**⟡ ChatGPT**

# ClipCard (Risk & Recheck): A Triggered, Blame-Safe Add-On for High-Risk Decisions Across Ops Patterns (Kanban, SBAR, ADR, OODA)

## 1. Abstract

**Background:** High-risk operational decisions—whether in software deployments, clinical handoffs, or policy actions—often bypass thorough scrutiny due to time pressure or complacency. Traditional safety checks (e.g. static checklists) can suffer from poor compliance and false security [1] [2]. This paper proposes **ClipCard**, a lightweight intervention triggered *only* when risk is high, to bolster decision rigor without burdening routine workflow. **Intervention:** A ClipCard is a one-page add-on attached to an existing task (Kanban ticket, SBAR report, Architecture Decision Record, OODA loop step) when **Impact × Uncertainty ≥ 18 (on 1–5 scales)** or when a safety/legal/irreversible red flag is present. The card mandates structured fields: a specific **Hazard** definition ("if X then Y fails [in context Z]"), a **Recheck plan** (with date or conditional "jump"), a named **Recheck Steward** responsible, an **Authority Window** limiting initial scope (e.g. rollout ≤10%), a **Two-Key approval or TTL (time-to-live)** for high-impact actions, explicit **Kill Criteria** (trigger threshold and abort action), and **Evidence Links** (snapshotted references). **Methods:** We outline an evaluation via a stepped-wedge cluster trial and time-series analysis across domains (DevOps, clinical operations, Trust & Safety), measuring outcomes like false approval rate (post-decision incidents), on-time rechecks, near-miss capture, and cycle time impact. **Anticipated Findings:** We hypothesize ClipCard will reduce costly reversions/incidents (false approvals) while ensuring ≥80% of scheduled rechecks occur on time and catching more near-misses, all without materially slowing low-risk throughput. **Conclusion:** If supported by empirical results, ClipCard offers a falsifiable, cross-domain safety mechanism that integrates into existing processes to improve high-risk decision outcomes in a blameless, systematic fashion.

## 2. Introduction

High-risk decisions in operational settings too often rely on ad-hoc judgment or incomplete checks. From software **Site Reliability Engineering (SRE)** to **clinical care**, teams face a dilemma: strict process safeguards improve safety, but applied universally they can bog down work and be circumvented. For example, the well-known **WHO surgical safety checklist** has saved lives, yet real-world compliance is only ~54% on some critical steps [1]. Researchers ask pointedly: if checklists are skipped nearly half the time, *do they truly make us safer*, or might patchy use create a *false sense of security* [3]? Indeed, studies have noted that when a checklist is treated as a mere symbolic formality (posted on a wall or inconsistently remembered), it becomes a "weak safety barrier" **prone to bypass** [4] [2]. Simply put, obvious safety protocols often fail because people adapt and cut corners under pressure.

Two dynamic forces underlie this problem. First, **normalization of deviance** can set in: if parts of a checklist or process are seen as low-value, people gradually omit them without immediate consequence, eroding the safety margin [2]. Second, the **Efficiency–Thoroughness Trade-off (ETTO)** constantly pressures teams to get work done faster at the expense of rigor [5]. As Hollnagel and others observe, organizations and

individuals "continually adjust between efficiency and thoroughness," and systems inherently "strive to be as efficient as possible" [6]. Over time, incremental workarounds to save effort can push operations to the **boundary of acceptable performance** [7] [8]. In such conditions, safety steps that are always required (for every change or patient or case) may be viewed as obstacles and skipped, especially if they rarely catch problems. Ironically, adding more routine rules can backfire—each new rule creates another opportunity for non-compliance, and widespread rule violations can become normalized if the rules are seen as overly hindering [9] [10]. In combination, a flawed safety routine plus the omission of other checks (on the assumption "the checklist covers it") can induce a **new risk**: the organization believes it is safe when in reality it has hollow safeguards [11].

**ClipCard** is introduced in this context as a novel solution: it aims to deliver thorough scrutiny *only when and where it's needed*. Rather than apply a heavy checklist to every task, ClipCard is a **triggered, micro-scale intervention** attached to a task or decision item when its risk warrants extra caution. ClipCard's design explicitly addresses why obvious checks fail. By triggering on a quantified risk threshold (e.g. **Impact × Uncertainty ≥ 18** on 5-point scales), it focuses attention on the riskiest 5–10% of cases, avoiding fatigue from over-use. By "clipping on" to existing workflows (like adding a card or form section in Kanban, SBAR, etc.), it preserves the familiar process for low-risk work, thus *not* inviting daily workarounds. And critically, ClipCard fosters a **blame-safe, learning-oriented culture**: it emphasizes capturing hazards and mitigations without assigning personal fault, aligning with *just culture* principles where errors are reported and addressed systemically rather than punished [12] [13]. We posit that this approach can achieve the best of both worlds – the **rigor of checklists** when needed, without their downsides when not.

**Intervention Overview:** A ClipCard is a brief template (on paper or digital) that is attached to a work item only under high-risk conditions. Each ClipCard requires the author or team to fill in six mandatory elements (with concise prompts): (1) a specific **Hazard** statement capturing the worst-case *if–then* failure mode ("If \<trigger/event> occurs, then \<outcome> fails in \<context>"), (2) a planned **Recheck or Jump** condition – either a date/time for follow-up or a conditional trigger (a "jump") that will prompt re-evaluation, (3) a designated **Recheck Steward** (name or role) responsible for ensuring the recheck happens, (4) an **Authority Window** limiting the initial impact scope or duration (for example, *only deploy to 10% of users* or *operate for 1 hour and then require confirmation*), (5) a **Two-Key approval or Time-To-Live (TTL)** mechanism for any irreversible or high-impact action – meaning either two independent approvers must turn "keys" to authorize, and/or the action has a TTL after which it automatically requires re-approval or shuts down, and (6) explicit **Kill Criteria**, i.e. predefined observable conditions that, if met, will trigger an immediate stop or rollback of the action. Additionally, the ClipCard includes an **Evidence Links** section to attach or reference supporting data – these could be snapshots of logs, analysis results, or prior incidents, archived in a stable form so that decisions can later be audited against what was known at the time. All of this is designed to fit in a single page or card. Figure 1 illustrates the ClipCard components and how a high-risk task flows with the ClipCard process (using a swimlane diagram of roles).

**Contextual Integration:** ClipCard is not a standalone approval gate but an *overlay* to existing operational patterns. It is designed to clip onto common frameworks: for software and DevOps teams using **Kanban** boards or issue trackers, a ClipCard can be attached to a ticket tagged high-risk, adding fields like "Hazard" and "Recheck due date" to the ticket. For architecture or engineering decisions, which often use **Architecture Decision Records (ADR)**, a ClipCard section can be appended to the ADR markdown when a decision has high uncertainty or impact, ensuring future revisits of that decision (ADR processes already value documentation of decisions and their context [14]; ClipCard augments this with risk management fields). In clinical or quality operations, where the **SBAR** (Situation-Background-Assessment-

Recommendation) format is standard for escalating critical issues [15] [16], a ClipCard would "clip" onto the SBAR communication for cases that are especially high-stakes (e.g. a complex patient transfer or policy decision), adding a clear hazard forecast and follow-up plan beyond the immediate recommendation. Even strategic and crisis decision processes like the **OODA loop** (Observe–Orient–Decide–Act) [17] can host ClipCard as a formal check between Decide and Act: when an OODA cycle involves a high-impact action (say, deploying a quick fix in a cyber incident or making a rapid policy call), a ClipCard is triggered to slow down just enough to double-check the key risks and exit criteria before taking action (Figure 2 shows examples of how ClipCards attach in Kanban, SBAR, ADR, and OODA contexts).

By integrating with these workflows, ClipCard aims to **keep low-risk flow light**—if Impact×Uncertainty is below the 18 threshold and no red flags are raised, teams proceed as usual (no extra form to fill, no approval hoops). But once the threshold is triggered, the presence of a ClipCard imposes a *ritualized pause*: the team must explicitly document what could go wrong, establish who will follow up, and put guardrails on the action. This not only surfaces critical thinking in the moment, but also creates an artifact for retrospective learning. Importantly, the ClipCard approach is **blame-safe**: the role of **Recheck Steward** is one of oversight and facilitation, not personal accountability for the outcome. If a hazard manifests despite the ClipCard, the focus remains on improving the process or criteria, not blaming the steward or author. This is akin to the "blameless post-mortem" culture in SRE, where detailed accounts of failures are encouraged without fear, to enable learning and prevention [12]. In fact, by having a named steward and a planned recheck, ClipCard institutionalizes the kind of follow-through that often only happens in post-incident reviews, thereby potentially catching issues in *near-miss* mode instead of after full failure.

**Contributions:** This paper defines the ClipCard concept and situates it relative to existing safety and reliability practices. We describe the theoretical rationale (drawing on safety science and human factors research) for why a triggered add-on might outperform always-on checklists. We then detail the ClipCard format and propose how to test it rigorously in real operations. Key contributions include: (a) a formal specification of ClipCard fields and triggering criteria (making the intervention reproducible), (b) an experimental design strategy (using **stepped-wedge cluster trials** and alternative A/B or Interrupted Time Series approaches) to evaluate ClipCard's effectiveness in different domains, (c) a set of operational **metrics** (false approvals, recheck rates, near-misses, etc.) and targets to quantify success, and (d) an open discussion of risks, biases, and implementation challenges (to ensure falsifiability – we explicitly consider how ClipCard could fail or be misused). Throughout, we maintain a neutral, empirical tone: ClipCard is presented as a testable intervention, not a panacea or silver bullet.

Before diving into related work, we state our core hypotheses for evaluation of ClipCard in practice:

- **H1 (False Approvals Reduction):** Implementing ClipCard for high-risk decisions will *reduce the rate of false approvals* – i.e. decisions that proceed but later result in costly incident or reversion – compared to baseline processes without ClipCard. We expect a measurable drop in post-decision incidents when ClipCards are used.
- **H2 (Recheck Compliance):** ClipCard will ensure follow-through on planned reviews, achieving at least **80% on-time rechecks** for flagged items (meaning the scheduled or conditional re-evaluations happen by the assigned time or trigger, at a far higher rate than ad hoc follow-ups in the control condition).
- **H3 (Near-Miss Capture and Efficiency):** Use of ClipCard will *increase the capture of near-misses* – situations where a potential failure is caught and corrected before harm – by prompting hazard identification and kill criteria. We target an **escalation/stop rate of ~10–25%** of ClipCard-tagged

actions (i.e. in about one out of 4–10 high-risk cases, the ClipCard process triggers an escalation or halt that avoids a bigger failure). At the same time, the median time to fill out a ClipCard (the overhead) is expected to be ≦6 minutes, keeping the burden low.

- **H4 (Minimal Impact on Low-Risk Flow):** ClipCard will *not materially slow down* work on decisions that remain low-risk. In other words, the cycle time for routine, non-escalated tasks or decisions will remain approximately the same with ClipCard in place, since those tasks seldom invoke a ClipCard. This hypothesis checks that the presence of the ClipCard mechanism doesn't create drag on overall throughput by over-triggering or encouraging unnecessary caution.

In the following sections, we review related safety frameworks (Section 3), explain the ClipCard methodology and planned evaluation in detail (Sections 4 and 5), address implementation and organizational factors (Sections 5 and 6), discuss ethical and validity considerations (Sections 7 and 8), and outline limitations and future directions (Sections 9 and 10). Figures and tables are included to illustrate key concepts (ClipCard flow, integration examples, authority window timeline) and to define metrics and checklists. Appendices provide concrete artifacts (templates, policy snippets, audit checklists, data schema, analysis plan pseudocode) to facilitate adoption or further research. By the end of the paper, practitioners should understand when ClipCard is appropriate, how to deploy it, and how to objectively determine if it improves safety and reliability in their domain.

## 3. Related Work

We situate ClipCard in the context of existing practices across various fields:

**Structured Communication Frameworks (SBAR):** The healthcare domain's SBAR format – *Situation, Background, Assessment, Recommendation* – exemplifies a simple checklist for information exchange during critical conversations [16] . SBAR ensures that, for example, a nurse calling a physician about a deteriorating patient systematically covers what is happening, pertinent history, their assessment, and what they need. ClipCard complements SBAR by adding a *risk-specific extension*: SBAR conveys the immediate clinical picture, while a ClipCard (triggered if the case is especially high-risk or complex) would attach details like a predicted hazard (e.g. "if labs not reviewed, patient could have X complication"), a plan to recheck the patient or labs after an intervention, and clear criteria for escalation (e.g. if no improvement in 1 hour, transfer to ICU). Unlike SBAR, which is a communication tool, ClipCard is a *decision accountability tool*. However, both aim to foster a culture of safety – SBAR by standardizing critical communication [15] , and ClipCard by ensuring critical thinking steps aren't missed when stakes are high. A key difference is **triggering**: SBAR is used for every critical conversation, whereas ClipCard is used *only* for the subset of cases meeting risk criteria, thus reducing cognitive overload in routine situations.

**Operational Review and "Stop-the-Line" Practices:** In Lean manufacturing and high-reliability organizations, workers are empowered to **stop the line** if they detect a serious quality or safety problem [18] . For instance, Toyota's production system has an *andon cord* that anyone can pull to halt production when a defect is found [19] . This approach ensures immediate attention to issues but can be disruptive if misused. ClipCard can be seen as a nuanced "stop-the-line" mechanism: instead of halting all operations, attaching a ClipCard effectively *pauses that particular decision process* for a focused scrutiny. It's less drastic – other work can continue in parallel – yet it introduces a checkpoint for the risky item. There are also parallels in **DevOps**: many organizations implement hold points or manual approval gates in deployment pipelines for production changes. A common pattern is requiring a second human approval for production pushes (a basic two-person rule) or limiting deployments to a small percentage (canary releases) initially.

ClipCard formalizes and combines these patterns (two-key approvals, authority windows for safe rollout percentages) in a domain-agnostic way. The difference is that ClipCard includes the forward-looking recheck: not only do we stop or slow the line to fix an issue, we also schedule to *come back later* and see if any latent issue is emerging. This ritualized recheck goes beyond many stop-the-line implementations, which fix the immediate problem but may not always follow up on downstream effects.

**Checklists and Quality Control:** The **Surgical Safety Checklist** and **aviation pre-flight checklists** are iconic safety tools. They demonstrate that simple checklists can significantly reduce errors (e.g., the WHO surgical checklist reduced post-op complications and mortality in trials). However, as discussed in the Introduction, checklist compliance in practice can wane, and their effectiveness depends on consistent use [20]. Researchers Rydenfält *et al.* (2014) characterize checklists as sometimes "symbolic" safety barriers that can be bypassed or become rote [21] [22]. ClipCard differs in being *triggered* and **adaptive**: it's not a fixed list for every operation, but a dynamic add-on when risk is above threshold. In essence, it is a *contingent checklist* – invoked by a risk metric or red flag. This is akin to the concept of "conditional pause" in some industries: e.g., NASA has "hold points" in launch sequences only if certain anomaly flags trip. Additionally, ClipCard's content is not a one-size-fits-all list of items to check (like "verify patient ID, mark surgical site" on a surgery checklist); instead, it is *content generated by the team for that scenario*. Each ClipCard is context-specific (the hazard and kill switch for a code deployment will look very different from that for a policy decision), which encourages active engagement rather than ticking through a generic list. In that sense, ClipCard draws from **cognitive forcing functions** – tools that force decision-makers to articulate their reasoning and contingency plans, thus reducing biases. It turns implicit assumptions ("this change *should* be fine") into explicit statements ("if my assumption is wrong, here's how we'll know and what we'll do").

**Two-Person and High-Authority Controls:** The idea of requiring multiple approvals or limiting single-person authority for critical actions is well-established. The **two-person rule** (also called two-man rule) is famously used in nuclear launch control: two officers must concur to authenticate a launch, to prevent a single rogue actor from causing catastrophe [23]. In technology and finance, dual control is recommended for high-risk transactions (for example, code changes to production at certain companies need a second approver, or large financial transfers require two authorizations). ClipCard incorporates this principle as the "Two-Key/TTL" element for high-impact actions. What's novel is combining it with a **TTL (Time-To-Live)** – effectively an expiration on the approval. For instance, a ClipCard might state that two separate leaders must approve a new policy rollout, and that approval is only valid for 24 hours; if the change isn't executed in that window, approvals reset. This prevents "authorization creep" where a plan approved under certain conditions is executed much later under different conditions without re-evaluation. We illustrate this in Figure 3 as a timeline: one person's approval (first key) starts a ticking clock (TTL); if the second key isn't turned in time or if conditions change, the plan must be revisited. The **authority window** concept similarly relates to known practices like **canary releases** or **blast radius limits** in software deployment – initially only affect a small subset (say 5–10% of users or systems) and only proceed to 100% if no issues are detected [24]. Many companies like Facebook, Google, etc., have internal tools to do gradual rollouts with automatic rollback triggers at, say, 5% error increase. ClipCard formalizes the expectation that any high-risk change defines such a window (e.g., "we will only expose 10% initially and auto-revert if error rate >X%"). In regulated contexts, **time-limited operations** are common (consider a pharmaceutical emergency use authorization that expires in a year if not renewed). ClipCard's authority window ensures that high-risk actions are not open-ended; it embeds a bias for reversible, testable increments.

**Adaptive Decision and Incident Frameworks:** The **OODA loop** (Observe–Orient–Decide–Act) from military strategy has influenced business and incident response by emphasizing agility and continuous feedback [17]

[25] . One lesson from OODA is that decisions are iterative and an entity can gain advantage by cycling faster and learning. ClipCard resonates with OODA in that it inserts a feedback checkpoint in the cycle. Specifically, after "Act," the ClipCard's recheck/jump forces an explicit Observation & Orientation on the results of that action at a set time, rather than assuming the loop will naturally be revisited. It's a way to guard against a common failure mode: once a decision is implemented, teams often move on and may miss signals that the decision is not producing the expected outcome (until a disaster happens). By calendaring a recheck or by linking it to an automatic alert (e.g., "if error rate goes above threshold, reconsider"), ClipCard closes the OODA loop proactively. In **Site Reliability Engineering (SRE)** and IT incident management, a related concept is the **blameless post-mortem** and continuous improvement cycle [26] [12] . ClipCard can be viewed as pushing some of that post-mortem reflection *to the pre-mortem phase* – essentially performing a pre-implementation risk review and scheduling a mini-post-mortem (the recheck) whether or not an incident occurs. This aligns with approaches like **Chaos Engineering** which intentionally trigger failures to learn system limits – except here, we anticipate possible failures and define the kill switch in advance.

In summary, ClipCard's novelty is not in inventing new safety concepts per se, but in **combining and operationalizing** them across domains in a triggered, lightweight form. It builds on the success of checklists, the caution of two-key rules, the empowerment of stop-the-line, and the learning focus of post-mortems. However, unlike broad safety programs or cultural mandates, it provides a concrete artifact and process that can be experimentally evaluated. Next, we describe in detail how ClipCard is structured and the research methods to test its efficacy.

# 4. Methods

This section describes the ClipCard intervention specification and outlines how we plan to evaluate it rigorously. We cover the ClipCard format and usage rules, proposed study designs (with one or two primary approaches and a backup design), the operational contexts in which we will trial ClipCard, the outcome measures and performance targets, and the analysis plan including statistical methods. We also discuss considerations for statistical power and sample size in these pragmatic trials.

## 4.1 Intervention Specification: ClipCard Format and Use

Each ClipCard is a structured mini-document (physical card or digital form) with required fields capturing the risk mitigation plan for a specific decision or action. The mandatory fields, as introduced earlier, are:

- **Hazard:** A concise, specific statement of the primary hazard or failure mode we are worried about. This should be phrased as "**If X, then Y fails [in Z]**" to force clarity. For example, *"If the new database deployment script fails on step 4, then customer data may become corrupted in the accounts table."* The hazard identifies the *triggering condition* (X), the *failure outcome* (Y), and optionally the scope or where it applies (Z). An effective hazard is concrete and testable (avoiding vague language). This is not merely a risk rating but a narrative of how things could go wrong.
- **Recheck/Jumps:** A plan for follow-up. This can take two forms (often both): a **dated recheck** (e.g. "Recheck results on 2025-11-01 by 17:00") and/or a **conditional jump** trigger (e.g. "If error monitoring alert `ALERT_DB_CORRUPTION` fires, do an immediate recheck"). The term *Jump* refers to a jump condition that might prompt earlier re-evaluation than the scheduled time. Essentially, the team commits to *when* and *under what conditions* they will revisit this decision to see if everything is as expected. Multiple conditions can be combined, but our guidelines recommend having *one composite jump condition* for simplicity (e.g., the OR of several alerts could be one jump trigger).

- **Recheck Steward:** The person (or role) responsible for carrying out the recheck and coordinating any further action. This is not necessarily the person who made the decision; it could be a separate reliability engineer, a team lead, or an on-call responsible individual. Naming a steward establishes clear ownership for the follow-up. We also encourage listing a backup steward in case the primary is unavailable. The steward is expected to ensure the calendar reminder or alert is in place, and that results of the recheck are noted.
- **Authority Window:** The predefined limits on the authority or scope of this action before additional review is needed. This typically includes a **scope limit** (e.g., affect ≤10% of users, process ≤100 records, run in a sandbox environment first) and/or a **time window** (e.g., the change will automatically revert or pause after 2 hours unless extended). The authority window serves two purposes: it reduces potential blast radius (if something goes wrong, the impact is constrained), and it creates a natural pause point for the recheck (e.g., if we set a 2-hour window, we know our recheck should happen by then to decide on full rollout). In practice, implementing an authority window might involve technical means (like a feature flag auto-disable, or a temporary access key that expires) or procedural means (like requiring a second check-in meeting before going from pilot to full scale).
- **Two-Key/TTL for High-Impact:** If the action is high-impact (by magnitude or irreversibility), the ClipCard should enforce a **two-key approval**—meaning two independent decision-makers must approve the action. For example, both the engineering manager and the SRE lead must turn their "keys" (which could be literal signatures or digital approvals) to allow a production deployment that could cause downtime. In addition, a **TTL (Time-To-Live)** on the approval may be specified: e.g., "Both approvals are required and expire after 24 hours." This ensures that if execution is delayed, fresh judgment is applied rather than relying on stale approval. Two-key approval is a classic risk control to prevent single-point human error [23], and when combined with TTL it adds temporal control. If an action is extremely urgent (e.g., emergency fix), TTL alone might be used to auto-revert unless continued approval is given. In summary, this field formalizes whether special authorization mechanics apply. If a two-key is required, the ClipCard lists who the two approvers are. If a TTL applies, it states the time limit or condition (e.g., "Approval valid for 60 minutes once given" or "Session will auto-logout after 10 minutes unless second approval entered").
- **Kill Criteria:** Predefined conditions under which the operation will be aborted or rolled back, along with the specific kill action. This is essentially an if-then for emergency termination. For example: "**Kill Criterion:** If more than 1% of transactions error out OR any data corruption is detected in monitoring, **Action:** immediately rollback the deployment to version v1.2 within 15 minutes." Kill criteria should be *observable and measurable* triggers (thresholds, specific alert signals, etc.), not subjective feelings. The paired action should be something that can be executed quickly (like initiating a rollback, calling a code red, reverting a policy decision, etc.). The presence of kill criteria is crucial: it forces the team to think "what would definitely tell us this has gone wrong?" ahead of time, and commit to not riding out a bad situation. This draws on the idea from experimental design where you set stopping rules to avoid causing harm.
- **Evidence Links:** References to evidence used in making the decision and evaluating risk, captured in a stable form. For instance, links to test results, simulation outputs, past incident reports, or policy documents that were considered. These links should ideally point to *snapshotted* or versioned resources (so later when someone looks back, they see what information was available at decision time, even if the source data changed). If linking to dynamic sources, one might attach a PDF or include a brief summary on the card. The purpose is transparency and auditability: if a ClipCard decision is later reviewed (in an audit or post-mortem), one can trace back the rationale and data

behind it. It's also helpful during the recheck – the steward can quickly revisit these links to verify assumptions (e.g., re-run a query, check if a cited document's context still holds).

**ClipCard Trigger Rule:** The default trigger for requiring a ClipCard is **Impact × Uncertainty ≥ 18** (on 1–5 scales for each, where 5 is highest). This risk scoring is inspired by common risk matrices in project management and engineering. For example, if a change is assessed as very high impact (5) and medium uncertainty (4), the product is 20, triggering a ClipCard. Conversely, even a highly uncertain change (5 uncertainty) with moderate impact (3 impact) is 15, which falls below 18 and might not trigger (unless a red flag is present). We choose 18 as a threshold to capture roughly the top 10–15% of risk scenarios – it requires at least one of the factors to be near max and the other at least mid-high (e.g. 5×4, or 4×5, or 3×5 =15 which is just below threshold, so actually one factor has to be 5 and the other ≥4, or both 5 and 5). This threshold can be tuned per organization (some may lower to 15 or raise to 20 depending on how risk-averse they are or how frequently they want ClipCards). Additionally, **automatic triggers** override the numeric score for certain categories: if an item raises a **safety, legal, regulatory, public trust, or irreversible consequence** flag, a ClipCard is mandated regardless of the score. For instance, a seemingly low-impact change that touches patient privacy might be flagged for legal/public risk and hence get a ClipCard by policy. These triggers ensure no critical issue slips through just because someone underestimated impact or uncertainty subjectively.

**Use and Workflow:** When a task or decision is identified as meeting the ClipCard trigger, the responsible person (e.g., the engineer, analyst, or clinician bringing up the change) starts a ClipCard. Practically, this might be done by clicking a "Add ClipCard" button in a tracking system which opens the template fields to fill, or by pulling a physical card and writing in pen if working analog. The key is that it becomes *required* to fill all fields before proceeding further in the decision process. The filled ClipCard is then attached or linked to the task. Depending on the environment, there may be gating: e.g., in a software deployment pipeline, a script might check that if risk score >=18, a ClipCard file is present and signed-off, otherwise deployment is blocked. In a hospital setting, a policy might be that a procedure with a ClipCard cannot be started until a supervisor has reviewed the ClipCard entries. These gates are minimal just to enforce usage when triggered.

Once the ClipCard is active, execution proceeds but within the constraints of the card (the authority window and kill criteria are effectively "in force"). The Recheck Steward makes sure any needed timers or alerts are set (for example, scheduling a calendar event for the recheck, or configuring a monitoring alert to ping if kill criteria are met – some organizations might integrate this into tooling so that filling a ClipCard auto-sets an alarm). The team carries out the decision or change. If a **kill condition** trips or a problem is noticed, they execute the kill action immediately (and later, in the post-review, credit the ClipCard for having pre-defined that threshold). At the designated recheck time (or jump trigger), the steward and relevant team members convene (even if briefly) to evaluate the outcome: Did the hazard materialize at all? Is everything within safe parameters? If all is well, they document that the recheck passed and the item can be fully closed out (and for example, the pilot can continue to full rollout beyond the authority window). If issues are found, either the kill action is taken or further mitigation steps are decided. Either way, the ClipCard then serves as an artifact in the incident repository or decision log.

Importantly, ClipCards should be logged in a way that the data can be aggregated for analysis (this is covered in Appendix D: a minimal CSV schema for logging ClipCard usage and outcomes). That data includes timestamps, risk scores, whether kill criteria were triggered, etc., to enable measuring our outcome metrics.

## 4.2 Study Design Options

To evaluate ClipCard's effectiveness, we consider experimental and quasi-experimental designs that balance rigor with feasibility in operational environments. Implementing a randomized controlled trial at the individual decision level is challenging (one cannot easily randomize *within* a team which risky decisions get a safety intervention without causing unfairness or confusion). Instead, we focus on **cluster-level and time-sequenced designs** where different groups or time periods adopt ClipCard versus control.

**Primary Design A – Stepped-Wedge Cluster Randomized Trial:** We propose a stepped-wedge design as our primary evaluation method. In a **Stepped-Wedge Cluster Randomized Trial (SW-CRT)**, all participating clusters (e.g., teams or departments) will eventually receive the intervention, but the order and timing of rollout are randomized [27] . We start with a baseline period where no cluster is using ClipCard (all using their normal decision processes). Then at predefined intervals (say, every 2 weeks or 1 month), a random subset of clusters crosses over to using ClipCard. This continues until all clusters have switched to the ClipCard condition. Throughout the study, data on outcomes are collected continuously for each cluster. By the end, every cluster has contributed both control data (before they got ClipCard) and intervention data (after adoption), at different times. This design is powerful for implementation research because it ensures every group gets the potentially beneficial intervention (important for ethical reasons if we have belief in ClipCard's value), and it allows within-group comparisons (each cluster serves as its own control in a sense). We will randomize the sequence of rollout to avoid bias (e.g., not always starting with the most enthusiastic team, etc.). We will also include a small number of **"hold-out" clusters** if possible, that adopt the intervention last, to serve as concurrent controls in the early phases.

The stepped-wedge approach is depicted schematically in Table 2 (though Table 2 in our list will detail gate implementation recipes). We anticipate using, for example, 6 clusters (teams) switching over in 6 steps (so study length might be 6 intervals plus baseline). Each cluster's outcomes (like false approval incidents per week, etc.) pre- and post-adoption will be analyzed with segmented regression or mixed-effects models (see Analysis section). The stepped-wedge design assumes the intervention effect is either immediate or gradual but one-directional once a cluster adopts (we won't withdraw it). We will have to watch for temporal trends (if overall risk or incident rates are changing over time due to other factors) – the analysis will adjust for underlying time trends.

**Primary Design B – Alternating Interventions (AB/BA crossover by time):** If clustering by team is not feasible (for example, if there are not many distinct teams or if work is organized by weekly sprints rather than stable teams), an alternative is a time-based alternating intervention. For instance, in a tech ops environment, we could do an A/B test by *week*: during **"odd" weeks**, any high-risk decisions use ClipCard, during **"even" weeks** they use the normal process (no ClipCard). This would be announced and enforced as an experiment. Over, say, 10 weeks, we'd have 5 weeks with ClipCard on, 5 off, alternating. This is a form of **repeated measures crossover**: each timeframe acts as either control or intervention, and teams likely experience both conditions repeatedly. The advantage is it controls for time-invariant differences (the same team, same type of work, just different weeks). We'd have to be careful about carryover effects – e.g., if a ClipCard created in an "on" week schedules a recheck in an "off" week, we need protocols for that. Possibly we treat each decision's process according to the week it *started*. We'd likely introduce a short washout or training period as well. The alternating-week design makes sense in environments with continuous homogeneous work. It might be less suited if decisions are rare or if work intensity fluctuates widely week to week (randomization could alleviate some of that). We could also randomize at a finer grain (day by day) but that might be too chaotic for participants; weekly is a reasonable cadence teams can plan around.

**Backup Design – Interrupted Time Series (ITS):** If randomization is infeasible (e.g., leadership decides to roll this out to everyone at once due to safety concerns), we will rely on an **Interrupted Time Series** analysis with segmented regression [28] [29]. This involves collecting outcome data for several periods before ClipCard implementation (baseline trend) and after implementation (post-intervention trend) and comparing the level and slope changes. For example, we'd measure the rate of false-approval incidents per month for 6 months pre-ClipCard and 6 months post-ClipCard introduction. The **segmented regression model** would estimate any immediate drop in incident rate (change in intercept at the "interruption") and any change in trend (slope difference before vs after). ITS is a robust quasi-experimental design especially suitable if the data is at an aggregate level and the intervention timing is known [30] [31]. We would strengthen it by checking for possible confounders that change around the same time and by possibly using a control series if available (e.g., if one similar department did not implement ClipCard, use their incident trend as a comparison). This is inherently less strong than randomization, but still can provide evidence of effect, especially if the effect size is large relative to variability and if no other major changes coincide with the intervention.

We may combine designs: for instance, do a stepped-wedge across multiple organizations (e.g., across 3 hospitals or 3 companies), and also analyze each with ITS. Or if the rollout has to be sequential for logistical reasons, stepped-wedge naturally occurs. If we have enough clusters, we might consider a parallel cluster randomized trial (half teams start using ClipCard, half not, then compare), but operationally we suspect a stepped-wedge is more palatable because everyone gets it eventually.

## 4.3 Contexts (Settings) for Evaluation

We plan to trial ClipCard in at least two distinct contexts to test its generality:

- **Software Reliability / DevOps (Kanban & ADR):** In a software engineering organization (such as an IT ops department or a SaaS product team), using agile Kanban boards and Architecture Decision Records. High-risk items here could include production changes, infrastructure configuration updates, or architectural decisions with big impact on reliability or security. We'll integrate ClipCard into tools like Jira or GitHub issues – for example, a Jira workflow condition that if a ticket is marked Impact=High and Uncertainty=High, a ClipCard must be attached before moving to "Done". ADRs for major design choices will have a section to fill if the design is flagged high risk (say novel technology adoption with uncertainty). In SRE incident management, if someone is about to apply a fix that might have side effects, they would do a quick ClipCard. This context will test ClipCard's effect on incidents (outages, rollbacks) and on the development velocity.
- **Clinical Operations / Healthcare (SBAR handoff & critical procedures):** In a hospital or clinic setting, focusing on processes like patient transfers, critical care decisions, or invasive procedures. We'll tie ClipCard to SBAR reports or procedural checklists. For instance, if during patient rounds a situation scores high on risk (maybe a complex case with uncertain diagnosis and potential for rapid decline), the team creates a ClipCard outlining what to watch for (hazard), who will follow up (steward, e.g., charge nurse), and what the threshold is to call a rapid response (kill criteria). We would measure outcomes like "near misses" (e.g., catching patient deterioration early), compliance with follow-up (did that patient get the planned recheck labs on time?), and adverse events or code blue rates. Because healthcare already has safety protocols, we'd ensure ClipCard fits as a supplement, not a replacement (perhaps in parallel with existing checklists).
- **Trust & Safety / Policy Decision-making:** In online platforms or corporate policy groups, decisions about content moderation, user safety actions, or policy changes can be high-risk (affecting public

perception, legal exposure, irreversible outcomes like banning a user or releasing a controversial feature). We'll pilot ClipCard in a Trust & Safety team where, for example, an unusual account suspension (impacting a high-profile user) triggers a ClipCard to enumerate potential fallout ("if our detection is wrong, we could erroneously suspend a politician – public outcry"), set a recheck (e.g., review case in 24 hours with a second team), and define kill criteria (e.g., "if >10% of flagged accounts in this sweep appeal successfully, pause the policy rollout"). Metrics here might include false positive/negative rates of enforcement, user harm incidents, etc., as well as operational metrics like how quickly issues are corrected.

- **(Optional) Security or Civic Domain:** We may also explore ClipCard in a security operations center (for incident response decisions) or even in civic/NGO project management for decisions that impact communities. These would broaden the evaluation but are optional if resources permit.

Each context will generate data for the defined outcomes, and we can perform context-specific analysis as well as pooled analysis (possibly via meta-analysis) to see if ClipCard's effects are consistent. The **multi-context approach** addresses external validity – ClipCard is meant to be domain-agnostic, so it's important to test that.

## 4.4 Outcomes and Metrics

We define clear outcome metrics to test the hypotheses (summarized in Table 1 with operational definitions):

**Primary Outcomes:**

- **False Approvals Rate:** This is the frequency of decisions that were approved/executed but later found to be wrong or harmful, requiring significant reversion or causing an incident. In software, this could be measured as the proportion of changes that had to be rolled back or caused a post-deployment incident (severity 1 or 2) within, say, 7 days of deployment. In healthcare, it might be procedures or discharges after which the patient had an adverse event that was preventable. We will specifically track *high-risk decisions* and see how many result in a bad outcome in ClipCard vs non-ClipCard scenarios. We expect ClipCard to lower this. This metric is akin to a "failure rate post-approval." It will be binary per decision (0 = no incident, 1 = incident) and we'll aggregate rates or use survival analysis (time to incident) if appropriate.
- **On-Time Rechecks:** Among the ClipCard-tagged decisions (or analogous high-risk decisions in control), did the planned follow-up happen on schedule? For ClipCards, this is straightforward: each card has a recheck due (time or condition), and we log whether it was completed by that time. For control cases (without ClipCard), we may only know if any follow-up occurred by examining logs or asking teams (which is harder). This metric basically measures adherence to the follow-through. On-time recheck is binary per decision (yes/no), and we'll compute the percentage. Our target is ≥80% with ClipCard. We'll also look at distribution of delay for those that were late.
- **Near-Miss Capture:** We will count "near-misses" – cases where a potential failure was averted due to the process. Concretely, a near-miss in ClipCard context could be identified by either the kill criteria triggering (which implies we caught a condition and took action *before* it became a bigger incident) or by the recheck uncovering an issue that was then mitigated. For example, if a ClipCard's hazard was "if X then Y fails" and at recheck we find early signs of X happening and fix it, that's a near-miss capture. In control, near-misses might be harder to quantify (likely underreported), but we may use proxy like any incident that was narrowly avoided or quickly fixed due to someone's vigilance. We

expect ClipCard to increase reported/observed near-misses, because it sets traps to catch them. We'll measure near-misses per 100 high-risk decisions or similar.

- **Time Overhead (Median Fill Time):** We will measure how long it takes to fill out a ClipCard, as an indicator of added friction. This can be tracked via timestamps (when a card is triggered vs when it's marked complete or approved). We aim for the median to be ≤6 minutes. This ensures ClipCard is indeed a "tiny" intervention, not a half-day risk assessment meeting. We'll also note variation; if some cases take 30 minutes, maybe complexity or confusion occurred – that's instructive.
- **Escalation Rate:** This refers to how often a ClipCard leads to an escalation or higher-level review. For example, requiring two-key approval inherently escalates to another approver; or a recheck might escalate an issue to management or specialist. We target an **escalation range of ~10–25%** of ClipCard cases – a healthy rate suggesting it's catching things that need more attention, but not so high as to overwhelm leadership or indicate over-triggering. If <10%, perhaps ClipCard is too stringent (most are false alarms), and >25% might mean we are applying it too late (almost every high-risk needs escalation, maybe threshold should be lower to catch earlier, or teams are deferring too much to ClipCard).
- **Low-Risk Cycle Time:** As a secondary outcome tied to H4, we'll monitor the overall cycle time or throughput for tasks that are not ClipCard-tagged, to ensure they remain steady. For example, average lead time of normal tasks in Kanban, or number of patients processed per week in a clinic. We expect no significant difference pre vs post ClipCard introduction in those metrics, indicating normal work wasn't slowed. If we see a slow-down, it could mean cognitive drag or indirect effects (like maybe teams become more cautious broadly, which could be good or bad).

**Secondary Outcomes and Process Metrics:**

- **ClipCard Utilization Rate:** How often is the intervention triggered as a fraction of all decisions? This checks that our threshold ~18 is hitting the expected fraction. If we find far more or fewer ClipCards than expected, that suggests calibration issues.
- **Quality of ClipCard Entries:** This isn't a direct outcome but for internal validity we'll audit a sample of ClipCards (say 10% randomly) for completeness and specificity (Appendix C's audit checklist provides criteria). Metrics could include: % of cards with a clearly measurable kill criterion, % that identified a non-generic hazard, etc. We might use this to ensure no significant drift or Goodhart's Law effect (people creating junk data to satisfy the requirement).
- **Team Satisfaction and Safety Climate:** We will gather qualitative feedback or survey data on how the team perceives ClipCard (does it make them feel safer, or burdened?). This is secondary but important for adoption likelihood. We might use a Likert scale survey for those who used ClipCard vs those who didn't, regarding their confidence in decisions, psychological safety, etc.
- **Error Types and Near-miss Types:** We will categorize what kinds of hazards were caught (e.g., design flaw, oversight, communication gap) to see *how* ClipCards are providing value. If, for example, many ClipCards catch communication issues (like someone forgot to tell X something and the recheck catches it), that helps refine training.

Table 1 in the Figures & Tables section will list these outcomes with precise definitions (e.g., how exactly we define "false approval" in each context, how we calculate on-time recheck percentage, etc.).

## 4.5 Analysis Plan

We outline the analysis approaches for the chosen study designs, focusing on estimating effect sizes with confidence intervals and hypothesis testing, while being mindful of the correlated nature of the data (since

decisions within a team may not be independent, etc.). We also describe how we will present results in an "overhead vs avoided loss" framework to make the value case.

**Effect Size Estimation:** The primary effect of interest (H1) is the reduction in false approval rate. We will calculate the difference in this rate between ClipCard and control conditions. For instance, if baseline false approvals are 10% of high-risk changes and with ClipCard it is 4%, that is an absolute reduction of 6 percentage points (or a relative risk ~0.4). We will report relative risk or odds ratio with 95% confidence intervals, as well as absolute differences. Given likely low frequencies, we might use a **Generalized Linear Model (GLM)** with a logit link (logistic regression) or even a mixed-effects logistic model if clustering is used, to estimate these differences. For each outcome, we'll tailor the model: e.g., for on-time recheck (binary), logistic; for near-miss count per team per month (count), possibly Poisson or negative binomial GLM (with cluster as random effect if needed); for fill time (continuous, likely skewed), we might use a median comparison or log-transform and use linear mixed model.

**Stepped-Wedge Analysis:** In a stepped-wedge CRT, the standard analysis is a mixed-effects model with fixed effects for time (to adjust for secular trend) and a fixed effect for intervention (ClipCard vs control), plus random effects for cluster. We will likely use a generalized linear mixed model (GLMM) appropriate for each outcome (e.g., logistic for incident occurrence). There are specialized methods for stepped-wedge (including Hussey & Hughes approach) that account for the multiple time points [27] . We will also perform **segmented regression** on the aggregated time series per cluster as a complementary analysis: essentially treating each cluster's data as an ITS and possibly combining via meta-analysis or pooled regression with interaction terms. Segmented regression in this context will model an intercept change and slope change once ClipCard is introduced in each cluster [28] [31] . Autocorrelation in time series will be checked (e.g., Durbin-Watson or using robust SEs; if needed, AR(1) terms can be included or use generalized estimating equations with a correlation structure).

**Alternating Weeks Analysis:** If we do the AB crossover by weeks, we can use simpler paired comparisons. One approach is to aggregate data by week (since that's the unit of intervention) – e.g., count incidents or near-misses that week – and then compare the distribution of those metrics between ClipCard weeks vs control weeks using paired statistical tests (each pair being two consecutive weeks or a more advanced method using time series regression with a sine wave for the alternating pattern). However, a more granular approach is also possible: include all individual decisions in a logistic regression with a predictor for "ClipCard active or not" and maybe a term for time trend or week number. Because the same team sees alternating conditions, we might include a random effect for week pair or use GEE (Generalized Estimating Equations) clustering by week or sprint to account for correlation. If the outcomes are measured per decision, we can cluster by week or by team as needed.

**Interrupted Time Series:** For an ITS design, we will use **segmented regression** as described. The model for an outcome $Y_t$ (like monthly incident count or rate) would be:

$$Y_t = \beta_0 + \beta_1 t + \beta_2 I(t \geq t_0) + \beta_3 (t - t_0)_+ + \epsilon_t,$$

where $t_0$ is the time of ClipCard introduction, $I(t \ge t_0)$ is an indicator for post-intervention period (so $\beta_2$ is the level change), and $(t - t_0)_+$ is time elapsed since intervention (0 before $t_0$, linear after; $\beta_3$ is slope change) [31] . We will check $\beta_2$ and $\beta_3$ for significance to see immediate and gradual effects. If multiple sites or units, we could extend this to a multi-site ITS with random slopes (to allow each site a different baseline and trend) and use a fixed intervention effect across

sites (or random effect if we think effect may differ by site). If data points are few and overdispersion is a concern, exact tests (like Fisher's exact for binary aggregated outcomes pre vs post) might be employed as a simpler check, though segmented regression is preferred for utilizing all data points [32].

**Statistical Tests:** Alongside effect size estimates, we will use two-tailed significance tests (likely $\alpha = 0.05$) for key hypotheses. For H1, that might be a Z-test on the logistic regression coefficient for ClipCard vs control. For H2 (80% on-time recheck target), we might test if the rate is significantly different from some baseline or from control. H3's near-miss increase and median fill time can be assessed by comparing distributions (maybe Wilcoxon rank-sum for fill time if not normal). However, given our focus on falsifiability, we prioritize confidence intervals and practical significance over p-values. A result like "false approval rate 10% → 4%, OR = 0.38, 95% CI [0.20, 0.72]" speaks to effect magnitude and uncertainty.

**Overhead vs Avoided Loss:** We will explicitly analyze the "cost" of ClipCard in terms of time and any delays, versus the "benefit" in avoided incidents or rework. This could involve computing something like: total hours spent filling ClipCards + extra approval time, etc., and comparing to estimated hours of downtime or rework saved by prevented incidents. For example, if each ClipCard took 5 minutes (0.083 hours) and 100 ClipCards were made, that's ~8.3 hours of effort. If it prevented 2 major incidents that would have each caused 4 hours of downtime and 10 hours of investigation, that's 28 hours saved, netting ~20 hours benefit. These back-of-envelope calculations help communicate ROI to decision-makers. We might present a table or graph of "Additional effort vs. avoided incident cost" under various assumptions. If available, we'll also use any financial or safety-critical impact metrics (like dollars saved, or harm prevented, etc.) to make a compelling case. However, we'll also note intangible benefits (improved knowledge sharing, confidence) that are harder to quantify but valuable.

**Meta-Analysis (if multi-site):** If we have multiple independent sites (e.g., different companies or hospitals) running similar experiments, we will analyze each separately and also do a **random-effects meta-analysis** on the effect sizes. For instance, each site yields a risk reduction estimate with CI, and we combine them to see the overall mean effect and heterogeneity. A Cochran's Q or I^2 might be reported to indicate variation across contexts. If one context shows no effect and another a strong effect, that points to context dependency and will be discussed.

**Interim Monitoring:** Although this is primarily a research study, if ClipCard appears to be causing harm or overwhelming burden, we would ethically need to adjust. We might plan an interim analysis after X months or after first Y ClipCards to check key metrics (particularly incident rates and recheck compliance). This will not be used with strict stopping rules like a medical trial, but as a safety check for the study itself.

**Data Visualization:** We will use run charts and time series plots to visually show trends. For example, a plot of monthly incident rate with a dashed line when ClipCard started, showing before/after trend (with segmented regression fit). Or cumulative count of near-misses captured over time in intervention vs control. Figures may include these to illustrate results.

## 4.6 Power and Sample Size Considerations

Because this is a pragmatic field trial, sample size is partly constrained by available events and duration. However, we attempt to estimate if we have enough power to detect hypothesized effects. For instance, say baseline false approval incident rate is 10% of high-risk decisions. We aim to detect a reduction to 5%. Using a two-sided test with $\alpha=0.05$, power 80%, for a simple comparison of proportions, we'd need on the

order of ~200 high-risk decisions in each group (control and ClipCard) to detect that difference (this is a rough estimate via power calculation for two proportions). In a stepped-wedge, each cluster contributes multiple observations; the calculation is more complex due to correlation within cluster and repeated measures. We can use formulas or simulations for stepped-wedge design. Typically, fewer clusters (say 6-8) can be compensated by more time points and events per cluster. We might simulate based on expected incidents per week etc. If not enough events occur, we might extend the study duration or include more clusters.

For recheck compliance (H2), if baseline follow-up rate is, say, 50% and ClipCard raises it to 90%, that difference is huge and easy to detect even with small N (maybe 30 decisions would show that with high significance). For near-misses (H3), since near-miss capture might be nearly zero in control (if not actively tracked) and maybe >0 with ClipCard, even a qualitative comparison might suffice. However, measuring near-miss reliably might need a definition and possibly requiring self-report in control which could undercount – so we acknowledge a limitation there.

We also consider the risk of **false positives** due to multiple outcomes – however our primary outcome is clearly false approvals. Others are secondary, so we won't adjust alpha for each, but we will interpret them in context (and possibly use a Bonferroni or Holm correction in a formal sense if needed).

We will provide guidelines for organizations on roughly how many decisions or how long they should run ClipCard trials to see significant results. For example, if incidents are rare (say 1 in 50 high-risk decisions), one might need on the order of 200–300 decisions observed to see a statistically significant halving of that rate. We'll encourage sites to monitor the metrics over time and possibly pool data with others for stronger inference.

Finally, we note that statistical significance is not the sole aim; even if some results are not "significant" due to low sample, the magnitude and confidence intervals will inform whether ClipCard is promising or not. If the point estimate shows improvement but CI is wide (due to low N), it suggests more data collection is needed. Conversely, if it shows negligible improvement with narrow CI, that falsifies our hypothesis robustly.

We now turn to practical considerations of implementation (Section 5), having laid out how the intervention is defined and will be rigorously tested.

## 5. Implementation Notes

Implementing ClipCard in real operations requires thoughtful integration into existing tools and routines with minimal disruption. In this section, we provide guidance and patterns for deployment, including where to insert ClipCard in various process flows, how to enforce it with minimal "gates," and how to leverage existing systems (calendars, version control, alerting) to support the ClipCard's functions. We also describe lightweight solutions to common implementation challenges (like ensuring evidence snapshots and scheduling rechecks).

**Placement Patterns in Workflow:** ClipCard should attach at the juncture between planning and execution – essentially at the "decision commit" point of a process. We highlight placement in four archetypal workflows: - *Kanban (Software Task Management):* When using a Kanban board or issue tracker, a natural

point is when a ticket is about to move from the "In Review" or "Ready for deploy" column to "Done." Before that final move, if the ticket is labeled high-risk (by our trigger criteria), a ClipCard must be completed. In Jira, one could implement this via a workflow transition screen or required fields: e.g., if a custom field "RiskScore" $\geq$ 18, then fields "Hazard, Steward, Recheck date, etc." become required to resolve the issue. The ClipCard content can live as a sub-task or attached form. Figure 2A (in our conceptual diagrams) might show a Kanban card with a red ClipCard icon indicating an attached risk card. This placement ensures normal low-risk tasks flow freely, but risky ones get an extra step right before completion/handoff. - *SBAR (Clinical Handoff):* In an SBAR communication, typically a nurse or resident composes the SBAR to call a senior doctor. If a case triggers a ClipCard, the SBAR form (often a template or structured page) would have an added section for ClipCard items after the Recommendation. Alternatively, one could initiate a separate "ClipCard form" linked to the patient case. The placement is right after the initial assessment/ recommendation – essentially saying: "We think we should do X (Recommendation), and since this is high risk, here's our ClipCard specifying how we'll double-check and who's watching." The attending physician might be required to sign off the ClipCard items as well, analogous to two-key for critical clinical decisions (like giving a high-risk medication). - *ADR (Architecture Decision Record):* When writing an ADR for a major design or infrastructure change, teams typically record context, decision, alternatives, etc. We advise adding a section titled "Risk & Recheck (ClipCard)" within the ADR if the decision meets the risk trigger. This section would enumerate the hazard ("if our assumption about scaling is wrong, the system could fail under load"), the steward ("e.g., SRE team will re-test 1 month after launch"), kill criteria ("if CPU goes above 80% on average in any week, roll back feature toggle"), etc. By embedding in the ADR, it stays as part of the permanent decision log. The placement is at decision documentation time – often ADRs are reviewed by an architecture council, and they could require that high-risk ADRs have a completed ClipCard section before approval. - *OODA/Incident Response:* In fast OODA cycles (like emergency management or military decisions), implementing a ClipCard requires discipline to pause between Decide and Act when needed. One practical way: during incident calls or war-room scenarios, designate someone as "Recheck Steward" at the start if the impact/uncertainty is high. That person will fill a quick ClipCard (could even be on a whiteboard visible to all: Hazard, Kill criteria, etc.) before the action phase. Alternatively, some incident management tools could prompt "Is this action high-risk? If yes, fill quick risk card." The ClipCard could be literally a sticky note on the ops dashboard listing what condition will stop the action. For example, in a cyber incident: Decide = disconnect a certain server from network; before acting, fill ClipCard: hazard "if wrong server, we cut off hospital's system," steward = incident commander, recheck in 10 minutes if attack stops, kill criteria "if any critical service fails after disconnect, plug back in within 30s." This is extreme rapid use, but shows adaptability.

**Minimal Gates and Automation:** We strive to enforce ClipCard usage in triggered cases without adding heavy bureaucracy. Some *lightweight gating recipes* (to be summarized in Table 2) include: - *Jira "Recheck due" field:* Configure a Jira or Trello automation that if a ticket is high-risk (perhaps using labels or a numeric field), it cannot be closed unless a "Recheck Due Date" and "Steward" fields are set (which implies a ClipCard exists). This is a soft gate; it makes sure those using the system fill in something. The presence of those fields effectively signals an attached ClipCard. - *CI/CD Pipeline Rule:* In deployment pipelines (e.g., Jenkins, GitLab CI, Azure DevOps), include a manual approval gate or script that checks for a ClipCard artifact. For instance, a pipeline could call an API or check a repository for a file named `CLIPCARD-<change_id>.json` with required fields completed. If not found or not approved by a second person, the pipeline fails before production deploy. This leverages existing **checks and approvals features** (like Azure DevOps checks [33] or AWS CodePipeline manual approvals). - *Handoff Roll-call:* In shift changes or handoffs (common in healthcare or 24/7 operations), include ClipCards in the handoff checklist. E.g., outgoing lead says: "These 3 ClipCards are active – we are awaiting recheck on X at 3 PM, Y has kill criteria

monitoring." Incoming lead acknowledges. This ensures continuity, so that ClipCards with future recheck times aren't forgotten at shift change. - *Alert-as-Jump integration:* Many monitoring systems (Datadog, CloudWatch, etc.) allow webhooks or custom actions when an alert triggers. We can integrate so that if a specific alert (tied to a ClipCard's kill criteria or jump condition) fires, it automatically pages the Recheck Steward or creates a task. Essentially, codify the ClipCard's jump into the alerting system. This might involve tagging alerts with an ID that maps to the ClipCard. Tools could even post a message like "Kill criteria threshold met for ClipCard #123: hazard XYZ – take action!" - *Calendar Wiring:* For scheduled rechecks, we strongly suggest using a shared team calendar or reminder system. Ideally, filling "Recheck Date/Time" on a ClipCard automatically sends a calendar invite to the steward (and perhaps relevant team members) at that time with the hazard context. If using Outlook/Google Calendar APIs, this can be automated. Otherwise, it's on the steward to set a reminder. But procedural policy can enforce it: e.g., "When you fill a ClipCard, also send a meeting invite for the recheck to the team alias." - *Evidence Snapshotting:* To avoid reliance on mutable data, teams can utilize version control or attachment features. For example, if linking to a query result or log file, one might attach a CSV export or screenshot to the ClipCard. If the ClipCard is digital, allow file attachments. For code or config states, referencing a specific git commit hash ensures we know what was deployed. In an ADR, one might freeze the referenced document version. We can provide a small script or guideline: "for important evidence links, consider appending `?version=<id>` or using archive services." Even printing a PDF of a web page into an evidence repository might be done for critical references (like a regulatory guideline page). - *Versioning ClipCards:* If a ClipCard needs updates (say during the process someone refines the kill criteria), how to handle version? A simple approach: treat the ClipCard as a living document until the task is closed, with edits tracked (if digital) or marked (if paper). After closure, it's archived as final version. If we needed to revert to an earlier plan, presumably it's recorded as changes in the card. This likely doesn't need a heavy system but worth noting.

**Training and Onboarding:** We recommend a brief training for teams adopting ClipCard. This can be a 1-hour workshop explaining why it exists (perhaps referencing cases where "lack of recheck" caused failure), walking through a sample ClipCard, and clarifying that it's not meant to assign blame but to catch issues. Role-play or simulation (maybe doing a mock ClipCard for a hypothetical scenario) can help teams become comfortable. Emphasize especially how to write a good Hazard and Kill Criteria, as these may be new skills (some might initially write vague hazards like "System might break" – training should push for specificity: "If memory usage >90% then response times degrade causing outage in region"). Appendix C's audit checklist can double as a quick guideline during training ("check your hazard is specific, kill is observable," etc.).

**Scaling and Tool Support:** For large organizations, implementing ClipCard might be aided by tool support: possibly an add-on to existing ticket systems or a form template. If none exists, starting with a simple document template (Markdown or Google Form) is fine. Over time, one could integrate it into risk management systems. The benefit of hooking into existing systems is ease of adoption and data capture (for evaluation). In the absence of fancy tools, even a shared Excel/Sheets listing ClipCards or a wiki page per ClipCard could work. The key is that it's accessible and reviewable.

**One-Page Cap Discipline:** We mentioned guardrails to avoid process creep: one such guardrail is keeping ClipCards to a single page or a few minutes to fill. Managers and safety officers should reinforce this: if people start creating 5-page risk assessments for each ClipCard, it defeats the purpose. Perhaps have a rule like "no more than 3 bullets under Hazard, 1 line each; kill criteria should be at most 2 conditions," etc. The idea is to encourage concise critical thinking, not exhaustive analysis (which there usually isn't time for in fast ops). The audit process (random checks of ClipCards, as in Section 6) will also ensure they remain short and focused.

**Notification Routing:** Who gets notified of what in ClipCard usage? Perhaps the Recheck Steward is always notified of triggers. If kill criteria triggers an alert, perhaps the on-call and steward and maybe a tech lead are all notified. Some organizations may want a summary of all active ClipCards accessible (like a dashboard). We can implement a simple dashboard listing active ClipCards, their next recheck times, and statuses (maybe integrated in project management tool or just a wiki page someone updates). This can be reviewed in weekly reliability meetings to see patterns (e.g., "we have 5 ClipCards currently open, 2 were closed last week after safe outcomes, 1 triggered a rollback").

In short, implementing ClipCard leverages many existing capabilities (task fields, pipeline approvals, on-call alerts, etc.) stitched together with a bit of process. By planning these minimal gates and wiring, we ensure ClipCards are neither forgotten nor onerous.

## 6. Risk, Bias & Validity

No intervention is without risks or potential biases. In adopting ClipCard, organizations must be vigilant about unintended consequences, and our evaluation must address validity threats.

**Goodhart's Law and Ritual Compliance:** There is a danger that teams treat ClipCard as just another bureaucratic requirement to satisfy ("fill the card to get the deploy through") without truly engaging in the spirit of hazard analysis. If metrics like "% of tasks with ClipCards" or "on-time recheck rate" are emphasized, people might game them (e.g., setting a trivial recheck time and checking a box without real analysis). This is a classic Goodhart's law scenario – when a measure becomes a target, it can cease to be a good measure. To counter this, we plan **quality audits** of a subset of ClipCards. For instance, for every 10 ClipCards, a safety officer or peer reviewer will audit 1 (10%) using a checklist: Is the hazard specific and meaningful? Is the kill criterion measurable and relevant? Was the steward actually a person who could act? Did they attach evidence? Each audit is pass/fail on those criteria. This can discourage perfunctory ClipCards, as teams know there's random checking. The audit results (anonymized) can be fed back as training: e.g., if many cards fail the "hazard specificity" check, we do a refresher on writing hazards. We do *not* tie audit fails to individual punishment; rather it's team-level coaching (to avoid blame). In summary, measure the measures to ensure they're not gamed.

**Selection Bias and Champion Effect:** If participation in ClipCard usage is voluntary, it's likely that teams with a safety mindset or better organization will use it more, and thus naturally have fewer incidents (even without ClipCard's effect). This could bias results (overestimating effect). In our study design, we mitigate this by using randomization (stepped-wedge ensures all clusters eventually use it, and order is random) and by measuring baseline performance of each cluster. We'll statistically adjust for any baseline differences. However, if adoption is uneven or if certain types of decisions get ClipCards while others (equally risky) do not due to human judgment, that's a selection bias. We explicitly instruct that the trigger rule is to be followed consistently (to avoid managers only sometimes applying it when they feel like it). In analysis, we can compare the risk profiles of ClipCard-tagged vs non-tagged decisions to check for any systematic differences beyond the threshold (maybe some borderline cases that should have had it but didn't, etc.). If found, we might do sensitivity analysis (e.g., excluding borderline cases or including an instrumental variable for risk score).

**Team vs Individual Attribution:** A validity (and ethical) point is that we measure many outcomes at team level (incidents, etc.), and improvements might be due to broader team learning, not just ClipCard. Conversely, a particularly diligent individual might drive ClipCard use and improvements, which might not

generalize. To handle this, we emphasize **team-level reporting** and discourage singling out individuals. For example, if one team's incident rate drops, we attribute to the process (ClipCard + team effort), not "Alice did ClipCards well." In analysis, the cluster (team) is our unit, which inherently pools individual behaviors. We avoid comparing individuals using ClipCards vs not, which would be confounded by personal skill. This also maintains blamelessness: we don't want individuals to feel an incident is a "failure of their ClipCard" – it's a team process failure if anything. In our retrospective reviews, we examine what could improve in the ClipCard process or thresholds, not "who wrote a bad card."

**Process Creep and Alert Fatigue:** There's a risk that ClipCard, if initially successful, starts getting applied to more and more scenarios (scope creep) or that kill criteria alerts become so frequent that they desensitize staff (alert fatigue). To guard against process creep, we set clear bounds: the ClipCard template is one page, and the trigger threshold defines its intended frequency. If we find that e.g. 50% of tasks are getting ClipCards, the threshold might be too low and we should adjust upward to restore selectivity. We also emphasize *threshold tuning*: if initial data shows that very few ClipCards actually escalate or catch issues, perhaps threshold could be slightly lowered to capture more (assuming capacity to handle more). Or if kill criteria rarely trigger because they were set too conservatively, we adjust guidance on setting meaningful thresholds. Essentially, treat the threshold and criteria as parameters to refine, not holy writ. On alert fatigue: if every ClipCard sets an alert and many are false alarms, that can degrade trust. Thus, writing good kill criteria (specific enough to not trigger too often on noise) is critical. Also, using multi-signal criteria can reduce noise (e.g., require two triggers to fire). We might implement a rule that if kill criteria alert fires and on investigation it was a false alarm, the team must refine either the threshold or the monitoring to be more accurate – continuous improvement on the signals. In evaluation, we will look at how often kill alerts fired vs actual issues to estimate false positive rate.

**Metric and Measurement Error:** Some outcomes (like "incident severity" or "near-miss") may involve subjective judgment in classification. If those judging are aware of intervention, bias could occur (observer bias). We will attempt to use objective measures where possible (like a rollback event either happened or not). For near-misses, we rely on reporting – we anticipate better reporting with ClipCard (which is partly the point), but that complicates comparing to control. The control might underreport near-misses, making ClipCard look beneficial partly just by surfacing information. We acknowledge this limitation: part of ClipCard's aim is indeed to reveal things that would have been hidden, so an increase in near-miss count is a success but also a measure of improved reporting. When interpreting near-miss data, we'll be careful: if near-misses go up, it could mean things that would have been minor incidents without ClipCard are now caught (good), or just that we now have a mechanism to count them. We might supplement this with qualitatively assessing if corresponding actual incidents go down, which would support that these near-misses were meaningful catches.

**Hawthorne Effect:** It's possible that teams knowing they are part of a safety experiment might improve their overall practices (with or without ClipCard). The presence of measurement can motivate temporary performance changes. To minimize this, we embed ClipCard as naturally as possible and conduct the study over a sufficient duration that novelty wears off. The alternating design (if used) helps because people can't slack off in control weeks lest something goes wrong, and they might be generally on alert. If anything, this could dilute measured effect (because control is better than usual), giving a conservative estimate.

**Context Validity and Domain Differences:** ClipCard may not be equally effective everywhere. For example, in extremely fast-paced environments (ER medicine, live site incident within minutes), filling a card might be impractical under urgent pressure. Our evaluation in multiple contexts will reveal where it works or not. It

could turn out that in domain X it prevents many issues, but in domain Y teams find it cumbersome and bypass it. Such differences are important validity considerations; they suggest moderators to the effect. We will analyze subgroups to see if effect size differs by context or by risk type. If significant, we refine the scope: maybe ClipCard is only recommended for certain domains or risk types where it proved its worth. Or adapt the intervention (maybe an ultra-light ClipCard variant for time-critical situations, like a verbal or mental ClipCard).

**Data Privacy and Sensitivity:** ClipCards, especially evidence links and narratives, might contain sensitive information (PII, confidential strategy, patient data). There's risk in how that's stored and who can access. This is more of an implementation risk than an evaluation bias, but it affects adoption. We advise that ClipCards be kept within the same access-controlled system as the work item. For example, if an ADR is internal to engineering, the ClipCard content stays in that doc; if a hospital uses a secure checklist app for SBAR, ClipCards reside there. For any analysis or cross-team data sharing, we will anonymize and aggregate data. The evaluation logging schema (Appendix D) is designed with minimal identifying info. One must be cautious if linking actual evidence snapshots – e.g., a database dump as evidence should not go in a public repo. Classification of information and following privacy guidelines (especially in clinical settings with HIPAA) is necessary. Our project will go through IRB or equivalent ethics review for any collection of sensitive operational data, ensuring compliance.

In summary, we proactively consider and mitigate biases: enforcing quality to combat superficial compliance, using rigorous design to handle selection effects, monitoring process load to avoid burnout, and keeping focus on learning over blame. These measures enhance the validity of both the intervention and the study results.

## 7. Ethics

Implementing an intervention like ClipCard and studying its effects involves ethical considerations on multiple levels: organizational ethics (avoiding a blame culture, respecting employee time and privacy) and research ethics (especially if human factors data or potentially patient-related data are collected).

Firstly, ClipCard is explicitly designed to be **blame-safe and fair**. From an ethical standpoint, we want to ensure it is used to improve system outcomes, not to assign liability to individuals. In training and policy (Appendix B provides a snippet), we clarify that the **Recheck Steward role is not a scapegoat** – being the steward does not mean if something goes wrong, it's "their fault." Instead, it's a role to facilitate safety. We encourage **role rotation** for stewards when feasible to distribute the responsibility and learning. For instance, different team members can take turns being steward for different ClipCards, so no one person is always "on the hook." This not only avoids burnout but also prevents a situation where, say, one cautious individual ends up steward for everything (which could inadvertently make them a target if something slips). By rotating, we emphasize team ownership of risk.

**Just Culture and Learning:** We embed ClipCard usage into a **Just Culture** approach [34] [13]. That means if a failure occurs even with a ClipCard in place, the analysis asks "Was our process sufficient? Were the conditions tricky? What can we improve?" rather than "Who screwed up filling the card?" If someone deliberately ignored a known kill criterion and caused harm, that might be an accountability issue, but ideally ClipCard's process-oriented nature guides people away from such negligent behavior. We will obtain organizational buy-in that data from ClipCards (like hazard predictions) will not be used to punish predictors

if something unforeseen happens. On the contrary, it should be praised that they attempted a prediction at all.

**Informed Participation:** In our study, we will inform all participants (teams, managers) about the trial's purpose and that their performance data will be collected, but in aggregate and for improvement. We'll likely go through an IRB (Institutional Review Board) or equivalent ethics review if, for example, we are doing this in a hospital with patient-related outcomes. Since ClipCard can affect patient safety decisions, any research in that environment must ensure patient care is not compromised. We position the study as a *quality improvement initiative* which often has a somewhat different oversight path than experimental research, but we err on the side of formal IRB if any doubt. We will require consent from staff if we interview or survey them, and we'll anonymize personal data.

**No Harm Principle:** We consider if ClipCard could introduce harm. For instance, could it delay emergency actions and thus harm a patient or system? We mitigate that by excluding extremely time-critical cases or by having an override: e.g., "If taking time for a ClipCard would clearly worsen outcomes (like during cardiac arrest), skip it." ClipCard is meant for high stakes but not when immediate action is needed to save a life in seconds – there's judgment here, and we will encode that in ethics guidelines: safety of people comes first, ClipCard second. In contexts like software, a few minutes delay rarely harms; in medical, it could, so clinicians are trained to use their judgment. The study will monitor for any unintended adverse consequences of ClipCard (like a delay leading to an issue).

**Data Security:** As mentioned, ClipCards may contain sensitive info. Ethically, we ensure that any logs or analysis data stripped of personal identifiers are handled securely. For example, if we log "Incident happened? Y/N" for a patient case, we do not include patient ID or name in our dataset, just maybe a case number that's internal. All data stored for analysis will be on secure servers. If publishing results, we only share aggregated or anonymized quotes.

**Transparency and Accountability:** We commit to sharing the results openly (with participating organizations and potentially publicly) regardless of outcome. If ClipCard does not work or causes slowdowns, ethically we should report that so others don't adopt a flawed practice blindly. Conversely, if it works, we share details so others can benefit. We avoid proprietary black-box approach; ClipCard is not something we "sell" but a practice we are studying collaboratively. This openness aligns with ethical research norms and the broader operations safety community's collaborative spirit.

**Retrospective Learning and IRB:** If doing retrospective analysis of incident records to identify baseline, we might be looking at historical data. Typically, that's allowed under quality improvement with data use agreements, but again we would check with IRB if any patient or user data is involved. Many corporate environments would treat this as internal process improvement not requiring external IRB, but we'd still follow ethical principles (confidentiality, minimal dataset needed for analysis, etc.).

In sum, the ethical approach is to ensure ClipCard *helps and does not hurt*, that participants are respected and not blamed, that data is handled responsibly, and that we remain transparent. If at any point ethical concerns arise (e.g., staff feeling punished by the new process or significant pushback), we will address them through training or modify the process. For example, if someone fears that writing a hazard might later be used against them ("you predicted it yet it happened, so you're liable"), we must actively dispel that and institutionalize that no such blame will occur – in fact, having predicted a hazard at least shows diligence, which should be protected. We might incorporate a policy that all ClipCard contents are protected

under the same rules as incident post-mortems (which in many companies/hospitals are privileged or at least not used for discipline). Appendix B's policy snippet explicitly covers that.

Through these ethical guardrails, we aim for ClipCard to be a positive safety intervention that reinforces a culture of safety and learning, rather than fear or bureaucracy.

# 8. Limitations & Threats to Validity

While we have attempted a thorough design, it's important to acknowledge limitations of ClipCard and of our study, to temper expectations and guide future improvements.

**Adoption Bias & Cultural Factors:** ClipCard's effectiveness likely depends on organizational culture. In a company with an established blameless, safety-first culture, ClipCard might thrive and be taken seriously; in a company with a rushed or blame-heavy culture, people might pencil-whip the cards or resent the extra step. So our results may not generalize to all contexts. Early adopters might be those already inclined to safety (hence their baseline might be better). Conversely, if forced on reluctant teams, they might comply minimally. We try to measure some cultural aspects via surveys, but it's hard to fully adjust for this. Thus, one limitation is that ClipCard as a process tool cannot fix deeper cultural issues; it must be accompanied by leadership support and psychological safety. In our study, sites volunteering are likely ones already motivated to improve safety, which could bias results upward (Hawthorne effect or just good culture effect).

**Context Dependence:** We will test across domains, but maybe our sample is still narrow (e.g., only one hospital, one software org). Differences in regulatory environment (aviation or pharma would have far more formal processes) could make ClipCard redundant or non-compliant with existing regs. For example, in nuclear ops, two-person rule is already mandatory, and adding ClipCard might conflict or add duplication. Or in heavily regulated clinical trials, you might need formal risk management that ClipCard can't replace. Our study doesn't cover all such environments, so we caution that ClipCard might be most useful in *medium-regulated* scenarios (like tech, healthcare ops not research, internal company decisions) and not in already extremely regulated ones (which have their own thorough processes) nor in completely ad hoc ones (where culture may not accept it). Additionally, scale matters: ClipCard is designed per decision, so if decisions are very frequent and high-risk (e.g., stock trading algorithms making micro-decisions – impossible to do ClipCard for each), it doesn't apply there either. We assume a human decision maker and time on the order of minutes to hours available.

**Measuring "Success" is Hard:** Some outcomes (like "incident avoided") are inherently counterfactual. If nothing bad happened after a ClipCard, was it because of ClipCard or it was going to be fine anyway? We can only infer from comparing rates. Even then, incidents are rare, so a statistically significant drop might need a large sample. We might not reach significance for that outcome in all contexts due to rarity. That doesn't necessarily mean ClipCard had no effect; it could mean insufficient data. The study duration is a limitation – if we run for a few months, maybe not enough incidents occur to judge conclusively. We could rely on near-misses as proxy, but near-miss counting has the bias explained. So one has to interpret results with caution. It may take longer-term adoption to truly see reduction in catastrophic failures (like an airline might not see differences in crash rates for years due to rarity). Instead, we focus on intermediate metrics like process adherence which are easier to detect changes in.

**Alert Fatigue and Noise:** We assume kill criteria can be set such that alerts or stops aren't too frequent. If we get it wrong, there's a risk ClipCard generates lots of "false alarms" – cases where kill criteria triggered

but actually it would have been fine. If that happens often, teams might start ignoring ClipCard signals (like crying wolf). We plan to monitor it, but our study timeframe might not capture the long-term fatigue if any. If initial results are good but a year later people start dismissing ClipCard fields as rote, that long-term erosion is beyond our current evaluation window. That's a threat to sustainability. It suggests follow-up studies or continuous monitoring will be needed beyond the initial introduction.

**Potential Slowdown of Innovation:** One could argue that requiring this rigorous process for high-risk endeavors might make teams avoid taking risks at all, potentially stifling innovation. For example, if engineers know a project will be labeled high risk and thus extra scrutiny, they might label it lower risk or not attempt it. This is difficult to measure (the decisions not taken). We have to rely on anecdotal feedback or observing if any high-impact initiatives were perhaps delayed. In our short-term evaluation, we likely won't see that strategic effect, but it's a theoretical concern. We would counter that ClipCard only asks for a small time investment to allow big moves safely, but perception matters. We can at least ask managers if they felt decisions were postponed due to ClipCard overhead. If any evidence of that arises, it's a negative side-effect to address (perhaps by streamlining further or making clear that leadership values the safety step so it shouldn't deter bold actions).

**ClipCard Quality Variability:** The tool is only as good as the content people put in. Some might fill excellent hazard analyses, others might be weak. Variation in ClipCard quality could muddy results: e.g., one site might see no benefit because their ClipCards were poorly done (thus didn't actually catch anything). We audit quality, but still, variability exists. In analysis we might find some hazard fields were empty or kill criteria not well thought out. If we include all data, effect could be diluted by those low-quality implementations. Perhaps a subset analysis on "cards that passed quality audit" vs those that didn't would be insightful. But realistically, not all orgs will do A+ job on it initially. That limitation means ClipCard might need iterative improvement and not guaranteed effective everywhere by default.

**Threats to Internal Validity:** If something else changes during our study, results could be confounded. E.g., if at the same time as ClipCard introduction a new **change management system** or a new **training program** was rolled out, any improvement might be partly or wholly due to that. We'll document any known concurrent initiatives and try to account for them (maybe via qualitative notes or including a covariate if quantifiable). In stepped-wedge, any global time effects (like seasonal variation in incidents or an overall safety trend) should be adjusted by the time fixed effects. But unknown factors always remain a possibility.

**Non-Compliance with the Study Protocol:** Teams might occasionally forget to create a ClipCard when they should (slip through the trigger). Or conversely, might create one even if criteria not met because they felt uneasy (that's fine, but then it's technically outside our trigger rule). Such deviations can affect our intention-to-treat vs per-protocol analysis. We likely will do an "as-observed" analysis (did it have a card or not) because randomization is at cluster level, so intention-to-treat = cluster assigned at time. But any non-compliance reduces effect size – if many high-risk decisions in control clusters still get some similar treatment or vice versa, groups aren't well separated. We mitigate by clear rules and monitoring logs (we can see risk score and whether a card was made to catch misses). But enforcement can't be 100% (especially in busy ops, someone might forget or decide it's borderline and skip it). Thus some dilution in effect is expected. We plan to capture that and possibly do secondary analysis excluding non-compliers (with caution as it could introduce bias).

In acknowledging these limitations, we maintain a falsification-friendly stance: we are prepared for the possibility that ClipCard's effect is smaller than hypothesized or only works under certain conditions. The

goal is to learn *where and when* it helps, and where it does not. By being upfront about these threats to validity, we improve the credibility of results and outline necessary caution for practitioners considering ClipCard.

# 9. Future Work

Our exploration of ClipCard opens several avenues for further development and research beyond the initial implementation and trial:

**Adaptive Thresholds and Machine Learning:** The current trigger threshold (Impact×Uncertainty $\geq$ 18) is a crude rule based on human-assigned scores. In the future, we could employ data-driven methods to refine when to trigger a ClipCard. For example, as we collect data on decisions, their context, and outcomes, we could train a model to predict risk more accurately or flag situations that historically led to incidents even if the subjective scores were low. This could lead to an **adaptive triggering system** that learns domain-specific risk patterns (perhaps lowering threshold for certain types of changes known to be tricky, and raising for others). Machine learning could also assist in suggesting likely hazards or kill criteria based on past ClipCards (like an AI helper that, given a task description, suggests "People in similar changes worried about X"). This moves towards a semi-automated ClipCard generation, reducing burden further.

**Automated Evidence Gathering:** Filling the evidence links could be streamlined by tools. Future iterations might integrate with monitoring and documentation systems so that when you initiate a ClipCard, the system automatically attaches relevant recent data – e.g., latest test results, system metrics, or a link to the last incident report on a similar component. Also, implementing a "snapshot" feature: one click to archive the current version of a document or configuration and link it. This ensures evidence quality and saves time. Research into how to automatically gather and summarize risk-related evidence for a decision would complement ClipCard. For instance, an algorithm could parse code changes and identify modules touched with past incidents (hence suggest evidence: "module X had 2 incidents last year [35] ").

**Integration with Human-in-the-Loop Decision Systems:** As organizations adopt more **automation and AI in decision-making**, we can integrate ClipCard principles there. For example, in a CI/CD pipeline with automated tests, one could automatically flag uncertain deployments (maybe based on test flakiness or lack of coverage) and require a ClipCard. Or consider an AI system that makes content moderation decisions – a ClipCard-like process could be triggered for borderline cases, bringing a human to review with hazard foresight ("If this AI is wrong, what harm?"). Essentially, exploring ClipCard in human-AI collaboration contexts, ensuring that even an AI's "decisions" that are high-risk get a human-governed risk check.

**Hazard Libraries and Domain Templates:** As ClipCards are used, patterns of hazards and kill criteria will emerge. We foresee creating **hazard libraries** per domain. For example, a library for cloud deployments might include common hazards like "If config file format is wrong, service won't start" or "If new feature flags all users, could overload DB" and typical kill switches for them. These can serve as checklists or suggestion lists when someone fills a new ClipCard (similar to how surgical checklists were refined from common failures). There's potential to standardize a core set of hazard scenarios (perhaps 20% of hazards cover 80% of cases in a domain). Teams could then pick from a menu or at least be reminded of them. This reduces reliance on individuals to brainstorm all risks and provides consistency. However, one must avoid just copy-pasting hazards without thinking; the library is a guide, not a substitute for scenario-specific analysis.

**Scaling Up to Program/Portfolio Level:** ClipCard currently addresses individual decisions. Future work could examine aggregating ClipCard insights to inform higher-level risk management. For instance, if over a quarter, 5 ClipCards in a project all cited "if user load > expected, system slows," that indicates a systemic risk (capacity issues) that should be addressed in product planning. So an extension is to roll up hazard frequencies and near-misses to portfolio dashboards. This allows proactive resource allocation (e.g., a pattern of kill criteria triggers might justify a project to improve that area). Essentially, linking ClipCard data to enterprise risk registers or OKR (Objectives & Key Results) processes ensures the learnings lead to structural improvements, not just case-by-case fixes.

**Tooling and User Experience Research:** We intend to refine the user experience of creating and managing ClipCards. Right now it might be a form or a markdown template. HCI (Human-Computer Interaction) research could help design an interface that guides users logically through hazard identification (maybe with prompts or analogies), and that integrates seamlessly with their workflow (e.g., voice input for clinicians who can't type while gloved, or mobile app for on-call engineers awakened at 3am). Also, exploring ways to visualize active ClipCards – perhaps an interactive timeline (like Figure 3) that shows where we are in the authority window and how long until recheck. UX improvements can greatly influence adoption.

**Verification in Additional Domains:** We would like to test ClipCard (or variants) in safety-critical industries like aviation, oil & gas, or construction. Each has its own systems (and many have existing check processes). For example, in aviation maintenance: a ClipCard could be triggered for an unusual repair procedure on a critical component. Or in construction, for an innovative build method that's uncertain. These industries have high stakes and often formal risk assessments (Job Safety Analyses, etc.), but ClipCard's concise triggered approach could complement if carefully integrated. Future collaborative studies with such industries could validate generalizability.

**Policy and Governance Use:** Another future angle is applying ClipCard concept to policy-making and governance decisions, not just operational ones. For instance, a city council making a high-impact decision (like deploying a new technology city-wide) could employ a ClipCard-like process: identify hazard ("if this fails, public trust erodes"), recheck ("public hearing after 6 months"), kill criteria ("if >20% negative feedback in first quarter, pause program"). It formalizes good governance practices. Research could evaluate if that improves outcomes or public satisfaction. However, adopting something like ClipCard in bureaucratic or political contexts might face different challenges (e.g., accountability and blame are politically loaded).

**Automated Compliance and Auditing:** In the future, we could develop tools that automatically audit ClipCards for completeness or even accuracy (to an extent). For example, natural language processing could parse hazard statements to ensure they have an "if…then" structure and maybe flag if it's too vague. Or consistency checks: if hazard mentions a threshold, kill criteria should align. This can assist stewards or quality officers in real-time ("Your hazard doesn't specify where 'fails' – please clarify location [Z]"). These little prompts ensure the card quality. Additionally, audit bots could periodically remind stewards of upcoming rechecks or escalate if a recheck was missed (closing the automation loop).

In summary, ClipCard is a starting point for more responsive and intelligent risk management in operations. Future work will refine triggers (maybe personalized to team risk appetite), integrate more automation (AI suggestions, auto-evidence), and extend the concept to broader decision governance. The ultimate vision is an organizational practice where high-risk decisions are routinely accompanied by a thoughtful plan (like ClipCard) that is in part generated by the system and in part by human judgment, constantly improved

through learning – *Entrogenic* in a sense, adapting across cycles of change and feedback. This would help systems become more resilient and self-correcting as they evolve.

## 10. Conclusion

In high-velocity operational environments, it is easy for critical decisions to slip through with optimistic assumptions and insufficient follow-up. **ClipCard** offers a pragmatic, minimal-friction mechanism to inject a dose of rigor and foresight precisely when it's most needed – and only then. By *triggering on high risk*, it avoids burdening routine work, preserving agility. By standardizing key elements (hazard, recheck, kill switch, etc.), it makes risk mitigation actions explicit and trackable. And by fostering a blame-safe, learning-oriented usage, it aims to improve outcomes without dampening the culture.

Our proposed evaluation suggests that ClipCard, if properly implemented, can reduce the rate of undetected false decisions (those that would later cause incidents), ensure that planned mitigations (like follow-ups) actually occur, and catch more problems early (near-misses), all for a modest investment of time on the part of the team. Notably, ClipCard is not a silver bullet; it will not eliminate all failures, nor is it suited for every situation. It helps most in that middle ground of decisions that are not trivial but not so obviously dangerous that they already receive heavy scrutiny. There will still be cases where unforeseen issues arise outside the scope of any ClipCard's hazard (because we can't predict everything), and conversely cases where ClipCard triggers but everything would have been fine (an acceptable inefficiency if the false-positive rate is kept reasonable). The goal is to tilt the balance: fewer catastrophic surprises, more controlled experiments.

Practically, organizations adopting ClipCard should do so with clear leadership support and alignment with existing processes. When ClipCard helps, it will often be "uneventful" – the lack of incident or the smooth rollback might not make headlines, but those silent saves are the value. When ClipCard is not to be used: if a situation is truly urgent (no time to fill anything) or very well understood and low risk, forcing a ClipCard would indeed slow things needlessly. Thus, knowing when *not* to use it is as important as when to use it. Our threshold and red flag criteria are initial guidance, but teams will refine them. If everything starts needing a ClipCard, the bar is too low; if incidents still sneak through un-carded, maybe the bar is too high or certain risks weren't recognized.

The broader implication of ClipCard is cultural: it reinforces the message that **it's okay to pause and think** even amid urgency, and that doing so is a mark of professionalism, not hesitation. It provides a documented way to say, "We are doing something risky, and here's how we manage that risk," which can be reassuring internally and even to external stakeholders (imagine being able to show an auditor or a client: for every big change, we have this safety card—demonstrating operational maturity). ClipCard's cross-framework nature also encourages cross-pollination of safety practices; an engineer sees how a clinician sets kill criteria for patient observation, or a clinician sees how an SRE uses two-key approvals, and both learn.

In conclusion, ClipCard represents a small intervention with potentially outsized impact, functioning as a safety net for high-impact decisions across diverse operational arenas. Our research is an initial step in evaluating its effectiveness. If the hypotheses hold true, ClipCard could become a recommended best practice in the toolkit of operations management, analogous to how checklists became standard in aviation and medicine—but more targeted and dynamically applied. If results are mixed, we will have learned valuable information about the conditions required for such interventions to work. Either way, the pursuit

of reliable, falsifiable approaches to risk management continues. We invite the operations research and safety communities to scrutinize, replicate, and iterate on these findings. Ultimately, the measure of success is when organizations routinely catch more issues early and recover from them gracefully, with ClipCard or similar mechanisms quietly empowering teams to make safer decisions without sacrificing speed or innovation.

When ClipCard is working, nothing dramatic happens—**and that is exactly the point**.

## Figures and Tables

*Figure 1: ClipCard components and trigger flow.* This swimlane diagram illustrates how a high-risk work item progresses with a ClipCard. The process begins when a task is flagged as high risk (Impact×Uncertainty ≥ 18 or red flag). In the "Owner" lane, the person responsible pauses to initiate a ClipCard, filling Hazard, Recheck, Steward, etc. Meanwhile, in the "Approver/Lead" lane, if two-key approval is required, a second person signs off (represented by parallel approval step). The "System/Process" lane shows automated gates (e.g., a CI pipeline check) that block progress until the ClipCard is completed. Once approved, the task moves to execution within the Authority Window (shown as a bounded region). The "Monitoring" lane shows kill criteria being watched (e.g., an alert set up). Time advances to the Recheck point, where the Steward (in their lane) conducts the recheck. Depending on outcome, either the task is confirmed safe and the ClipCard is closed, or a kill action is triggered (flow looping to a rollback step). The diagram emphasizes coordination between roles and the timeline of triggers.

*Figure 2: Placement examples across Kanban, SBAR, ADR, and OODA.* Four mini-scenarios demonstrate ClipCard integration: - *(a) Kanban:* A screenshot of a Kanban board with one card marked with a red "Risk" label. Clicking it shows a ClipCard form pop-up with fields filled. The card cannot move to "Done" until the form is complete. - *(b) SBAR:* An SBAR report template where a fifth section "R (Risk Recheck)" is added. It contains: Hazard ("Patient may have adverse reaction if…"), Steward (charge nurse), Recheck (vitals recheck in 1 hour), and Kill Criteria (e.g., "if blood pressure drops below X, call Rapid Response"). - *(c) ADR:* Excerpt from an Architecture Decision Record. After listing context and decision, a heading "Risk & Recheck" lists hazard ("If throughput need > estimate, new service could crash"), plan ("Load test again 2 weeks post-launch, steward: PerfEngineer"), window ("gradual rollout to 20% traffic"), etc. - *(d) OODA:* A depiction of an OODA loop with an added checkpoint between Decide and Act labeled "ClipCard?". In a military exercise scenario, after deciding on a course of action, the commander quickly fills a ClipCard: hazard ("If assumption about enemy position wrong, our flank is exposed"), recheck (scout report in 30 min), kill criterion ("if encounter resistance beyond Level Y, retreat signal"). These examples show ClipCard's adaptability and how it fits naturally into existing structures, adding a small layer for risk management.

*Figure 3: Authority window and two-key timeline.* A timeline diagram demonstrating a high-impact software deployment with two-key approval and TTL: - T0: Developer requests deployment, fills ClipCard. Manager (Key 1) approves at T0. - T0+15min: SRE Lead (Key 2) approves. Deployment begins to 10% users (authority window enforced by system). - A TTL clock of 2 hours starts at the moment of Key 2 approval. - During the authority window, monitoring is active. At T0+1h: no errors, team does recheck (all good). They decide to continue rollout. - At ~T0+1.5h: Kill criterion triggers (simulated in diagram by an error spike). The timeline shows an immediate rollback action (vertical event line) before wider impact. - The TTL expires at T0+2h with no further approval; since rollback was already done, it simply prevents re-deployment. - The figure also shows what would happen if all went well: after authority window, at T0+2h with no issues, the system would auto-allow scaling to 100% (or require a quick renewal approval). The figure highlights how time-

bounding and dual approval create a controlled experiment window for a risky action and ensure automatic safety if conditions degrade or if time runs out without confirmation.

*Table 1: Outcome metrics and operational definitions.* This table enumerates each key metric, how it's defined, and how it's measured: - **False Approval Rate:** Defined as % of high-risk decisions that resulted in a significant incident or rollback post-implementation. Measured via incident reports linked to decisions (e.g., deployment caused outage). Calculation: (# of ClipCard-tagged actions with incident / total # of ClipCard-tagged actions) vs same ratio for non-ClipCard baseline. - **On-Time Recheck:** % of ClipCards where the planned recheck was completed by the scheduled time or trigger. Measured by comparing the "Recheck due" timestamp in ClipCard to actual completion log. In control, approximated by whether any follow-up was done within equivalent timeframe. - **Near-Miss Capture:** Number of near-misses caught per 100 high-risk decisions. Near-miss defined as an event that met kill criteria or hazard signs observed that could have led to incident but didn't due to intervention. Measured from ClipCard data (kill triggers, recheck findings) and incident logs. - **Median ClipCard Fill Time:** Median minutes from trigger of a ClipCard to it being completed (all fields filled and approved). Measured via timestamps in the tracking tool. - **Escalation Rate:** % of ClipCards that required escalation (involvement of higher authority or triggering of kill/stop). Measured by presence of two-key or kill activation in each card. - **Low-risk cycle time:** Average duration of tasks that did not trigger ClipCard, before and after implementation. Measured in days or hours from task start to finish for tasks with risk score < 18, to see if any change. Definitions ensure clarity on what each metric exactly captures, aiding reproducibility.

*Table 2: Implementation gate examples (Jira, CI, SBAR, Alerts).* This table provides concrete "recipes" for enforcing ClipCard in different systems: - **Jira:** *Rule:* If custom field "Risk = High" on issue, require filling "ClipCard" subtask before closing. *Implementation:* Jira workflow condition or validator script that checks a subtask of type "ClipCard" exists and has all required fields. Also, an Automation rule sends reminder to steward on due date. - **CI Pipeline (e.g., GitLab):** *Rule:* Pipeline queries a ClipCard API with change ID; if risk flagged and no ClipCard, it fails. *Implementation:* Include a manual approval stage requiring 2 approvers for "Production" environment. Also use environment scope limiting rollout (like a canary stage). - **SBAR Process:** *Rule:* Nurse must indicate if a case is "high risk" when preparing SBAR; if yes, attach ClipCard form. *Implementation:* Additional section on SBAR form with mandatory fields if HighRisk box checked. The supervisor on call must review it. - **Alert integration:** *Rule:* For each ClipCard kill criterion, set up a corresponding alert in monitoring. *Implementation:* e.g., if kill criterion is "error rate >5%", create an alert in Datadog with that threshold, configured to page the steward and possibly auto-roll back via script (if supported). This table acts like a quick-start guide for teams to apply ClipCard practically.

*Table 3: ClipCard audit checklist (1-minute audit criteria).* A checklist of yes/no questions an auditor or team lead can use to quickly evaluate a completed ClipCard: 1. **Hazard specific?** (Y/N) – e.g., does it clearly state a trigger and outcome, not generic. 2. **Measurable Kill Switch?** (Y/N) – is there a clear threshold and action. 3. **Steward named and backup?** (Y/N) – is one or more persons identified. 4. **Authority window defined?** (Y/N) – either percent or time limit given. 5. **One Jump/Recheck plan?** (Y/N) – at least one recheck or trigger, not left blank. 6. **Evidence provided?** (Y/N) – at least one link or attachment included. 7. **Recheck scheduled?** (Y/N) – an actual date/time or condition noted, presumably also put on calendar. Each "No" is a flag to follow up and improve that aspect. A fully compliant ClipCard should tick all boxes. This table can be printed as a reference or even a form to fill for internal quality audits.

# Appendices

## Appendix A: ClipCard Template (Markdown & JSON Example)

**Markdown Template (for manual use or wiki):**

```
**ClipCard ID:** <unique ID or link to task/ADR>

**Hazard:** If <trigger> then <failure> [in <context>].

**Recheck Plan:** <Date/time or condition for recheck>; specify "Jump" triggers
if any.

**Recheck Steward:** <Name (Role)> (Backup: <Name>)

**Authority Window:** <Scope/Duration limits> (e.g., "initial rollout 10%, 1
hour max without further approval").

**Two-Key/TTL:** <Required approvers> and TTL <expiry>; or "N/A" if not high-
impact.

**Kill Criteria:** If <observable threshold> then <immediate action>. (Can list
multiple if needed).

**Evidence Links:**
- [Link 1 description](URL or reference)
- [Link 2 description](URL or reference)
```

*Example (Markdown filled):*

```
**ClipCard ID:** PROD-117-CC1 (Linked to Change Request PROD-117)

**Hazard:** If the new payment service fails to handle currency conversion
properly, then transactions will be processed incorrectly (potential double-
charges).

**Recheck Plan:** Recheck logs and transaction accuracy at 17:00 on 2025-11-10
(4 hours post-deploy). Jump: if error alert "PaymentErrors>50/min" triggers, do
immediate check.

**Recheck Steward:** Jane Doe (Site Reliability Engineer) (Backup: John Ops)

**Authority Window:** Deploy to 5% of users for first 2 hours; auto-pause
deployment if not explicitly continued.

**Two-Key/TTL:** Two-key required – Approver 1: Tech Lead, Approver 2: SRE on-
```

call. Approval valid for 4 hours only.

**Kill Criteria:** If >1% of payment transactions error OR any instance of double-charge detected, then immediately rollback service to previous version.

**Evidence Links:**
- [Load Test Results](link-to-load-test-report.pdf)
- [Design Doc Section on Currency Logic](link-to-google-doc#section5)
- [Previous Incident INC-2045 Postmortem](link-to-Incident-2045)

**JSON Template (for tool integration):**

```
{
  "id": "PROD-117-CC1",
  "linked_item": "PROD-117",
  "hazard": "If new payment service miscalculates currency conversion, then transactions may double-charge customers.",
  "recheck": {
    "due_time": "2025-11-10T17:00:00Z",
    "jump_trigger": "alert:PaymentErrors>50min"
  },
  "steward": {
    "primary": "Jane Doe (SRE)",
    "backup": "John Ops (SRE)"
  },
  "authority_window": {
    "scope_limit": "5% users",
    "time_limit": "2h",
    "auto_pause": true
  },
  "two_key_ttl": {
    "required": true,
    "approvers": ["Tech Lead", "OnCall SRE"],
    "ttl_hours": 4
  },
  "kill_criteria": [
    {"condition": "error_rate > 1%", "action": "rollback service"},
    {"condition": "any double_charge_event", "action": "rollback service"}
  ],
  "evidence_links": [
    {"name": "Load Test Results", "url": "https://docs.company.com/perf/loadtest123.pdf"},
    {"name": "Design Doc - Currency Logic", "url": "https://docs.company.com/designs/payment#currency"},
    {"name": "Incident INC-2045 Postmortem", "url": "https://wiki.company.com/Incident-2045"}
```

```
    ]
  }
```

This JSON corresponds to the markdown example above. In practice, a system could store ClipCards as JSON for easy querying (e.g., to automate checks or compile stats). Field names are chosen to be self-explanatory. Some fields (like `jump_trigger` or structured kill criteria) might be optional arrays, etc., depending on schema. The above shows how a conditional trigger and multiple kill criteria can be represented.

## Appendix B: Blame-Safe Policy Snippet

*Excerpt to include in Safety/Ops Policy Manual:*

**Risk Recheck ("ClipCard") Policy:** To promote a culture of safety and continuous improvement, our organization adopts a blame-free approach to high-risk decision checks:

- **Blame-Free Hazard Recording:** Team members are expected to candidly document potential hazards and failure modes on ClipCards without fear. Identifying a hazard is considered a positive action for the team's safety, not an admission of incompetence. No punitive action will be taken for hazards that are later realized, nor for hazards that in hindsight were missed on a ClipCard – these are learning opportunities.

- **Recheck Steward Role:** A Recheck Steward is assigned for each high-risk action to ensure follow-up. The Steward's duty is to coordinate rechecks and escalations *as a facilitator*, not to bear personal responsibility for the outcome of the underlying decision. Accountability for decisions remains collective at the team level (and with leadership for systemic issues), consistent with our Just Culture principles. Stewardship will be rotated or shared to distribute experience and avoid burnout.

- **Authority Windows:** All high-risk changes must define an authority/impact window (limited blast radius and/or duration). Teams have the authority to act within that window. Extending beyond the window requires additional review or approval as specified in the ClipCard. This ensures reversibility – any change should be revertible within the defined safe window if issues arise. Leadership supports immediate rollback or halt actions if kill criteria are met, without needing further approval (pre-authorization is granted by the defined criteria).

- **Two-Key Approval & TTL:** Certain irreversible or broadly impactful actions (as defined by risk criteria) require two independent approvers. The approvers must each review the ClipCard and consent to proceeding. To prevent indefinite risk, such approvals carry a TTL (time-to-live) after which the action must be re-approved if not executed or if it's ongoing. This ensures conditions are re-evaluated if delays occur. It is the responsibility of the team to obtain second approval timely; however, if obtaining a second approval in an emergency would cause undue delay and harm, the team lead may proceed under emergency protocol and document rationale later – this is only for true exigent circumstances.

- **Kill Criteria and Stop Work:** Every ClipCard shall include clear kill criteria. If a kill criterion is triggered, the team is empowered to **stop the line** or roll back immediately – and must do so without fear of repercussion. Management will back decisions made under pre-defined kill criteria.

Stopping a process to prevent harm is a protected action. A post-action review will focus on how to improve criteria or process, not on penalizing the individuals who executed the stop.

- **Team-Level Incident Review:** Incidents or near-misses that occur on ClipCard-tagged decisions will be reviewed in our blameless post-mortem process. The review looks at process, context, and system factors. For example, if a hazard was overlooked or a kill trigger not acted upon, we examine why (e.g., was the trigger too hard to detect? Did the team have alert fatigue? etc.) rather than blaming the steward or author. Findings will be used to update training and possibly adjust ClipCard templates or thresholds.

- **Confidentiality and Use of Data:** ClipCard documents may contain sensitive information. They are to be stored in access-controlled systems. Data from ClipCards can be used in aggregate for safety analytics and reporting, but not to evaluate individual performance. Any personal identifiers in ClipCard content will be redacted in widely shared reports.

This policy reaffirms our commitment that safety interventions like ClipCard are learning tools, not compliance checklists to punish. Management and all staff should encourage thorough ClipCard completion and honor the intent behind it: making our operations safer together.

*(End of Policy Excerpt)*

## Appendix C: 1-Minute ClipCard Audit Checklist

Use this quick checklist to audit a filled-out ClipCard for completeness and quality. An auditor (could be a team lead, QA engineer, or safety officer) should be able to go down this list in about a minute per card:

1. **Hazard Specificity:** Does the hazard statement clearly identify a trigger and outcome?
2. *Pass:* "If X event happens, then Y failure occurs in Z scope." (Concrete and testable)

3. *Fail:* Vague or missing failure mode (e.g., "Things might go wrong.").

4. **Kill Criteria Measurability:** Are the kill criteria quantifiable or observable, with a definitive action?

5. *Pass:* Contains a threshold and an action ("If error rate >2%, do rollback").

6. *Fail:* No clear trigger point (e.g., "monitor closely and decide if it looks bad").

7. **Steward/Backup Assigned:** Is one person (and optionally a backup) named as Recheck Steward?

8. *Pass:* "Steward: Alice (QA Lead), Backup: Bob."

9. *Fail:* No name or just a team with no individual responsible.

10. **Authority Window Defined:** Does the card define an initial safe limit (percentage, duration, or scope)?

11. *Pass:* "Deploy to 10% for 1 hour" or "Gradual increase over 24h with check."

12. *Fail:* No mention of any limit, implying full rollout at once with no checkpoint.

13. **Jump/Recheck Present:** Is there at least one scheduled recheck time or conditional jump trigger?

14. *Pass:* Specifies a date/time or condition (alert) for re-evaluation.

15. *Fail:* Recheck section empty or just "N/A." (Every high-risk should have something here.)

16. **Evidence Linked:** Did they link or attach supporting evidence (logs, test results, analysis)?

17. *Pass:* Contains references to documents or data used.

18. *Fail:* Section blank or says "none." (No evidence might be okay if truly nothing applicable, but usually there is something to reference).

19. **Calendar/Alert Set (verification):** [This may require follow-up beyond the card] Is there indication that the recheck is scheduled (e.g., "calendar invite sent" note) or that alerts are configured for the kill criteria?

20. *Pass:* Mention of a calendar event or an alert ID, or auditor can confirm event in calendar system.
21. *Warn:* Not explicitly mentioned – auditor may need to remind steward to set it.

**Audit Instructions:** If any of 1–5 fails, the ClipCard is considered incomplete/insufficient and should be sent back for revision immediately (before proceeding if possible). If 6 fails, it's not an immediate show-stopper but the auditor should ask if evidence could be added for completeness. Item 7 is a verification step – if missing, auditor should prompt the steward to ensure follow-ups are in place (this item might be recorded but not count as pass/fail for card content quality).

For audit tracking, you can record the ClipCard ID, auditor, date, and any "fail" items with comments. For example: - ClipCard ID: PROD-117-CC1 – **Fail 1:** Hazard too vague ("system might break") – fixed with more specific hazard. - ClipCard ID: CLIN-05-CC2 – **Fail 5:** No recheck time given – team added "recheck next shift." - etc.

Regular audits (e.g., 10% of ClipCards) help maintain process integrity and provide feedback for training. Celebrate passes to reinforce good practices!

## Appendix D: Minimal CSV Schema for Evaluation Logging

To evaluate ClipCard, teams should log each high-risk decision (whether a ClipCard was used or not in control cases) in a structured way. Below is a minimal schema for a CSV (or database table) to collect the needed data for analysis, with definitions:

**Columns:**

- `decision_id` – *Unique identifier for the decision instance.* (Could be a ticket number, incident ID, ADR ID, etc., or a synthetic ID if needed.)

- `team` – *Cluster/team name or ID.* (Allows grouping by team/cluster in analysis. Use consistent names, or "Site A", "Site B".)
- `date` – *Date of decision completion.* (Could be date of deployment, date of SBAR action, etc. Use ISO format YYYY-MM-DD or timestamp if needed.)
- `risk_impact_score` – *Impact rating 1–5.* (As assessed prior to decision. If not explicitly rated in control cases, an approximate rating by retrospective analysis or by a risk matrix could be used.)
- `risk_uncertainty_score` – *Uncertainty rating 1–5.* (Likewise. If unknown for past, can leave blank or estimate.)
- `clipcard_used` – *Boolean (Y/N).* (Y if a ClipCard was filled for this decision, N if not. In intervention phase high-risk should be Y; in control, presumably N unless done voluntarily.)
- `two_key_required` – *Boolean (Y/N).* (Was two-key approval applied? Essentially whether the ClipCard had two_key_ttl.required true. If clipcard_used = N, leave blank or N.)
- `two_key_obtained` – *Boolean (Y/N).* (If two-key was required, was it actually obtained before action? Normally yes, but log if any exceptions.)
- `ttl_hours` – *Numeric.* (If TTL set, how many hours; else blank.)
- `authority_window` – *Text/number.* (Some representation of scope limit, e.g., "10%/1h" or could split into two columns like `window_percent` and `window_duration`.)
- `kill_triggered` – *Boolean (Y/N).* (Did any kill criteria actually trigger during execution? Essentially was a kill action taken.)
- `recheck_due` – *Datetime.* (Scheduled recheck time or condition. If condition, could note as text like "alert X or 2025-11-10 17:00, whichever first".)
- `recheck_done` – *Datetime.* (When the recheck was actually performed. Or blank if not done yet or missed.)
- `recheck_on_time` – *Boolean (Y/N).* (Y if `recheck_done` <= `recheck_due` (within tolerance, say within a few minutes or before condition escalated), N if late or not done.)
- `outcome` – *Categorical.* (Outcome of the decision: e.g., "Success" (no issues), "Incident" (false approval – major issue occurred), "NearMissHandled" (kill triggered or issue caught and mitigated), "RolledBackByKill", "RolledBackByDecision" (rolled back not due to trigger but manual decision on recheck perhaps). We can simplify to success vs incident vs near-miss.)
- `incident_id` – *Link to incident (if any).* (If outcome was Incident or NearMissHandled, a reference number to the incident report or problem ticket.)
- `fill_time_minutes` – *Numeric.* (Time spent preparing ClipCard, if recorded. Can be calculated from timestamps – if we have start and finish times for card prep.)
- `escalation` – *Boolean (Y/N).* (Did this case require escalation beyond the immediate team? e.g., management sign-off or involvement due to issues. This might overlap with two_key or kill_triggered, but could include cases where steward had to call in extra help even if kill not triggered.)
- `notes` – *Text.* (Optional freeform for any context, like "pilot only, not rolled out" or "special waiver given". This won't be used in quantitative analysis but helpful for qualitative review.)

**Definition Notes:** - Each row is one decision event. If a decision had multiple phases (like partial rollout then full), we consider the initial decision and its follow-up under one entry (with outcome reflecting if any phase failed). - Include all high-risk decisions encountered, so that in control periods clipcard_used mostly N. If possible, also log some medium risk decisions to see if any incident happened there and if threshold missed it (this could be separate analysis). - Ensure privacy: do not include personal names in these logs (we use team and role info instead). `decision_id` should be an internal code or anonymized if data leaves the company. - If a kill trigger fired and prevented an incident, we might mark outcome as

"NearMissHandled", whereas if no kill trigger but an incident happened, "Incident". If kill trigger fired but too late to avoid impact, might still be "Incident" (meaning our kill was reactive but didn't avoid all harm). - For recheck timing, if jump trigger happened and they rechecked then (earlier than scheduled), mark that as recheck_done time and that would be on time (since condition happened). - `fill_time_minutes` : if not automatically logged, teams can estimate or leave blank. Could be derived if we capture creation and completion timestamps of ClipCard artifact.

This schema can be adjusted based on tooling (for example, if using a database, booleans could be 0/1, etc.). The key is to capture enough to test H1–H4: clipcard vs not, outcomes, recheck performance, overhead, escalation frequency.

We will provide templates or scripts to collect this from whatever systems teams use (e.g., Jira queries, form outputs). Consistent logging is crucial for credible analysis.

## Appendix E: Example A/B and ITS Analysis Plans

*(This appendix provides pseudo-code and model formulas for analyzing data from the study designs, to ensure transparency and reproducibility of the analysis.)*

**E.1 Week-level A/B Analysis Plan (Alternate Weeks Design):**

Assume we have data aggregated by week, with a binary indicator whether ClipCard was active that week (1 = ClipCard used for high-risk decisions, 0 = baseline process). We measure outcomes like total incidents in that week, number of high-risk decisions, etc.

Pseudocode (in R-like syntax):

```
# Data: weekly_data with columns: week_index, clipcard_active (0/1),
# high_risk_decisions, incidents, near_misses, avg_fill_time, median_fill_time,
etc.
# We also have possibly time trend or seasonality to consider.

# Model for incident probability per decision:
glm_incident <- glm(cbind(incidents, high_risk_decisions - incidents) ~
clipcard_active + week_index,
                    family = binomial, data = weekly_data)
summary(glm_incident)
# This gives an estimate of log-odds reduction when clipcard_active=1.
# We include week_index to adjust for any linear trend over time (if needed).

# Alternatively, use GEE to account that weeks might be correlated (though here
each week is independent enough).
library(gee)
gee_incident <- gee((incidents/high_risk_decisions) ~ clipcard_active, id = 1,
                    family = binomial, data = weekly_data)
# (id=1 just means all in one cluster if we assume independence anyway)
```

```
# Model for on-time recheck rate (only meaningful when clipcard_active=1):
# We might separate analysis: compare recheck rate in ClipCard weeks vs some
proxy in control (control may not have formal rechecks, so this might be skipped
or assumed 0).
mean_on_time_recheck_clip <-
mean(weekly_data$on_time_recheck_rate[weekly_data$clipcard_active==1])
# Compare with target 0.8 via binomial test or just descriptive.

# Paired analysis (if alternating regularly):
# We can pair each ClipCard week with the preceding control week:
paired_data <- data.frame(
    pair_index = floor((weekly_data$week_index-1)/2) + 1,
    diff_incidents = NA
)
# populate diff_incidents = incident_rate(CC) - incident_rate(noCC) per pair
for(p in unique(paired_data$pair_index)) {
    cc_week <- weekly_data$incidents[2*p-1] /
weekly_data$high_risk_decisions[2*p-1]  # assuming odd weeks ClipCard
    base_week <- weekly_data$incidents[2*p] /
weekly_data$high_risk_decisions[2*p]
    paired_data$diff_incidents[p] = cc_week - base_week
}
t.test(paired_data$diff_incidents, mu=0, alternative="less")
# one-sided test if expecting CC to have less incidents -> diff < 0.
```

Additionally, for near_misses, since control likely near_miss=0 (no mechanism to catch), any non-zero in CC weeks is improvement but tricky to test. We can just report total near_misses captured vs baseline.

**E.2 Stepped-Wedge/Cluster Analysis Plan:**

We will use a mixed model or GLMM. For example, in R using `lme4`:

```
# Data at decision-level:
# columns: cluster_id, period (could be time index), clipcard_used (0/1),
incident (0/1 outcome per decision), etc.
glmm <- glmer(incident ~ clipcard_used + period + (1|cluster_id), family =
binomial, data = decision_data)
summary(glmm)
# This yields an odds ratio for incident when ClipCard used vs not.

# Or use glmer with random intercept for cluster and maybe random slope for
period if clusters differ in baseline trend:
glmm2 <- glmer(incident ~ clipcard_used + period + (1|cluster_id) + (0+period|
cluster_id), family=binomial, data=decision_data)
```

```
# For count of incidents per cluster-period (if decision-level data too
granular):
agg <- aggregate(incident ~ cluster + time_window + clipcard_phase,
data=decision_data, FUN=sum)
# clipcard_phase might be 0 or 1 if that cluster at that time is using ClipCard
agg$n_decisions <- aggregate(decision_id ~ cluster+ time_window,
data=decision_data, FUN=length)$decision_id
glmm_rate <- glmer(cbind(incident, n_decisions-incident) ~ clipcard_phase +
time_window + (1|cluster),
                    family=binomial, data=agg)

# The stepped-wedge can also be analyzed with random cluster effects and fixed
time effects explicitly:
sw_mod <- glmer(cbind(incident, n_decisions-incident) ~ clipcard_phase +
factor(time_window) + (1|cluster),
                family=binomial, data=agg)
# factor(time_window) adjusts non-parametrically for each period.

# For on-time recheck, which is cluster-level metric:
# Use something like:
recheck_mod <- glmer(cbind(on_time_rechecks, total_rechecks-on_time_rechecks) ~
phase + (1|cluster),
                        family=binomial, data=recheck_data)
# But since control has no rechecks, this might be just descriptive as well.
```

We will check model fit, possible overdispersion (use quasibinomial if needed).

**E.3 Interrupted Time Series (ITS) Plan:**

We use segmented regression on aggregated time series (all clusters combined by week, or perhaps as a GEE with multiple sites):

```
library(segmented)
# Suppose we have monthly_data: months 1..N, incidents_per_decision (rate) or
something.
model <- glm(incident_rate ~ time + post, data = monthly_data, family =
gaussian)
# where post = 0 before ClipCard intro, 1 after (level shift)
# and perhaps include time*post for slope change.
model2 <- glm(incident_rate ~ time + post + I(time*post), data=monthly_data)
# Or use segmented package:
seg_model <- segmented(model, seg.Z = ~time, psi = clip_intro_time)
summary(seg_model)
# The segmented model directly estimates slopes pre and post.

# Alternatively, use linear regression with Newey-West standard errors to
```

```
account for autocorrelation:
library(sandwich)
library(lmtest)
its_model <- lm(incident_rate ~ time + post + I(time*post), data=monthly_data)
coeftest(its_model, vcov=NeweyWest(its_model, lag=3))
```

We will include a check for autocorrelation (Durbin-Watson or plotting residuals ACF) and possibly adjust.

If multiple sites:

```
# treat each site as an ITS and do meta analysis:
site_results <- data.frame(site=unique(data$site), level_change=NA,
slope_change=NA, se_level=NA, se_slope=NA)
for(s in unique(data$site)) {
  site_data <- subset(data, site==s)
  m <- lm(rate ~ time + post + I(time*post), data=site_data)
  coefs <- coef(m); vc <- vcovHC(m, type="HC1")  # robust SE
  site_results$level_change[site_results$site==s] <- coefs["post"]
  site_results$se_level[site_results$site==s] <- sqrt(diag(vc))["post"]
  site_results$slope_change[site_results$site==s] <- coefs["I(time * post)"]
  site_results$se_slope[site_results$site==s] <- sqrt(diag(vc))["I(time *
post)"]
}
# Then meta analysis:
library(meta)
meta_level <- metagen(TE = level_change, seTE = se_level, data=site_results,
comb.fixed=FALSE, studlab=site)
print(meta_level)
```

We will document actual code and make dataset and code available (if possible, anonymized) for reproducibility as part of our reproducibility checklist.

These analysis plans ensure that anyone can follow the described approach to verify results. They cover the main hypotheses; additional analyses (like effect on near-miss counts or cycle time) would use analogous methods (e.g., Poisson regression for counts, t-tests for time differences, etc.).

1 2 3 4 5 6 7 8 9 10 11 20 21 22 35 Safety checklist compliance and a false sense of safety: New directions for research – SafetyInsights.org
https://safetyinsights.org/2024/10/09/safety-checklist-compliance-and-a-false-sense-of-safety-new-directions-for-research/

12 13 26 34 Etsy Engineering | Blameless PostMortems and a Just Culture
https://www.etsy.com/codeascraft/blameless-postmortems

14 Architecture decision record - Microsoft Azure Well-Architected Framework | Microsoft Learn
https://learn.microsoft.com/en-us/azure/well-architected/architect-role/architecture-decision-record

15 16 SBAR Tool: Situation-Background-Assessment-Recommendation | Institute for Healthcare Improvement
https://www.ihi.org/library/tools/sbar-tool-situation-background-assessment-recommendation

17 25 OODA loop - Wikipedia
https://en.wikipedia.org/wiki/OODA_loop

18 "Stop the Line" in a Hospital - Lean Blog
https://www.leanblog.org/2009/03/stop-line-in-hospital/

19 Toyota Production System
https://www.toyota-europe.com/about-us/toyota-vision-and-philosophy/toyota-production-system

23 AWS Gurus, Battle Tested Processes, On Your Team | Trek10
https://www.trek10.com/blog/enforcing-two-person-rule-aws-codepipeline

24 Tesla Robotaxi launch is a dangerous game of smoke and mirrors
https://news.ycombinator.com/item?id=44300727

27 Stepped-wedge trial - Wikipedia
https://en.wikipedia.org/wiki/Stepped-wedge_trial

28 29 30 31 The use of segmented regression in analysing interrupted time series studies: an example in pre-hospital ambulance care | Implementation Science | Full Text
https://implementationscience.biomedcentral.com/articles/10.1186/1748-5908-9-77

32 Segmented regression analysis of interrupted time series studies in …
https://pubmed.ncbi.nlm.nih.gov/12174032/

33 Understand release gates, checks, and approvals - Azure Pipelines
https://learn.microsoft.com/en-us/azure/devops/pipelines/release/approvals/?view=azure-devops