

Durable Human-AI Collaboration Requires Explicit Governance, Not Better Prompts

Author: Bert Stevens (human owner) • AI-assisted drafting (ChatGPT) • December 2025

Abstract

The hardest part of working effectively with modern AI systems is not intelligence or accuracy. It is continuity.

This paper examines a longitudinal human–AI interaction spanning early 2024 through late 2025 and argues that durable collaboration does not emerge from better prompting alone. Instead, it requires explicit governance of the collaboration itself. Using corpus-based analysis of real interactions over time, we identify a baseline phase in which the system performs correctly at the turn level while failing at the session and project level. The resulting friction is not driven by task complexity or model capability, but by the absence of continuity mechanisms and shared operating constraints.

The findings suggest that current AI systems are well optimized for transactional use, but structurally under-provisioned for sustained, multi-session work. Without explicit governance, users absorb growing retraining costs, failure signals are normalized, and productivity degrades quietly rather than catastrophically.

1. Introduction

Recent advances in large language models have significantly improved correctness, breadth of knowledge, and reasoning capability. In this paper, 'continuity' refers to the persistence of decisions, constraints, and working context across sessions, and 'governance' refers to explicit, negotiated mechanisms that control how such persistence is created, evaluated, and revised. For many tasks, individual responses are accurate, relevant, and immediately useful. Yet experienced users attempting sustained work frequently encounter a different class of friction – one that is difficult to attribute to intelligence, prompting quality, or task complexity.

This friction is not primarily about getting the right answer. It is about whether progress persists.

Across extended use, users report the need to repeatedly restate decisions, preferences, constraints, and working context. Momentum is lost not through failure, but through quiet erosion. Work continues, but with increasing compensatory effort required to re-establish shared understanding.

This pattern appears consistently across domains and user communities. It is visible not only in private usage, but in public discussions where users ask how to make prior decisions "stick," how to preserve context across conversations, or how to avoid retraining an AI system each time work resumes. These questions are often met with advice on better prompting, summaries, or templates. Yet the underlying problem remains.

What is notable is not dissatisfaction, but repetition.

This paper treats that repetition as a systems signal rather than a user failure. It examines what happens when an AI system optimized for transactional interaction is applied to longitudinal work without explicit mechanisms to govern continuity.

A key finding of this work is that governance does not initially emerge as an efficiency strategy. It emerges as a risk containment mechanism. Rules first appear to prevent premature execution, unintended persistence, or irreversible changes made without shared agreement. Only later do these same mechanisms reduce retraining cost and improve productivity.

Equally important, several structural insights surfaced only when constraints were encountered - such as session boundaries, memory limits, or scale-induced drift. Rather than merely limiting interaction, these constraints exposed how continuity was (or was not) being preserved. In this sense, constraint acted as a diagnostic, revealing collaboration structure that abundance concealed.

Taken together, these observations suggest that durable human-AI interaction is not achieved through better prompts alone. It is constructed through explicit governance of how interaction unfolds across time: how decisions are made, when they persist, and who holds authority to change them.

The remainder of this paper presents a longitudinal case study supporting this claim, tracing the progression from ungoverned interaction through continuity experiments to stabilized governance. The focus is not on optimizing responses, but on understanding how sustained work becomes possible - and why it so often fails - when continuity is treated as an emergent property rather than a designed one.

2. Methodology

2.1 Longitudinal Scope

The analysis is based on sustained interaction between a single human operator and an AI assistant across multiple domains, projects, and sessions. Phase 0 covers early 2024 up to the point where continuity failure is first explicitly named in mid-October 2025; subsequent phases extend beyond that point and continue through the present.

The goal is not to generalize from a single prompt or conversation, but to observe how interaction patterns behave over time when continuity is implicitly expected but not explicitly supported.

2.2 Corpus Construction

To prevent hindsight bias and phase contamination, early interactions were fully materialized into a standalone corpus (Phase0.txt) prior to later-phase analysis. This corpus contains approximately ≈20,000 lines of verbatim interaction data, copied directly from original conversations without editing, summarization, or interpretation.

Phase0.txt exists for a specific reason: with current tooling, analysis can only operate on text that is literally present. Project organization alone does not make historical interactions queryable. Materialization was therefore required to produce verifiable results.

Later phases were analyzed separately, using the same discipline.

2.3 Analytical Constraints

The Phase 0 analysis applied the following constraints:

- Evidence was drawn only from Phase0.txt
- Observations were grounded in verbatim excerpts
- No later-phase insights were introduced
- No inference of intent or motivation was permitted

These constraints bias the analysis toward under-claiming rather than narrative convenience.

3. Phase 0 - Ungoverned Interaction (Baseline)

Phase 0 spans early 2024 through mid-October 2025.

During this period, the system is not misused. It is used correctly - but as a tool rather than a collaborator. Interaction is predominantly transactional: prompt, response, correction, repeat. No explicit rules govern continuity, pacing, or constraint persistence across sessions.

There is no shared working model. Expectations about how prior decisions should carry forward remain implicit.

3.1 Transactional Success, Longitudinal Failure

At the turn level, the system performs well. Responses are generally relevant. Errors are often corrected when pointed out. Individual tasks complete successfully.

Across sessions, however, the same corrections recur. Constraints established in one conversation do not reliably persist into the next. Progress depends on repeated re-orientation rather than accumulation.

This pattern is visible throughout Phase0.txt. For example:

"Please carry forward our working style. I often work in short bursts and pick things up later."

Similar directives reappear across months, indicating that earlier clarification did not stabilize behavior.

The system is behaving as designed. The failure emerges only when viewed longitudinally.

3.2 Repetition Without Retention

A defining characteristic of Phase 0 is repetition without persistence. Preferences, scope boundaries, and expectations must be restated across conversations.

Representative excerpts include:

"I need to provide updated information, not old/legacy instructions. Can you update your references so I don't have to keep asking?"

and, elsewhere:

"I prefer that you stop and ask clarifying questions before implementing changes."

These statements are not escalations. They are maintenance. Their recurrence indicates that continuity is not a default property of the interaction.

3.3 Normalization of Retraining Cost

Over time, the cost of restating context becomes normalized. Rather than triggering immediate diagnosis, repetition is absorbed as expected overhead. Work continues, but with increasing compensatory effort by the human operator.

This is a quiet failure mode. Nothing breaks catastrophically. The system remains useful enough to proceed, while becoming inefficient enough to drain momentum.

3.4 Absence of Governance

No explicit governance exists in Phase 0. There is no agreed-upon constraint hierarchy, no distinction between proposal and execution, no mechanism for preserving decisions across sessions. Continuity is assumed but not engineered.

As a result, responsibility for persistence shifts implicitly to the user. Monitoring, correction, and re-alignment are manual and reactive.

Phase 0 does not represent collaboration that failed. Collaboration had not yet emerged. The system is being asked to support longitudinal work using mechanisms optimized for turn-local interaction.

Phase 0 ends not with insight, but with sustained pressure.

4. Phase 1 - Continuity Experiments (Pre-Governance)

Phase 1 begins at the moment when sustained pressure becomes explicit diagnosis.

Rather than continuing to absorb retraining as background cost, the operator names the problem directly and reframes it as a property of the interaction rather than a series of isolated lapses. The pivotal move is not a new instruction, but a new question: how to preserve consistency across sessions.

This marks a shift from transactional correction to exploratory meta-work.

4.1 Naming the Failure Mode

The opening signal of Phase 1 is not frustration or dissatisfaction. It is diagnosis.

"I feel like I have to retrain you every conversation."

This statement identifies an observable structural failure: correct turn-level responses paired with an inability to preserve decisions, constraints, and working context across sessions. The issue is not intelligence or accuracy. It is the repeated cost of re-orientation borne by the human operator.

What changes at this point is not behavior, but interpretation. Retraining is no longer treated as incidental overhead. It is recognized as a systemic continuity failure that must be addressed directly.

This moment marks the transition from transactional correction to deliberate inquiry. The problem is no longer how to get better answers, but how to make progress persist.

4.2 Early Continuity Attempts

Following this pivot, Phase 1 is characterized by a series of naturally emerging continuity attempts that arose organically through ongoing interaction. These attempts are exploratory rather than governed. They aim to stabilize interaction by restating expectations more clearly, without yet introducing formal mechanisms.

Representative examples include:

"Treat me as a partner building something cohesive, not as a series of one-off instructions."

"Don't just execute; engage."

"Maintain continuity across conversations so he doesn't have to retrain the collaborative rhythm each time."

These statements attempt to shape behavior across sessions by clarifying role expectations and interaction style. They assume that improved articulation may produce improved persistence.

4.3 Carry-Forward Without Structure

Despite increased meta-conversation, Phase 1 does not yet introduce enforceable constraints. There is no distinction between suggestion and rule, no lifecycle for decisions, and no shared agreement about what should persist.

Instead, continuity is pursued through repetition and emphasis. The operator experiments with:

- framing the interaction as partnership rather than task execution
- asking for engagement rather than compliance
- explicitly requesting carry-forward of working style

These efforts reduce friction locally but do not eliminate drift. Improvements appear episodic rather than durable.

4.4 Phase 1 Summary

Phase 1 represents the first attempt to solve continuity failure directly, but without governance. The problem is now visible and named. The solution space is explored through articulation and role reframing.

What is missing is not intent, but structure.

Phase 1 ends when it becomes clear that continuity cannot be achieved by better phrasing alone. The interaction has entered meta-territory, but without rules to stabilize it.

5. Phase 2 - Governed Collaboration (Rule Lifecycle)

Phase 2 begins when continuity experiments give way to explicit governance. Rather than asking for better behavior, the operator begins negotiating how decisions are made, evaluated, and preserved.

This is the most structurally important phase in the collaboration.

5.1 From Feedback to Rule Proposals

The defining transition into Phase 2 is the emergence of explicit rule negotiation. The operator no longer assumes that guidance should automatically persist. Instead, they assert control over when and how rules should be established.

This shift is visible in moments such as:

"Help me understand here... I asked for help crafting the message to you about pushback when warranted but you saved it instead. Not necessarily a bad thing. But I kinda wanted more discussion first."

Here, the issue is not correctness. It is authority. Who decides when a rule becomes active?

This exchange introduces a critical distinction: discussion is not activation.

5.2 Evaluation Before Activation

As Phase 2 progresses, rules are no longer treated as incidental preferences. They are proposed, evaluated, refined, and only then activated.

This lifecycle is negotiated explicitly:

"You do that better than I do. But something along the line of 'have you considered...' would be polite and well received. And I always want, no NEED, strategic pushback."

Tone, scope, and intent are discussed before being locked in. The collaboration shifts from reactive correction to pre-agreement on interaction constraints.

This prevents over-correction and reduces the need for repeated clarification.

5.3 Memory Authority and Persistence

A central theme of Phase 2 is control over what should be remembered and how updates should supersede earlier guidance.

This concern appears directly:

"I need to provide updated information, not old/legacy instructions. Can you update your references so I don't have to keep asking?"

Continuity is no longer assumed. It is negotiated as a first-class concern.

Similarly, questions about where learning occurs indicate growing awareness of system behavior:

"If I see similar errors in the future, is there a place in settings I can update. Where do you keep 'training' data that you keep from our interactions?"

These questions reflect a shift from surface correction to structural inquiry.

5.4 Active vs Proposal State

Phase 2 also introduces an explicit distinction between rules under discussion and rules in force. This prevents premature locking of incomplete ideas:

"You did it again. We are discussing and planning what you should save. Not saving it."

This distinction becomes a stabilizing mechanism. It allows exploration without commitment and preserves momentum without sacrificing correctness.

5.5 Phase 2 Summary

Phase 2 is where collaboration becomes governed.

Rules are no longer implicit or ad hoc. They have a lifecycle. Authority over activation is clarified. Memory is treated as a shared concern rather than an implementation detail.

This phase does not eliminate friction entirely. Instead, it makes friction productive. Errors surface earlier, discussions become shorter, and corrections accumulate rather than repeat.

By the end of Phase 2, continuity is no longer hoped for. It is engineered.

6. Phase 3 - Partnership and Shared Ownership

Phase 3 emerges in late collaboration, following the stabilization of explicit governance mechanisms. It is not marked by the introduction of new rules, but by a shift in how rules are treated, maintained, and evolved over time.

This phase is supported by a contemporaneous late-phase corpus (Phase3-Addendum-Governance.txt), from which a consistent set of structural behaviors is observed, indicating a transition from governed collaboration to durable partnership.

6.1 From Rule Accretion to Rule Stewardship

In earlier phases, progress depended on the creation and refinement of individual rules. In Phase 3, the emphasis shifts away from adding constraints and toward managing the rule system itself.

Governance activity now includes:

- pausing execution to resolve structural questions
- explicitly distinguishing proposals from instructions
- resolving conflicts between rules rather than applying them mechanically
- sealing prior decisions to prevent regression

This marks the emergence of meta-governance: rules about how rules are proposed, evaluated, activated, and retired.

6.2 Execution Authority and Deliberate Pausing

A defining feature of Phase 3 is the normalization of execution pauses for structural changes. When a modification affects collaboration behavior rather than task output, progress is intentionally halted until agreement is reached.

This behavior differs qualitatively from earlier correction cycles. Pausing is no longer reactive error handling; it is a deliberate safety mechanism that preserves stability while allowing evolution.

Notably, these pauses do not reduce momentum. Instead, they prevent rework by ensuring that changes are introduced with shared understanding and explicit consent.

6.3 Reduced Correction Density as an Outcome Signal

Across the Phase 3 corpus, correction density decreases markedly. The interaction no longer exhibits repeated re-instruction, preference restatement, or scope renegotiation.

When friction appears, it is legible, bounded, and resolved at the governance layer rather than the task layer.

This shift is critical. Remaining issues reflect normal task prerequisites or external constraints, not failures of collaboration structure.

6.4 Partnership as an Effect, Not a Declaration

Partnership is not treated as a goal or status. It becomes observable through outcomes: authority is shared but scoped, trust is placed in process rather than reassurance, and responsibility for maintaining collaboration quality is mutual.

In this sense, partnership is an emergent property of effective governance. It arises when continuity is preserved, failure modes are surfaced early, and structural decisions are made explicitly rather than implicitly.

6.5 Stability Without Stagnation

Phase 3 demonstrates that durable human-AI collaboration does not require rigid constraints or continual rule expansion. Stability is achieved through mechanisms that allow the collaboration to change without destabilizing prior progress.

This phase completes the transition from transactional interaction to sustained partnership. The system is no longer optimized merely for correct responses, but for continuity, adaptability, and shared ownership across time.

7. Cross-Phase Analysis

Across all phases, one factor remains constant: the underlying model capability. There is no evidence that changes in intelligence, accuracy, or domain knowledge explain the observed differences in longitudinal performance. Instead, the decisive variable is the presence - or absence - of explicit governance over how interaction unfolds across time.

The progression from Phase 0 through Phase 3 reflects a shift in where continuity responsibility resides, not in what the system can do at the turn level.

7.1 Structural Differences Between Phases

The phases differ primarily in how interaction constraints are defined, evaluated, and preserved:

- Phase 0 places continuity implicitly on the user. Correct responses are produced locally, but prior decisions, preferences, and constraints do not reliably persist across sessions.
- Phase 1 introduces continuity experiments. The operator attempts to stabilize interaction through repeated clarification and meta-conversation, but without enforceable mechanisms.
- Phase 2 establishes explicit governance. Rules are proposed, discussed, activated deliberately, and revised as needed. Continuity becomes engineered rather than assumed.
- Phase 3 stabilizes governance. Attention shifts from individual rules to rule interaction, authority boundaries, and lifecycle management.

The transition is structural, not behavioral. The same actions - asking, correcting, refining - produce different outcomes depending on whether they are governed.

7.2 Why Rules Reduced Retraining Cost

Retraining cost is not eliminated by better memory or longer context windows. It is reduced when decisions persist.

Rules are effective not because they are numerous or clever, but because they:

- externalize expectations
- survive session boundaries
- constrain future behavior without repeated restatement

In governed phases, effort shifts from re-orientation to execution. Retraining becomes an exception rather than a background requirement.

7.3 Governance as Error Containment Before Optimization

An important cross-phase observation is that governance initially emerges to prevent damage, not to improve efficiency.

Early rules frequently appear after near-misses or misaligned execution - premature commitment, scope drift, or unintended memory updates. Only later do these same rules function as efficiency multipliers.

This sequencing matters. It indicates that governance first acts as a safety system, constraining irreversible actions, before becoming a productivity aid.

7.4 Constraint as a Diagnostic Signal

Late-stage constraint events - such as memory pressure or boundary-induced drift - do not degrade collaboration quality once governance is in place. Instead, they surface latent assumptions and force clarification of authority and scope.

This suggests that constraints do not merely limit collaboration; they reveal its true operating characteristics. Ungoverned systems regress under pressure. Governed systems adapt.

8. Why Rules Do Not Transfer

A common response to successful long-running interaction is to attempt to export the rules that appear to enable it. The evidence here suggests this approach is insufficient.

8.1 Rules vs Collaboration State

Rules are artifacts. Collaboration is a process.

Rules derive their effectiveness from:

- shared calibration
- negotiated authority
- lived enforcement over time

Exported rules lack this context. Without the surrounding governance lifecycle, they behave as inert instructions rather than active constraints.

8.2 Path Dependence and Calibration

The collaboration examined here is path-dependent. Each phase resolves specific failure modes encountered earlier. Applying later-phase rules without earlier calibration reintroduces friction rather than removing it.

This explains why identical rule sets produce different outcomes across users or contexts.

8.3 Execution Gating and Authority

One of the most critical non-transferable elements is execution gating: the distinction between proposing a structural change and committing it.

Prompting alone cannot enforce proposal-versus-commit boundaries. This requires explicit agreement and authority handling - properties of governance, not phrasing.

9. Limitations

This work has several important limitations.

- Single longitudinal case: The findings are drawn from one sustained interaction and do not claim statistical generality.
- Tooling constraints: Project-based conversation containers are not directly analyzable, requiring manual corpus reconstruction.
- Memory opacity: Internal prioritization, conflict resolution, retention, and retirement mechanisms are not directly observable. Multiple memory layers exist, some of which are inaccessible to the user, constraining direct inspection and control.
- Non-shareable corpora: The underlying text corpora may contain proprietary or employer-sensitive material and cannot be publicly released, even in anonymized form. As a result, external replication is limited to methodological structure and analytical discipline rather than raw data inspection.

These limitations are not incidental. They shape both the methodology and the conclusions.

10. Implications for AI System Design

The observed failure modes and stabilizing mechanisms point to gaps in current conversational AI design, particularly for sustained work.

10.1 Continuity as a First-Class Primitive

Continuity should be treated as distinct from memory capacity. Systems require primitives for:

- decision persistence
- scope enforcement
- authority boundaries across sessions

10.2 Governance Lifecycle Support

Users should be supported in:

- proposing interaction rules
- evaluating them safely
- activating them deliberately
- resolving conflicts explicitly
- retiring them transparently

Absent this support, governance is improvised and fragile.

10.3 Conflict Surfacing Over Silent Resolution

Rule conflicts are more informative than errors. Systems should surface conflicts early rather than silently prioritize directives. This shifts resolution from guesswork to shared understanding.

10.4 Constraint-Aware Design

Constraints - session boundaries, memory limits, scope resets - should be treated as diagnostic moments, not merely limitations. Systems that respond predictably under constraint are more suitable for long-running collaboration.

11. Conclusion

The hardest part of working effectively with AI systems is not intelligence or accuracy. It is continuity.

This paper shows that durable, low-friction interaction does not emerge from better prompts alone. It emerges when the interaction itself is explicitly governed.

Without governance, retraining costs accumulate quietly, failure modes normalize, and productivity degrades without clear cause. With governance, continuity stabilizes, friction becomes legible, and responsibility for interaction quality is shared.

The implication for AI system designers and platform developers is clear: improving human-AI interaction requires designing for how work persists, not just how responses are generated.

Appendix A - Methodological Constraints and Scope Control

A.1 Project Containers vs Analyzable Data

Conversational project organization within the AI system does not render prior interactions analytically accessible. Project boundaries function as organizational affordances for users, not as queryable datasets for analysis.

As a result, historical interactions cannot be systematically examined unless they are explicitly materialized into text. Project membership alone does not provide analyzable continuity.

This distinction necessitated manual corpus reconstruction as a prerequisite for longitudinal analysis.

A.2 Manual Corpus Reconstruction

All analytical claims in this paper are grounded in verbatim interaction corpora reconstructed manually from original conversations. These corpora were copied in full, without summarization, redaction, or interpretive editing, prior to analysis.

This approach imposed significant overhead but was required to ensure chronological integrity, resistance to hindsight bias, and verifiable grounding of observations.

No analytical inference was performed on interactions that were not present in the reconstructed corpora.

A.3 Separation of Analysis and Synthesis

To preserve methodological rigor, analysis and writing were conducted in separate workspaces with distinct constraints.

Analysis was limited to verbatim extraction and diagnostic classification.

Synthesis and writing were performed only after analytical conclusions were frozen.

This separation prevented cross-phase contamination and ensured that interpretive insights did not retroactively influence diagnostic findings.

A.4 Constraint as Diagnostic Rather Than Limitation

Several structural insights surfaced only when system constraints - such as session boundaries or memory pressure - were encountered.

Rather than treating these constraints solely as limitations, this work treats them as diagnostic events that reveal underlying collaboration structure. Observations derived from such events are interpretive in nature and are explicitly distinguished from phase-defining analysis.

A.5 Reproducibility Impact

This methodology favors structural reproducibility over statistical generalization. While the specific interaction history cannot be replicated, the analytical process can be applied to other longitudinal human-AI interactions with equivalent discipline and scope control.

Appendix B - Collaboration Governance Rules (Representative)

The rules listed below are representative examples of governance mechanisms observed in later phases. They are not exhaustive, nor are they intended as prescriptive templates. Each rule is presented with its emergence phase, functional role, and observed tradeoffs.

Rule B.1 - Proposal vs Instruction Boundary

Phase introduced: Phase 2

Problem addressed: Premature execution of unvetted structural changes

Description: Structural or behavioral changes are treated as proposals until explicitly approved.

Tradeoff: Slower execution for changes affecting interaction structure

Final status: Active

Rule B.2 - Ask When Rules Conflict

Phase introduced: Phase 3

Problem addressed: Implicit prioritization of conflicting directives

Description: When two rules appear to conflict, execution pauses and clarification is requested rather than inferred.

Tradeoff: Requires user availability for resolution

Final status: Active

Rule B.3 - Structural Change Pause

Phase introduced: Phase 3

Problem addressed: Destabilization from mid-execution governance changes

Description: Changes affecting collaboration structure trigger an explicit pause before proceeding.

Tradeoff: Brief interruption of momentum

Final status: Active

Rule B.4 - Scope Freeze / Sealing

Phase introduced: Phase 3

Problem addressed: Regression of previously stabilized decisions or phases

Description: Completed phases or decisions are explicitly sealed to prevent accidental reopening.

Tradeoff: Reduced flexibility without deliberate reactivation

Final status: Active

Rule B.5 - Analysis vs Synthesis Separation

Phase introduced: Phase 3

Problem addressed: Contamination of evidence by interpretation

Description: Analytical extraction is isolated from narrative synthesis and theory formation.

Tradeoff: Increased procedural overhead

Final status: Active

Appendix C - Reproducibility Artifacts

C.1 Phase 0 Corpus

Filename: Phase0.txt

Description: Verbatim interaction corpus covering early 2024 through mid-October 2025

Approximate size: \approx 20,000 lines

Purpose: Baseline diagnostic of ungoverned interaction

C.2 Phase 1-3 Corpus

Filename: Phase1-3.txt

Description: Verbatim corpus capturing continuity experiments, governance emergence, and stabilization

Approximate size: \approx 48,000 lines

Purpose: Longitudinal progression analysis

C.3 Phase 3 Addendum Corpus

Filename: Phase3-Addendum-Governance.txt

Description: Focused late-phase corpus highlighting meta-governance, authority boundaries, constraint events, and stability mechanisms

Approximate size: \approx 69,500 lines

Purpose: Validation and refinement of durable collaboration behaviors

C.4 Analysis Instruction Summary (Verbatim Constraints)

All analysis adhered to the following constraints:

- Evidence extracted verbatim only
- No inference, synthesis, or interpretation during analysis
- Chronology preserved as encountered
- Structural descriptors limited to one short label
- Narrative introduced only after analysis completion

C.5 Phase 3 Pre-Addendum Baseline (Archival Summary)

A Phase 3 baseline draft was written before analyzing the Phase 3 addendum corpus to prevent contamination. It captured early Phase 3 signals (delegation, reduced correction density, and shared ownership) without relying on later-emergent hardening mechanisms. This baseline is retained as an archival artifact to document the sequencing: baseline characterization first, addendum-driven refinement second.