This more strongly supports the conclusion of the previous work that institutional limits on output volume itself are necessary.

Verification cost also depends strongly on the time required to detect the first breakdown contained in an output. Low-quality outputs that contain obvious factual errors or outdated information can be rejected at an early stage of verification, consuming relatively little cognitive cost. As AI capabilities improve, however, such easily detectable breakdowns decrease. Outputs that appear coherent on the surface and formally appropriate, but contain errors that only an expert can detect after reading the entire text carefully, maximize verification cost. In other words, AI capability improvements may increase  $L$  by raising the difficulty of detecting breakdowns, rather than lowering the verification component of  $L$ .

It is important to distinguish the task of confirming whether an output is correct from the response work that occurs after confirmation, because AI capability improvement acts asymmetrically on these two components. Intuitively, one might think that improved AI model accuracy reduces the burden of verification. In practice, however, accuracy improvement operates less by making verification easier than by inducing its omission.

For outputs from a model that is 99% accurate, it is extremely difficult, both cognitively and economically, for supervisors to maintain independent verification item by item. Moreover, because accuracy improvement increases the absolute number of correct outputs, the response burden increases. If a model with 30% accuracy produces 100 outputs, 30 outputs are correct and proceed to downstream processing. If a model with 99% accuracy produces 100 outputs, 99 outputs are correct, and 99 outputs proceed to downstream processing. Verification is required for all 100 outputs in both cases. Therefore, accuracy improvement increases the number of downstream cases from 30 to 99 while not eliminating the need to verify all items. Furthermore, because errors from a 99% accurate model appear only rarely, if the human checking regime has become nominal, the condition discussed in the previous work —  $E_{\text{detected}} \ll E_{\text{actual}}$  — becomes more severe.

### *B-2 On AI-Use Disclosure Regimes*

Regimes requiring disclosure of whether AI was used are being introduced in many domains. However, if there is no means of verifying the truth of the disclosure, supervisors cannot exclude the possibility that AI was involved. Even for a work product declared to have been produced without AI, supervisors must read it on the assumption that AI may have been involved.

This uncertainty creates an additional triage burden before the content of the output can be evaluated: the supervisor must triage the origin of the output and the evaluation criteria to be applied. In an environment where AI use has become widespread, even work products created without AI must be verified on the assumption that they may be low-quality AI-generated artifacts. The diffusion of AI therefore increases supervisory cost even when AI was not used.

These are examples of a vulnerability common to self-disclosure regimes. For penalties against false disclosure to function effectively, there must be an independent means of verifying the content of the disclosure. However, no reliable method has been established for post hoc distinguishing between generative AI output and human-created work.

Liang et al. (2023) evaluated seven widely used GPT detectors and showed that more than half of essays written by non-native English speakers were misclassified as AI-generated, while simple prompt manipulation reduced detection rates from 100% to 13%. OpenAI's own AI text classifier, released in 2023, correctly identified AI-generated text only 26% of the time and was discontinued because of insufficient accuracy.

Thus, approaches that automatically determine the origin of content using AI lack reliability in both false positives and false negatives. They cannot be used as the foundation of a human supervisory regime.

### *B-3 Why Triage Is Costly*

Why triage consumes substantial human cognitive resources can be explained by findings in cognitive science that have accumulated independently of AI research. The observation that conscious processing concentrates resources on selection among multiple candidates was systematized in Baars's (1988) global workspace theory. Among the countless processes running in parallel in the brain, only those that compete with one another and require integrated judgment are elevated into the conscious workspace, where they compete for limited bandwidth.

Dehaene and Naccache (2001) extended this theory neuroscientifically and showed that conscious access involves global activation of distributed cortical networks. The basic structure — that conscious processing is a scarce resource and is intensively consumed when conflicts among multiple candidates must be resolved — has been reproduced repeatedly in subsequent research (Dehaene et al., 2017).

When parallel processes are resolved automatically, this scarce resource is not consumed. Schneider and Shiffrin (1977) demonstrated that repeatedly trained information processing becomes automatic processing that does not require consciousness, while processing in novel and ambiguous situations becomes controlled processing that requires conscious resources. Kahneman's (2011) distinction between System 1 and System 2 restates this distinction for a general audience. The important question is what activates System 2, and the answer is ambiguity that cannot be resolved by automatic response.

Within decision-making contexts, depletion of finite decision-making resources has also been observed. Baumeister et al. (1998) showed that the quality of judgment declines in subjects who are repeatedly required to make decisions or exercise self-control, and named this phenomenon ego depletion. Replications of this effect vary, and debate continues regarding effect size (Hagger et al., 2016). Nevertheless, the general direction — that there is an upper bound on the number of high-quality judgments that can be made in a day — is independently supported by research on decision fatigue (Vohs et al., 2008; Danziger et al., 2011). Judgment resources are not infinite. Once consumed, they require time to recover.

What these findings show is the following proposition: conscious judgment is a scarce resource, and what triggers its consumption is not the amount of information, but the presence of ambiguity. Even if the amount of information increases, resources are not consumed if it can be processed by automatic response. Conversely, even if the amount of information is small, conscious resources are drawn upon when multiple interpretations coexist and cannot be resolved automatically.

When an AI-generated text may be a fact, an inference, an opinion, or a fictional narrative, and none of these interpretations can be rejected automatically, the reader must place multiple interpretations in the conscious workspace and choose among them. As a trigger for conscious processing, this satisfies a classical condition. The fact that generative AI outputs consume substantial cognitive resources is not an AI-specific phenomenon. Rather, AI outputs systematically generate the conditions described by consciousness research since Baars.

It is also logically derived from the design of generative AI that automatic response is difficult to establish, because generality and human-likeness are central to its commercial value. The ambiguity of triage is therefore not a transient problem caused by the immaturity of current AI. It is

a property inherent in the design of general-purpose AI. As AI capabilities improve and output style becomes closer to human speech, triage ambiguity increases rather than decreases.

The information sources that human cognitive systems have processed over tens of thousands of years can be roughly divided into the natural environment, other animals, and other humans. Allocation of conscious resources toward these sources has been adjusted by selection pressures on an evolutionary time scale. Generative AI supplies outputs that are statistically very similar to speech from other humans, but are not produced by other humans, at a rate whose marginal cost approaches zero. Put differently, it continuously fires triggers that ignite conscious processing at a rate exceeding human processing capacity. The finitude of conscious resources is an evolutionarily fixed constraint. Since that constraint is not removed by improvements in AI capability, the accumulation of triage cost inevitably collides with finite resources.

#### *B-4 Human Supervisors as Discriminators*

The development process of commercial interactive AI can be characterized as adversarial training between a generator and a discriminator (Goodfellow et al., 2014). The generator produces outputs. The discriminator evaluates whether those outputs are plausible as human speech, and the evaluation result is fed back into updates of the generator. RLHF (Ouyang et al., 2022) is a large-scale implementation of this structure through human comparison evaluation and reward models.

The role of discriminator is played, at the training stage, by human evaluators who provide preference data, and at the deployment stage, by human users who use the service. The generator is updated with each model generation and moves closer to more human-like output by reflecting the history of this adversarial process. Human cognition, which plays the role of discriminator, is updated only on an evolutionary time scale and is nearly fixed on the time scale of product cycles.

As the adversarial process repeats across generations, the cues that allow the discriminator to distinguish the generator's output as "something other than human speech" diminish. This is an explanation from the generator side of the pathway, discussed in the previous section, through which triage cost rises over time. The important point is that this asymmetric adversarial process is not an accidental technical side effect. It is part of the very definition of capability improvement in current AI.

The return of triage and criterion invocation from the device side to the user side is already observed in AI evaluation. In generational comparisons of physical products, agreement forms when evaluation can be reduced to a hierarchy of physical quantities. The superiority of an old television and the latest 8K television can be judged by external criteria such as resolution, contrast ratio, and response time. Even in comparisons among video games, once factors such as map size, number of missions, and resolution are decomposed, opinions converge within each factor, allowing discussion to move on to what individuals personally value.

By contrast, in generational comparisons of AI models, the release of a new model is repeatedly accompanied by the coexistence of polar evaluations such as "AGI has arrived" and "it is worse than the previous generation." Benchmarks have lost reliability for comparison because of contamination in training data and leakage of evaluation sets (Sainz et al., 2023; Balloccu et al., 2024). Different use cases require different axes for evaluating intelligence, and users themselves must determine which axis should be invoked.

Indistinguishability does not eliminate triage. It increases triage cost. Even when discrimination is impossible, empirical evidence shows that supervisors continue to switch evaluation criteria depending on whether an output is AI-derived (Longoni et al., 2022; Altay & Gilardi, 2024). Therefore, the attempt at triage itself continues. And as discrimination accuracy approaches chance level, the cognitive resource consumption per attempt increases.

If supervisors abandon triage and process all outputs under a uniform criterion, they must choose either to treat all outputs as AI-generated and apply maximum verification cost to all items, or to treat all outputs as human-generated and omit verification. The former maximizes the judgment component of L for all items. The latter means abandoning supervision in high-loss domains.

### *B-5 The Cognitive-Load Externalization Function of Information Infrastructure*

Many information infrastructures used by humans have functioned by moving individual cognitive load outside the individual. Writing moved memory retention from the brain to material objects (Goody, 1977; Ong, 1982). Money moved comparison of the value of goods from individual negotiation to a common unit. Double-entry bookkeeping moved the consistency check of transactions from memory and mental arithmetic to mechanical verification in ledgers. Legal systems moved the judgment of dispute resolution from the parties to procedures and precedents. Scientific publication and peer review moved the evaluation of the truth of claims from individual persuasiveness to reproducibility and mutual examination. Internet search indexes moved memory of where information is located from individuals to search engines.

What these systems share is a design direction in which cognitive load, instead of being completed within the individual, is distributed through external materials, procedures, and institutions. Hutchins (1995) organized this phenomenon as distributed cognition, and Clark and Chalmers (1998), under the concept of the extended mind, presented a framework for treating cognitive systems that include external devices as the unit of analysis.

A common property of these information infrastructures is that they are designed to suppress users' triage cost. When looking up a dictionary entry, the order of entries is fixed by an external criterion, alphabetical order. The user does not have to re-determine "this is the headword," "this is the definition," or "this is an example." Because such designs process triage in advance through devices or institutions, conscious resources can be concentrated on judgment and response.

Generative AI output runs counter to this externalization. As discussed above, because generative AI output does not fix boundaries among semantic types on the device side, triage returns to the user side. This is not a design feature of individual applications. It is a design consequence that cannot be avoided as long as generality and human-likeness are placed at the center of commercial value.

Triage costs that had historically been externalized must therefore be processed again by individual cognitive resources when the new information source is generative AI output. The diffusion of generative AI also spreads into externalized information infrastructures themselves. For information sources such as search-engine results, encyclopedia entries, preprints, and news articles — sources for which triage had previously been processed by the device or institution side — the contamination of generative-AI-derived outputs forces users to determine item by item whether a given item originates from generative AI. Triage that had been omitted under the premise that the source itself was reliable returns at the level of each item.

Multiple experiments have repeatedly observed that this triage is processed before other evaluation axes. Longoni, Fradkin, Cian, and Pennycook (2022) showed that the credibility evaluation of a news article declines when the same content is labeled as "AI-generated." Altay and Gilardi (2024), in a preregistered experiment (N = 4,976), reported that attaching an "AI-generated" label to a headline reduces both perceived accuracy and intention to share, regardless of whether the content is true or false and regardless of whether it was actually written by a human or by AI. These results show that determining whether content is AI-derived is processed independently of, and often prior to, the evaluation axis of the content itself.

### *B-6 Consequences of the Hollowing-Out of Human Oversight*

When  $V \times L$  exceeds  $C_{max}$ , the institution does not necessarily stop immediately. Rather, in many cases, supervisors can be expected to come to work as before, approve items as before, and report processing counts as before. From the outside, the institution appears to be operating normally. However, what is maintained in this state is not substantive oversight, but a flow of formal approval.

Supervisors are not sufficiently verifying individual outputs. In that case, it becomes rational for supervisors to rely on external signals rather than the content itself in order to minimize their own responsibility risk. In academic peer review, proxy indicators include institutional affiliation, existing

reputation, the author's country, institutional email addresses, and past citation records. In law, administration, and pharmaceutical review, similarly, the applicant's institutional status, representative's qualifications, and relationships with existing organizations become substitutes for content evaluation.

This occurs because supervisors are more likely to face clear responsibility if they allow AI slop to pass, or if they reject a valid application from a famous institution. By contrast, if they reject a legitimate low-signal application without reading it, the loss is likely to be invisibilized.

### *B-7 Limits of Content Judgment*

In current AI governance debates, human oversight almost always means content judgment. Humans judge whether the output text is accurate, whether generated code contains vulnerabilities, or whether a drafted document contains legal defects.

However, as this chapter has shown, as long as content judgment remains the foundation of governance, if AI capability continues to grow beyond the sum of human capability improvement and institutional reinforcement, many high-loss domains cannot avoid the collapse of oversight. Delegating content judgment to AI itself does not solve this problem.

Nor is this problem limited to hallucination. For example, even when an AI-generated summary consists only of facts contained in the original document, it may reverse the intention of the original document through the choice of what to retain and what to omit. Moreover, measures including source disclosure, retrieval-augmented support, cross-checking by multiple LLMs, output constraints, and all forms of prompt engineering may reduce errors through improvements in models or operating methods, but they cannot reduce errors to zero.

If AI is asked to judge whether an output is correct, the question of who verifies the correctness of that AI judgment arises recursively. Therefore, what is needed is not to improve content-judgment capability, but to bypass it.

## **Appendix C. Consistency with the Mythos Release Problem**

In April 2026, Anthropic restricted access to the preview version of Claude Mythos to approximately forty institutions, under Project Glasswing, on the grounds that the model's capabilities in the cybersecurity domain had improved beyond expectations. Regardless of whether the reason for this measure was the model's autonomous threat potential or the exceedance of human supervisory capacity, the result was a flow-rate limitation.

The fact that cybersecurity was the first domain to reach lockout was predictable in advance from the ordering implied by the B variable, because the physical-space constraint in this domain is almost zero ( $B \approx 0$ ). This paper positions the measure as the first visible case of Weakest-Link Lockout. However, the theoretical claim of this paper does not depend on the consequences of this single case.

## **Appendix D. Additional Discussion of the Physical Waiting Path**

This appendix collects additional analysis of the physical waiting path that was moved from the main text in order to maintain focus on the four design invariants as the central contribution.

### *D-1 Relationship to Computational Proof-of-Work*

The physical waiting path also belongs to the lineage of computational proof-of-work proposed by Dwork and Naor (1993) as a defense against spam email. Traditional anti-machine screening mechanisms, represented by CAPTCHA (von Ahn et al., 2003), have functioned by imposing cognitive tasks that are low-cost for humans and high-cost for AI. Yet these mechanisms have been repeatedly bypassed as AI capabilities improve, and many implementations are now circumventable (Searles et al., 2023). Because the physical waiting path requires physical time expenditure rather than a cognitive task, it realizes a form of proof-of-work that is robust against improvements in AI capability.

### *D-2 Cooling-off Period as an Alternative*

As an alternative to the physical waiting path, one could imagine imposing on an account a prohibited period for applications or submissions, namely a cooling-off period, in proportion to the exceedance multiplier of the cognitive cost score. Under this design, however, the optimal strategy for AI-enabled mass applicants would be to incur a large penalty at one institution and then immediately move to another. Preventing this evasion would require a massive surveillance infrastructure for sharing KYC and penalty information across institutions, creating both the construction and operating costs of that infrastructure and the secondary problem of concentrated power in the hands of infrastructure administrators.

### *D-3 Capital-Rich Entities and Proxy Labor*

A capital-rich corporation may distribute the burden of physical waiting among employees, but the personnel cost per employee, including wages during waiting time, scales with the number of employees. Even an actor that can generate 100,000 outputs without concern for API fees faces costs incomparable to API fees if it must employ 1,000 employees and require each of them to wait physically 100 times. Even if humanoid robots become practical, the purchase and maintenance costs of 1,000 robots are physical costs of an entirely different order from mass API calls in information space. This is an institutional operation that artificially raises  $B$ , the physical-space constraint variable defined in the previous work.

The purpose of the exceedance pathway is not to make outsourcing or proxy use impossible. Its purpose is to ensure that each threshold-exceeding application remains attached to a non-zero, identity-bound, per-application marginal cost. The physical waiting path is one implementation of this principle, but any mechanism satisfying the four invariants would serve the same institutional function.

### *D-4 False-Positive Criticism*

For the false-positive criticism against composite flow as applied to non-AI humans to hold, the cost of false positives under flow limitation must exceed the quality-degradation cost borne by all applicants under no flow limitation. This presupposes that the probability that LLMs do not increase  $V \times L$  is at least 50%. The empirical data of Kusumegi et al. (2025) suggest that this premise is not supported.

**Availability Statement:** The simulation has also been implemented as an HTML file that can be run in a browser and is publicly available at UTIE Research Institute (<https://utie-instruments.com/utie-research-institute.html>). Parameters can be changed and the simulation can be rerun freely.

## References

- AI & Society (2025). Exploring automation bias in human-AI collaboration: A review and implications for explainable AI. AI & Society. <https://doi.org/10.1007/s00146-025-02422-7>
- Altay, S. & Gilardi, F. (2024). "AI-generated" label reduces perceived accuracy and sharing of true and false news, regardless of actual AI-generation or political concordance. Preregistered experiment (N=4,976).
- Alter, A. L. & Oppenheimer, D. M. (2009). Uniting the tribes of fluency to form a metacognitive nation. *Personality and Social Psychology Review*, 13(3), 219-235.
- Anthropic (2026). Project Glasswing: Claude Mythos Preview limited access announcement. April 2026.
- Baars, B. J. (1988). *A Cognitive Theory of Consciousness*. Cambridge University Press.
- Bainbridge, L. (1983). Ironies of automation. *Automatica*, 19(6), 775-779.
- Balloccu, S., Schmidtke, D., Lango, M. & Dusek, O. (2024). Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source LLMs. Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (EACL 2024).

- Bastani, H. & Cachon, G. P. (2025). The human-AI contracting paradox. SSRN Working Paper. <https://doi.org/10.2139/ssrn.5962739>
- Baumeister, R. F., Bratslavsky, E., Muraven, M. & Tice, D. M. (1998). Ego depletion: Is the active self a limited resource? *Journal of Personality and Social Psychology*, 74(5), 1252-1265.
- Beck, U. (1986). *Risikogesellschaft: Auf dem Weg in eine andere Moderne*. Suhrkamp. [English edition: *Risk Society: Towards a New Modernity*, trans. M. Ritter, Sage, 1992.]
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877-1901.
- Carnat, I. (2024). Human, all too human: Accounting for automation bias in generative large language models. *International Data Privacy Law*, 14(4), 299-314.
- Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of cognitive science. *Behavioral and Brain Sciences*, 36(3), 181-204.
- Clark, A. & Chalmers, D. J. (1998). The extended mind. *Analysis*, 58(1), 7-19.
- Colbrook, M. J., Antun, V. & Hansen, A. C. (2024). The consistent reasoning paradox. Working paper.
- Cowan, N. (2001). The magical number 4 in short-term memory: A reconsideration of mental storage capacity. *Behavioral and Brain Sciences*, 24(1), 87-114.
- Danziger, S., Levav, J. & Avnaim-Pesso, L. (2011). Extraneous factors in judicial decisions. *Proceedings of the National Academy of Sciences*, 108(17), 6889-6892.
- Dehaene, S. & Naccache, L. (2001). Towards a cognitive neuroscience of consciousness: Basic evidence and a workspace framework. *Cognition*, 79(1-2), 1-37.
- Dehaene, S., Lau, H. & Kouider, S. (2017). What is consciousness, and could machines have it? *Science*, 358(6362), 486-492.
- Elish, M. C. (2019). Moral crumple zones: Cautionary tales in human-robot interaction. *Engaging Science, Technology, and Society*, 5, 40-60.
- European Parliament (2024). Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence (Artificial Intelligence Act). *Official Journal of the European Union*.
- Friston, K. (2010). The free-energy principle: A unified brain theory? *Nature Reviews Neuroscience*, 11(2), 127-138.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. & Bengio, Y. (2014). Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27, 2672-2680.
- Goodhart, C. A. E. (1975). Problems of monetary management: The U.K. experience. *Papers in Monetary Economics*, Reserve Bank of Australia.
- Goody, J. (1977). *The Domestication of the Savage Mind*. Cambridge University Press.
- Hagger, M. S., Chatzisarantis, N. L. D., Alberts, H., Anggono, C. O., Batailler, C., Birt, A. R., ... & Zwieneberg, M. (2016). A multilab preregistered replication of the ego-depletion effect. *Perspectives on Psychological Science*, 11(4), 546-573.
- Horowitz, M. C. & Kahn, L. (2024). Bending the automation bias curve: A study of human and AI-based decision making in national security contexts. *International Studies Quarterly*, 68(2), sqae020.
- Hutchins, E. (1995). *Cognition in the Wild*. MIT Press.
- Kahneman, D. (1973). *Attention and Effort*. Prentice-Hall.
- Kahneman, D. (2011). *Thinking, Fast and Slow*. Farrar, Straus and Giroux.
- Kücking, F., Hübner, U., Przysucha, M., Rupp, T. & Saranto, K. (2024). Automation bias in AI-decision support: Results from an empirical study. *Studies in Health Technology and Informatics*, 317, 298-304.
- Kusumegi, K., Yang, X., Ginsparg, P., de Vaan, M., Stuart, T. & Yin, Y. (2025). Scientific production in the era of large language models. *Science*, 390(6779), 1240-1243. <https://doi.org/10.1126/science.adw3000>
- Liang, W., Yuksekgonul, M., Mao, Y., Wu, E. & Zou, J. (2023). GPT detectors are biased against non-native English writers. *Patterns*, 4(7), 100779.
- Longoni, C., Fradkin, A., Cian, L. & Pennycook, G. (2022). News from generative artificial intelligence is believed less. *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, 97-106.

- Naito, H. (2026). The supervision paradox: AI capability growth necessitates usage contraction in high-loss domains. Preprints.org, 2026030150 (v4). <https://doi.org/10.20944/preprints202603.1507.v4>
- Oakley, B., Mouloua, M. & Hancock, P. A. (2003). Effects of automation reliability on human monitoring performance. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, 47(1), 188-192.
- Ong, W. J. (1982). *Orality and Literacy: The Technologizing of the Word*. Methuen.
- OpenAI (2026). GPT-5.4-Cyber and Trusted Access for Cyber program announcement. April 2026.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35, 27730-27744.
- Parasuraman, R. & Manzey, D. H. (2010). Complacency and bias in human use of automation: An attentional integration. *Human Factors*, 52(3), 381-410.
- Parasuraman, R. & Riley, V. (1997). Humans and automation: Use, misuse, disuse, abuse. *Human Factors*, 39(2), 230-253.
- Parasuraman, R., Molloy, R. & Singh, I. L. (1993). Performance consequences of automation-induced "complacency." *International Journal of Aviation Psychology*, 3(1), 1-23.
- Park, J., Kim, M. & Han, S. (2026). The enrichment paradox: Critical capability thresholds and irreversible dependency in human-AI symbiosis. arXiv:2603.24391.
- Perrow, C. (1984). *Normal Accidents: Living with High-Risk Technologies*. Basic Books.
- Radford, A., Narasimhan, K., Salimans, T. & Sutskever, I. (2018). Improving language understanding by generative pre-training. OpenAI Technical Report.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D. & Sutskever, I. (2019). Language models are unsupervised multitask learners. OpenAI Technical Report.
- Sainz, O., Campos, J. A., García-Ferrero, I., Etxaniz, J., de Lacalle, O. L. & Agirre, E. (2023). NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark. *Findings of EMNLP 2023*.
- Salop, J. & Salop, S. (1976). Self-selection and turnover in the labor market. *Quarterly Journal of Economics*, 90(4), 619-627.
- Schneider, W. & Shiffrin, R. M. (1977). Controlled and automatic human information processing: I. Detection, search, and attention. *Psychological Review*, 84(1), 1-66.
- Searles, R. et al. (2023). An empirical study & evaluation of modern CAPTCHAs. *Proceedings of the 32nd USENIX Security Symposium*.
- Shepperd, J. A., Klein, W. M. P., Waters, E. A. & Weinstein, N. D. (2013). Taking stock of unrealistic optimism. *Perspectives on Psychological Science*, 8(4), 395-411.
- Singh, I. L., Molloy, R. & Parasuraman, R. (2009). Automation-induced monitoring inefficiency: Role of display location. *International Journal of Human-Computer Studies*, 40(4), 687-702.
- Spence, M. (1973). Job market signaling. *Quarterly Journal of Economics*, 87(3), 355-374.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł. & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 5998-6008.
- Vohs, K. D., Baumeister, R. F., Schmeichel, B. J., Twenge, J. M., Nelson, N. M. & Tice, D. M. (2008). Making choices impairs subsequent self-control: A limited-resource account of decision making, self-regulation, and active initiative. *Journal of Personality and Social Psychology*, 94(5), 883-898.
- von Ahn, L., Blum, M., Hopper, N. J. & Langford, J. (2003). CAPTCHA: Using hard AI problems for security. *Advances in Cryptology - EUROCRYPT 2003*, 294-311.
- Warm, J. S., Parasuraman, R. & Matthews, G. (2008). Vigilance requires hard mental work and is stressful. *Human Factors*, 50(3), 433-441.
- Weinstein, N. D. (1980). Unrealistic optimism about future life events. *Journal of Personality and Social Psychology*, 39(5), 806-820.
- Weinstein, N. D. (1982). Unrealistic optimism about susceptibility to health problems. *Journal of Behavioral Medicine*, 5(4), 441-460.
- Xu, Z., Jain, S. & Kankanhalli, M. (2024). Hallucination is inevitable: An innate limitation of large language models. arXiv:2401.11817.

- Jakesch, M., Hancock, J. T. & Naaman, M. (2023). Human heuristics for AI-generated language are flawed. *Proceedings of the National Academy of Sciences*, 120(11), e2208839120. <https://doi.org/10.1073/pnas.2208839120>
- Chein, J. M., Martinez, S. A. & Barone, A. R. (2024). Human intelligence can safeguard against artificial intelligence: individual differences in the discernment of human from AI texts. *Scientific Reports*, 14, 25989. <https://doi.org/10.1038/s41598-024-76218-y>

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.