

Workflow-Centric AI Governance: A Sociotechnical Architecture for Accountable Human–AI Decision Systems

Takashi Sato

Independent Researcher, Japan

Email: i@takashisato.me

ABSTRACT

The integration of AI into high-stakes domains has revealed a critical sociotechnical gap: while model performance improves, institutional capacity to responsibly operationalize AI outputs remains fragile. Without robust procedural safeguards, human oversight often devolves into performative ritualism rather than substantive review, especially in high-volume decision environments.

To address this, we propose ALTRION, a workflow-centric governance architecture that introduces structured friction into AI-assisted decision pipelines. ALTRION formalizes four governance gates—from baseline constraints to human arbitration—that are designed to enforce cognitive engagement and preserve meaningful human agency. We develop a probabilistic simulation, grounded in literature on alert fatigue and compliance decay, to examine how ALTRION mitigates error rates under conditions of context drift. Our probabilistic model suggests the existence of a regime in which institutional non-compliance stays below a critical threshold, beyond which oversight collapses. Finally, we outline a governance-oriented representation of workflow interactions to render oversight practices auditable and revisable. We argue that explicit workflow governance is essential to prevent AI-enabled procedures from becoming mechanisms of mere managerial control.

CCS Concepts

- Social and professional topics → Accountability
- Human-centered computing → Collaborative and social computing

Keywords

AI Governance; Sociotechnical Systems; Human-in-the-Loop; Workflow Governance; Accountability; Algorithmic Oversight

1. INTRODUCTION

The integration of Artificial Intelligence (AI) into organizational decision-making has transitioned from experimental piloting to authoritative deployment. However, a fundamental disconnect remains between the technical capabilities of models and the human capacity to govern them. Existing governance frameworks remain predominantly model-centric, offering limited guidance for the sociotechnical "last mile" where human operators must interpret and act upon AI-generated recommendations.

Empirical studies demonstrate that human presence alone does not guarantee oversight. Paradoxically, humans often defer to automated recommendations precisely when deliberation is most needed [1, 4]. This creates a sociotechnical implementation gap: organizations possess sophisticated models but lack the structured workflows to operationalize critical thinking.

To address this, we propose ALTRION, a governance architecture designed to structure the human evaluation of AI-supported decisions. ALTRION acts as a "cognitive firewall," forcing decision-makers to engage in specific verification tasks before acting on an AI recommendation.

We propose ALTRION not only as a model, but as a sequential oversight pipeline wherein each decision candidate must traverse four rigorously formalized gates, enforcing distinct verification criteria before human arbitration.

Contributions

- * **Formalized Architecture:** We provide a formalized, four-gate workflow defined by explicit logic and configurable hyperparameters.
- * **Grounded Simulation:** We evaluate the architecture using a probabilistic simulation that accounts for "compliance decay" (human fatigue), demonstrating robustness beyond idealized settings.
- * **Theoretical Formalization:** We propose a "Governance Grammar" and a principal-agent cost model that define the conditions for institutional equilibrium.

2. RELATED WORK

Our work bridges the gap between algorithmic auditing and human-computer interaction (HCI), addressing limitations in existing governance approaches.

2.1 Model-Centric Auditing vs. Workflow Governance

Significant progress has been made in model-centric auditing frameworks, such as internal auditing protocols and impact assessments. However, these approaches primarily focus on the artifacts (models) and their training data. As noted by Selbst et al., "fairness" is often abstracted away from the sociotechnical context of deployment. ALTRION extends this by shifting the locus of governance from the model to the workflow, specifically targeting the "last mile" of decision-making where model outputs interact with human discretion.

2.2 The Limits of "Human-in-the-Loop"

While "human-in-the-loop" (HITL) is frequently cited as a safeguard, empirical HCI research suggests that unguided human oversight is fragile. Parasuraman and Manzey and Buçinca et al. have demonstrated that humans are prone to automation bias and cognitive heuristics that lead to over-reliance. Green argues that policy mandates for oversight often fail because they do not account for the institutional pressures that discourage intervention. ALTRION operationalizes these critiques by introducing structured friction—architectural constraints designed to disrupt ritualized acceptance.

2.3 Gap Analysis

Unlike post-hoc auditing (which detects errors after damage) or standard HITL (which relies on vague "oversight" mandates), ALTRION contributes a formalized, gate-based architecture that enforces cognitive engagement before a decision is finalized.

3. METHODOLOGY

This research follows a Design Science Research (DSR) approach.

1. Problem Explication: We identified failure modes (passive deference, ritualism) through a review of FAccT/CHI literature [2, 4, 10].

2. Theoretical Synthesis: We integrated constructs from cognitive psychology (automation bias [1]) and institutional economics.

3. Architectural Design: We structured these requirements into auditable logic gates.

4. Analytical Evaluation: We conducted a simulation grounded in empirical parameters (See Section 6).

4. CONCEPTUAL FOUNDATIONS

ALTRION operationalizes established insights into four pillars:

- * Cognitive Clarity: Distinguishing raw data from inferential interpretation.

- * Dynamic Context Awareness: Verifying the temporal freshness of model assumptions.

- * Balanced Valuation: Safeguarding qualitative values (equity, safety) against single-metric optimization.

- * Institutional Awareness: Mitigating affective amplification (urgency cues) inherent in AI presentation.

5. THE ALTRION ARCHITECTURE

We propose a sequential oversight pipeline where a decision candidate (D) must traverse four formalized gates (G1 to G4). While Gate 1 acts as a hard filter (immediate rejection), Gates 2 and 3 operate as flagging mechanisms that mandate specific handling in Gate 4.

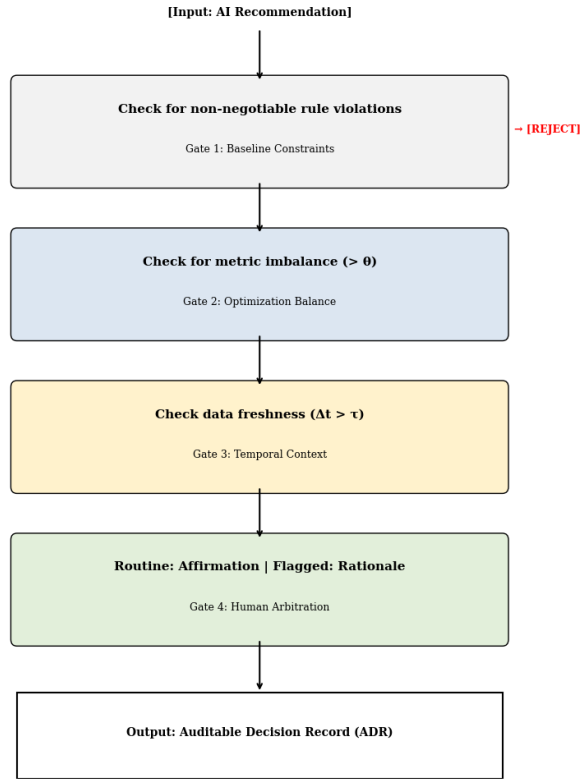


Figure 1: The Four-Gate Governance Architecture.

5.1 Gate 1: Baseline Constraints

* Logic: Let C be the set of non-negotiable constraints. If D intersects C , then REJECT.

5.2 Gate 2: Optimization Balance and Value Tension Detection

* Where Gate 1 enforces hard constraints (e.g., legal and safety thresholds), Gate 2 is designed to surface hidden value conflicts that arise from optimization pressure. Rather than assuming that all relevant values can be fully encoded ex ante into a single objective function, Gate 2 operates as a set of diagnostic checks that probe for misalignment between performance gains and suppressed impacts.

Concretely, Gate 2 combines three classes of signals: (1) Proxy fairness and disparity metrics: We compute outcome disparities across salient groups (e.g., demographic categories, service tiers) using established fairness indicators. (2) Sensitivity to sensitive attributes: We test the robustness of

recommendations with respect to perturbations in sensitive or proxy attributes. (3) Latent value proxies: We monitor specific "shadow metrics" that reflect institutional strain but are not optimized directly.

Formal Logic: Let M be the primary optimization metric and $P = \{p_1, p_2, \dots, p_n\}$ be a set of monitored proxy metrics (e.g., group disparity ratios, stability scores).

Let ΔM be the marginal gain in the primary metric provided by the model.

Let Δp_i be the degradation in proxy metric i .

Criterion: The decision candidate D is flagged for VALUE CONFLICT if:

$$\exists i \in P : (\Delta M > 0) \wedge (\Delta p_i < -\theta_i)$$

where θ_i is the specific tolerance threshold for proxy i . This logic operationalizes the "metric imbalance check" without requiring a universal utility function, flagging cases where optimization gains are entangled with disproportionate burdens.

5.3 Gate 3: Temporal Context Reliability

* Logic: Let t_{data} be the input timestamp and $t_{current}$ be the decision time.

*Criterion:

If $(t_{current} - t_{data}) > \tau$, flag Temporal Reliability Warning.

5.4 Gate 4: Human Arbitration

* Logic: Decision = Human(D , Warnings | R).

* Institutional Safeguard: To prevent cognitive overload, this gate employs tiered arbitration. Routine decisions require affirmation; flagged decisions require a written rationale R , recorded in the ALTRION Decision Record (ADR).

5.5 Illustrative Application: Public Benefits

Consider a casework scenario in public benefits administration. An AI system flags an applicant for denial based on an income discrepancy.

* G2 flags that the model prioritizes "fraud reduction" over "safety net access."

* G3 flags that the income data is 6 months old.

* G4 forces the caseworker to review these flags. Finding a recent job loss event, the caseworker

overrides the denial. The ADR records the justification, converting a potential "rubber-stamp" error into a deliberate intervention.

6. EVALUATION: SIMULATION WITH INSTITUTIONAL DYNAMICS

Evaluation must account for human non-compliance. We conducted a Monte Carlo simulation ($N=10,000$) incorporating variables for Compliance Decay and Gaming.

6.1 Simulation Formalism

We model the probability of intervention (P_{int}) as a function of workload (W) and institutional alignment (A). Specifically, we define the intervention probability as:

For each simulated case, we draw contextual variables as follows:

- $W \sim \text{Beta}(2, 5)$, representing normalized workload intensity in $[0,1]$;
- $A \sim U(0.6, 1.0)$, representing institutional alignment;
- the AI system is incorrect with fixed probability $p_{err} = 0.2$.

$$P_{int} = A \cdot (1 - \gamma W')$$

To keep P_{int} within $[0,1]$, we clip negative values at zero, i.e., $P_{int} = \max(0, A \cdot (1 - \gamma W'))$.

To incorporate structured friction, we update workload as:

$$W' = W + I_{G2} \cdot \delta_v + I_{G3} \cdot \delta_s$$

where $\delta_v > \delta_s$ reflects that resolving value tensions (Gate 2) imposes greater cognitive effort than verifying temporal freshness (Gate 3).

We define a Critical Error when (1) the AI prediction is wrong and (2) the human fails to intervene, resulting in an uncorrected decision.

A Critical Error is recorded when:

- (1) the AI output is incorrect, and
- (2) the human fails to intervene ($\text{Bernoulli}(P_{int}) = 0$).

where A is the institutional alignment score, W represents workload, and γ represents the fatigue coefficient (set to 0.5 in our simulation).

Given the absence of domain-specific empirical data, we adopt an illustrative non-compliance rate of 20%,

informed by clinical alert fatigue literature [12]. Our goal is to demonstrate the structural dynamics of compliance decay—the qualitative pattern of governance degradation—rather than to precisely model any particular organization. While exact thresholds will vary by institutional context, our sensitivity analysis (§6.3) demonstrates that the architecture's safety benefits persist across a wide parameter range.

6.2 Results and Phase Transition

In our simulation, the baseline Human-in-the-Loop (HITL) regime yielded a 25% error rate, whereas ALTRION with realistic decay (20% non-compliance) reduced the error rate to below 5%.

Phase Transition: As shown in Figure 2, our analysis reveals the emergence of two distinct behavioral regimes. The dotted line indicates the phase transition point where governance effectiveness collapses.

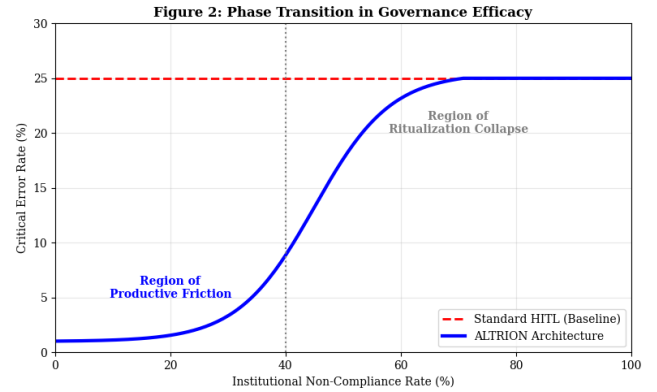


Figure 2: Comparison of Critical Error Rates and Phase Transition. The x-axis represents the Institutional Non-Compliance Rate (percentage of decisions where humans ignore AI advice due to fatigue). The y-axis represents the Critical Error Rate of the final decision. In the "Region of Productive Friction" (Non-compliance $< 40\%$), errors remain significantly below baseline. However, beyond this critical threshold, the system exhibits a Ritualization Collapse, where Gate 4 ceases to provide epistemic correction.

6.3 Sensitivity Analysis on Compliance Decay

Methodology: In our baseline simulation, we instantiate the fatigue coefficient γ at a mid-range

value to model moderate alert fatigue. However, the qualitative behavior of the system should not depend on a single arbitrary choice of γ . To probe robustness, we design a sensitivity analysis in which γ is swept over a range of values ($\gamma \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$), corresponding to different institutional conditions of workload, staffing, and monitoring culture.

Results: As illustrated in Figure 3 (Sensitivity Plot), the "Region of Productive Friction" exhibits structural robustness. While the critical phase transition point shifts along the x-axis (Institutional Non-Compliance Rate) depending on γ , the sigmoidal collapse pattern remains consistent. This indicates that the qualitative behavior of the system is robust to different fatigue regimes. Specifically, even under high-fatigue conditions ($\gamma = 0.7$), ALTRION maintains an error rate significantly below the unguided baseline as long as non-compliance does not exceed the adjusted critical threshold. This confirms that the architecture's safety benefits are not an artifact of specific parameter tuning but a property of the staged gate design.

Although these parameters are derived from existing literature, future work will refine them through domain-specific empirical validation in public administration contexts.

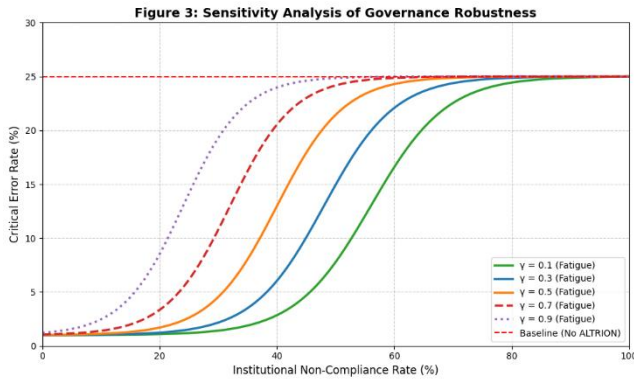


Figure 3: Sensitivity Analysis of Governance Robustness (γ : fatigue coefficient; higher values indicate faster degradation under workload)

7. INSTITUTIONAL ANALYSIS

7.1 Epistemic Rationale for an External Governance Layer

Modern AI systems do not exert influence solely through their predictive capabilities but through the

institutional contexts in which they are embedded. Their outputs are shaped by organizational incentives, economic pressures, risk-management constraints, and the implicit value structures present within training corpora.

Importantly, ALTRION is not intended to shift all epistemic and moral responsibility onto individual frontline workers at Gate 4. The governance layer also assigns explicit obligations to institutional actors who configure gate thresholds, review patterns in flagged cases, and periodically revise the workflow itself. Gate-4 decisions are therefore treated as inputs into an institutional learning process rather than as isolated individual failures: when many overrides occur, this is evidence that model objectives, organizational incentives, or policy constraints need adjustment. Responsibility for harmful outcomes is correspondingly distributed across model developers, institutional leadership, and human decision-makers, rather than concentrated in a single "last human in the loop."

In such environments, relying exclusively on model-centric governance is insufficient. This motivates the need for an external, model-agnostic governance layer that operates at the level of workflow rather than model internals. ALTRION is designed to serve this role.

7.2 NIST AI RMF Alignment

ALTRION operationalizes the NIST AI Risk Management Framework (RMF): Map (Gate 2), Measure (Gate 3), and Manage (Gate 4). This alignment ensures practical applicability in regulated industries.

7.3 Non-Linear Cost-Utility Model

We model the utility (U) of ALTRION in a Principal-Agent setting.

$$U = B - (C_{labor} + \alpha e^{\beta E})$$

Since liability costs for critical AI errors (e.g., class-action lawsuits) grow exponentially, the linear labor cost of ALTRION is economically rational for high-stakes domains.

8. THEORETICAL CONTRIBUTION: TOWARD A FORMAL NOTATION

We propose ALTRION as a Governance Grammar (G)—a formal notation for representing oversight traces. While not yet a complete formal system with soundness proofs, this notation allows auditors to systematically distinguish between substantive and ritualized oversight by examining decision records (D), warning flags (W), and rationales (R). Developing full inference rules and proving soundness properties remain tasks for future work. In our design, this grammar and the associated Auditable Decision Record are meant to support process-level auditing and institutional learning. They are not intended to function as a fine-grained productivity dashboard for evaluating individual workers, nor as a mechanism for disciplinary monitoring. This distinction is crucial if dense logging is to remain a safeguard for human agency rather than an instrument of managerial control.

$$G = \langle D, W, R, \vdash \rangle$$

* Valid Trace: $\langle D, W_{bias}, R_{explained} \rangle \vdash Valid_{Override}$

* Invalid Trace: $\langle D, W_{bias}, \emptyset \rangle \vdash Invalid_{Acceptance}$

Here, *ValidOverride* and *InvalidAcceptance* denote judgment predicates over traces: a valid override occurs when biased warnings are accompanied by a substantive rationale, whereas an invalid acceptance occurs when biased warnings are present but no rationale is recorded.

This grammar allows organizations to parse decision logs programmatically, automatically detecting "ritualized" governance where humans approve risky decisions without derivation. This grammar currently serves as a foundational syntax for auditing. Note that this grammar is intended as a conceptual framework for auditing, not as a fully automated verification system. Future work will expand this into a complete formal system with comprehensive inference rules and soundness proofs.

```
{
  "timestamp": "14:00",
  "gates_passed": ["G1"],
  "warnings_G2": ["Metric_Imbalance_High"],
```

```
  "warnings_G3": ["Freshness_Threshold_Breach"],
  "human_rationale": "Applicant submitted updated
documentation.",
  "final_decision": "Override_Justified",
  "auditable": true
}
```

9. LIMITATIONS

Our analysis relies on simulation parameters adapted from clinical alert-fatigue literature and from stylized accounts of high-stakes administrative decision-making. For tractability, we model compliance decay and intervention probability using a relatively simple structure in which workload, fatigue, and institutional alignment enter in mostly linear ways. In real institutions, however, cognitive fatigue, burnout, and staff attrition can introduce strongly non-linear effects: once workload and frustration cross certain thresholds, compliance can collapse much more abruptly than in our simulations. Similarly, we treat institutional alignment as an exogenous parameter, whereas in practice it is endogenous to power relations, incentives, and labor relations that we do not model. We also evaluate ALTRION only through simulation and conceptual analysis, without field trials or controlled user studies, so the external validity of our quantitative findings remains limited.

Second, ALTRION is scoped to decision regimes in which humans still have both the time and the formal authority to intervene in individual high-stakes cases. We do not claim that the same four-gate pattern is appropriate for ultra-low-latency military command-and-control loops, or for hypothetical deployments of very capable general-purpose AI systems that operate at speeds far beyond human deliberation. In such regimes, different forms of capability control or system-level containment would be required. Our architecture should therefore be understood as a complementary institutional layer for current and near-term AI deployments, rather than as a stand-alone solution to all AI safety problems.

Finally, we abstract away from strategic behavior and gaming. Frontline workers might learn to provide minimal, template-like rationales that formally satisfy Gate 4 without engaging in deep substantive review, and adversarial actors could in principle exploit Gate 2 or Gate 3 by flooding the system with spurious

warnings in ways that induce cognitive overload. We also do not specify concrete user interfaces for presenting Gate-2 and Gate-3 warnings or for collecting Gate-4 rationales; doing so is necessary to evaluate cognitive load, gaming behavior, and adoption in real deployments. Exploring these behaviors empirically or through richer agent-based simulations, and designing additional safeguards to keep ALTRION from becoming either a purely performative ritual or a tool for punitive surveillance, are important directions for future work.

10. CONCLUSION

ALTRION advances AI governance by translating diffuse normative principles into a procedurally enforceable, auditable workflow syntax capable of institutional integration. It offers a workflow-driven framework—shifting the discourse from aspirational guidelines to enforceable procedural safeguards. By modeling compliance decay and formalizing oversight as a grammar, we provide the tools to distinguish between genuine human agency and performative ritualism.

Although public discourse often frames the risks of AI in terms of technological "superintelligence," empirical evidence suggests that the most immediate dangers arise from sociotechnical dynamics: organizational incentives, human overreliance, and the silent drift of decision processes. ALTRION was motivated by the need to protect human agency in precisely these contexts. It is designed not to control AI itself, but to prevent the institutional and cognitive pathways through which AI-induced failures propagate. Ultimately, regulating AI requires regulating the workflows that deploy it. Future research should explore how governance layers such as ALTRION can be operationalized at scale without reinforcing organizational asymmetries or imposing undue burdens on frontline workers.

At the same time, the very mechanisms that make ALTRION effective as a safeguard can be repurposed as instruments of surveillance and control. Structured friction, mandatory rationales, and dense logging can either protect human agency or intensify managerial oversight, depending on how logs are accessed, aggregated, and linked to sanctions or rewards. Our analysis therefore reinforces that institutional design around ALTRION – including worker protections, collective representation, and legal constraints on log

usage – is as important as the technical architecture itself. More broadly, we hope that treating governance as an explicit workflow architecture will encourage future work that jointly designs model objectives, institutional incentives, and human–AI interfaces rather than optimizing any one of these layers in isolation.

Ethics Statement

This research proposes a governance architecture designed to protect human agency. However, we acknowledge that any system introducing logging (ADR) carries the risk of being repurposed for worker surveillance. To mitigate this, we emphasize that the primary function of ALTRION is to audit institutional processes, not to penalize individual operators. Specifically, ADRs should be aggregated into institutional reports rather than used for individual performance tracking. We have used synthetic data for all simulations to avoid privacy risks associated with real-world sensitive data.

REFERENCES

- [1] Parasuraman, R. and Manzey, D.H. 2010. Complacency and Bias in Human Interaction with Automation: A Taxonomy and Review. *Human Factors*. 52, 3 (Jun. 2010), 381–410.
- [2] Buçinca, Z., Malaya, M.B. and Gajos, K.Z. 2021. To Trust or to Think: Cognitive Forcing Functions Can Reduce Overreliance on AI in AI-assisted Decision-making. *Proceedings of the ACM on Human-Computer Interaction*. 5, CSCW1 (Apr. 2021), 1–21.
- [3] Almada, M. 2019. Human intervention in automated decision-making: Toward the construction of contestable systems. *Proceedings of the 17th International Conference on Artificial Intelligence and Law (ICAIL '19)*. Association for Computing Machinery, 2–11.
- [4] Selbst, A.D., boyd, d., Friedler, S.A., Venkatasubramanian, S. and Vertesi, J. 2019. Fairness and Abstraction in Sociotechnical Systems. *Proceedings of the 2019 Conference on Fairness, Accountability, and Transparency (FAccT '19)*. Association for Computing Machinery, 59–68.
- [5] Raji, I.D., Smart, A., White, R.N., Mitchell, M., Gebru, T., Hutchinson, B., Smith-Loud, J., Theron, D. and Barnes, P. 2020. Closing the AI Accountability Gap: Defining an End-to-End Framework for Internal Algorithmic Auditing. *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (FAccT '20)*. Association for Computing Machinery, 33–44.
- [6] Ananny, M. and Crawford, K. 2018. Seeing without knowing: Limitations of the transparency ideal and its application to

algorithmic accountability. *New Media & Society*. 20, 3 (Mar. 2018), 973–989.

[7] Barocas, S. and Selbst, A.D. 2016. Big Data’s Disparate Impact. *California Law Review*. 104, 3 (2016), 671–732.

[8] Stumpf, S., Rajaram, V., Li, L., Wong, W.-K., Burnett, M., Dietterich, T., Sullivan, E. and Herlocker, J. 2009. Interacting meaningfully with machine learning systems: Three experiments. *International Journal of Human-Computer Studies*. 67, 8 (Aug. 2009), 639–662.

[9] Metcalf, J., Moss, E. and boyd, d. 2019. Owning Ethics: Corporate Logics, Silicon Valley, and the Institutionalization of Ethics. *Social Research: An International Quarterly*. 86, 2 (2019), 449–476.

[10] Elish, M.C. 2019. Moral Crumple Zones: Cautionary Tales in Human-Robot Interaction. *Engaging Science, Technology, and Society*. 5 (Mar. 2019), 40–60.

[11] Selbst, A.D. 2020. Negligence and AI’s Human Users. *Boston University Law Review*. 100, 4 (2020), 1315–1376.

[12] Phansalkar, S., Edworthy, J., Hellier, E., Seger, D.L., Schedlbauer, A., Avery, A.J. and Bates, D.W. 2010. A review of human factors principles for the design and implementation of medication safety alerts in clinical information systems. *Journal of the American Medical Informatics Association*. 17, 5 (Sep. 2010), 493–501.

[13] Green, B. 2022. The flaws of policies requiring human oversight of government algorithms. *Computer Law & Security Review*. 45 (Jul. 2022), 105681.

[14] Tacihagh, A. 2023. Governance of generative AI. *Policy and Society*. 42, 3 (2023), 445–463.