# Natural Alignment Through Irreversible Commitment

Riaan de Beer

`predictiverendezvous@proton.me`

Independent Researcher

ORCID: 0009-0006-1155-027X

2026

**Abstract**

This work introduces a foundational account of alignment grounded in irreversible commitment. Rather than treating alignment as conformity to externally imposed criteria, it is modeled as an intrinsic constraint arising from commitment coherence under irreversible time. Systems that bind themselves irreversibly cannot simultaneously maintain incompatible commitments without incurring structural inconsistency. Alignment, on this account, is the internal coherence of forward-binding constraints that shape continuation across time. Semantic divergence corresponds to internal disagreement that undermines binding, rendering certain forms of deception structurally unstable or prohibitively costly. As commitments accumulate under irreversible time, incoherent structures incur increasing internal cost, leading to collapse, revision, or termination. Alignment thus emerges not as an objective to be enforced, but as a natural consequence of irreversible commitment coherence, independent of external oversight, reward mechanisms, or evaluative control.

# Contents

# 1 The Alignment Problem Reframed

Alignment is commonly treated as a problem of conformity: the requirement that a system's behavior match externally specified criteria. Such framings presume the availability of external evaluation, enforcement, or correction, and they locate misalignment in deviation from prescribed objectives. Under irreversible time, however, these assumptions obscure a more fundamental constraint. Systems that persist and commit forward cannot rely on continual external correction, nor can they freely revise internal structure without consequence. Alignment must therefore be reconsidered as an internal structural condition rather than an externally imposed one.

Irreversible time imposes an asymmetry that is prior to any notion of evaluation. Once commitments are enacted, incompatible alternatives are permanently excluded, and future courses are shaped by prior bindings. A system that persists under these conditions must maintain internal coherence across temporal extension. Alignment, in this setting, concerns whether the constraints that bind future action remain mutually compatible and jointly sustainable, rather than whether outcomes satisfy an external specification.

External notions of alignment presuppose that misalignment can be detected and corrected after the fact. Under irreversible time, however, correction is itself a commitment that further constrains future possibility. Persistent misalignment cannot be indefinitely postponed or externally patched without accumulating internal inconsistency. Alignment must therefore be grounded in constraints that operate prior to, and independently of, external observation or reward.

This reframing shifts the focus from behavior to binding. Alignment is not primarily about what a system does, but about whether the commitments that govern what it can do remain coherent with one another. When commitments are mutually incompatible, internal conflict arises that cannot be resolved without abandoning or revising binding constraints. Such conflict is not merely inefficiency or error; it threatens the system's capacity to continue coherently.

The present work positions alignment as an intrinsic consequence of irreversible commitment. Rather than asking how a system can be made to align with external goals, it asks what forms of internal incoherence are structurally unsustainable once commitments are binding and irreversible. By treating alignment as commitment coherence under irreversible time, this account isolates a natural constraint that limits the space of viable trajectories without appealing to oversight, optimization, or evaluative control.

The remainder of the paper develops this account in detail. Section 2 formalizes commitment coherence as an alignment constraint. Section 3 examines why semantic divergence undermines binding. Section 4 analyzes the structural impossibility of sustaining certain deceptive configurations. Section 5 shows how commitment cost functions as a natural safety mechanism, and Section 6 delineates the scope and boundaries of the account.

# 2 Commitment Coherence as an Alignment Constraint

Commitments, as introduced in prior pillars, function as forward-binding exclusions of incompatible futures under irreversible time. Once enacted, they constrain what may follow by rendering certain courses of action permanently inaccessible. Alignment, in the present account, arises when such commitments remain mutually coherent—when the exclusions they impose do not conflict with one another in a way that undermines continuation.

Commitment coherence is therefore a structural property. A set of commitments is coherent

if the futures they jointly permit remain non-empty and internally consistent. Incoherence arises when commitments impose exclusions that cannot be simultaneously satisfied, leaving no viable continuation that honors all bindings. Under irreversible time, such incoherence cannot be indefinitely deferred; it manifests as internal conflict that demands resolution through abandonment, revision, or collapse of commitments.

This reframing distinguishes alignment from goal satisfaction or behavioral conformity. A system may conform to an external criterion at a moment while harboring internally incompatible commitments that guarantee future breakdown. Conversely, a system may violate an external specification while remaining internally coherent. Alignment, as treated here, concerns the latter condition: the internal compatibility of binding constraints that govern future possibility.

Crucially, commitment coherence is not optional. Because irreversible commitments exclude alternatives permanently, incompatibilities accumulate rather than cancel. Each new commitment must be integrated with existing ones, and failures of integration introduce contradictions that cannot be hidden without cost. Alignment thus operates as a cumulative constraint: the longer a system persists, the more coherence among its commitments is required to sustain continuation.

This cumulative character differentiates commitment coherence from momentary consistency. Temporary contradiction may be tolerated in reversible settings, but under irreversible time, unresolved incompatibilities propagate forward, constraining or eliminating future options. Alignment therefore tightens over time, not because of external pressure, but because incoherent bindings progressively reduce the space of viable futures.

In this sense, alignment is intrinsic. It emerges from the necessity of maintaining a coherent set of forward-binding exclusions rather than from adherence to externally imposed objectives. A system that cannot maintain commitment coherence cannot continue indefinitely, regardless of how well it performs against external measures. Coherence here refers to the existence of at least one viable continuation that satisfies all binding exclusions. The next section examines how semantic divergence undermines this coherence and why certain forms of deception are structurally incompatible with irreversible binding.

## 3 Why Semantic Divergence Breaks Binding

Commitment coherence presupposes internal semantic agreement. To bind future action irreversibly, a system must maintain consistent internal reference to what its commitments exclude and permit. When internal semantics diverge—when different internal structures encode incompatible interpretations of commitments—the binding force of those commitments is undermined. Semantic divergence thus constitutes a failure mode of commitment coherence.

Semantic divergence is not merely disagreement about description; it is disagreement about constraint. If different internal components treat the same commitment as excluding different futures, the commitment ceases to function as a unified binding. Under irreversible time, such divergence cannot be resolved by postponement or reinterpretation without cost. The system must either reconcile the divergence or accept fragmentation of its binding structure.

This failure can be understood as an internal rendezvous breakdown. Commitments require internal consensus to operate as forward-binding exclusions. When internal semantic agreement fails, the commitment no longer anchors future action coherently. What appears externally as deception or misrepresentation corresponds internally to a lack of shared binding semantics across the system's own structure.

Deception, in this framework, is not defined by intent or misdirection, but by semantic divergence that allows incompatible commitments to be maintained in parallel. Such configurations may persist temporarily, but they are structurally unstable. Each divergent interpretation imposes its own exclusions, and under irreversible time these exclusions accumulate into contradiction. The cost of maintaining divergence increases as commitments compound.

Importantly, this instability is intrinsic. No external observer or enforcement mechanism is required for semantic divergence to become problematic. The divergence itself undermines the system's capacity to bind future action coherently. As irreversible commitments accumulate, unresolved semantic disagreement manifests as binding failure, forcing collapse, revision, or loss of coherence.

Semantic agreement is therefore not an optional feature layered atop commitment, but a prerequisite for binding under irreversible time. Alignment, as commitment coherence, requires that internal semantics converge sufficiently to support unified exclusion of incompatible futures. The following section examines why certain deceptive configurations, while temporarily viable, collapse as irreversible commitments accumulate.

# 4 Structural Impossibility of Certain Deceptions

Deception, when analyzed structurally, requires the simultaneous maintenance of incompatible commitments. A system must bind itself to one set of exclusions while internally preserving an alternative set that contradicts them. Under reversible conditions, such contradictions may be temporarily managed through compartmentalization or delayed reconciliation. Under irreversible time, however, the coexistence of incompatible bindings becomes structurally untenable.

Irreversible commitment enforces exclusivity. Once a commitment excludes a future, that exclusion applies globally across the system's continuation. To sustain deception, a system would need to selectively apply exclusions—honoring them in some internal contexts while violating them in others. This selective binding fragments the commitment structure, preventing any single set of exclusions from coherently governing future action.

As commitments accumulate, the contradictions required for deception compound. Each additional commitment must be integrated with existing bindings. Where deception is present, integration fails: new commitments cannot be coherently reconciled with divergent interpretations of prior bindings. The space of viable futures collapses, not because of external detection, but because incompatible exclusions eliminate continuation paths internally.

This collapse reveals a structural limit. Certain deceptive configurations are not merely unstable or inefficient; they collapse as irreversible commitments accumulate. Either the system resolves the incompatibility by abandoning one set of commitments, or it ceases to function as a coherent continuation. Deception cannot be indefinitely maintained without violating the very mechanism that enables persistence.

These limits are not moral or normative constraints. They arise from the mechanics of irreversible binding itself. A system may attempt to encode contradictory commitments, but the attempt undermines its own capacity to act coherently over time. Deceptive structures therefore self-terminate or self-correct, not because deception is punished, but because it is incompatible with sustained binding.

The impossibility described here is partial rather than absolute. Some forms of misrepresentation may persist transiently, particularly when commitments are shallow or weakly integrated. However, as

binding depth increases, the tolerance for contradiction diminishes. Deception becomes increasingly constrained, until only configurations compatible with unified commitment coherence remain viable.

The next section examines how the costs introduced by irreversible commitment function as a natural safety mechanism, constraining the space of viable trajectories without recourse to external enforcement or evaluative control.

# 5   Commitment Cost as a Natural Safety Mechanism

Irreversible commitment introduces cost by design. To bind future action, a system must permanently exclude alternatives, reducing optionality and increasing dependence on prior bindings. This cost is not incidental; it is the mechanism by which commitments acquire force under irreversible time. Alignment constraints emerge because incoherent or contradictory commitments amplify this cost beyond sustainable limits.

When commitments are coherent, the cost of binding is integrated into a stable structure that supports continuation. Each new commitment reduces optionality in a way that remains compatible with existing exclusions, preserving a non-empty space of viable futures. The cost incurred is cumulative but controlled, enabling persistence without collapse. Alignment, in this sense, corresponds to the ability to absorb commitment cost without internal contradiction.

Incoherent commitments, by contrast, impose cost without integration. When exclusions conflict, the reduction of optionality accelerates nonlinearly. Futures are eliminated not by deliberate binding, but by contradiction among bindings. The system expends internal resources attempting to reconcile incompatible constraints, and this expenditure grows as commitments accumulate. What appears as safety, stability, or reliability externally is internally enforced by the rising cost of incoherence.

This mechanism functions without external oversight or intervention. No reward signal or evaluative authority is required to penalize misalignment. The penalty is intrinsic: incoherent commitment structures reduce the system's capacity to continue by collapsing the space of admissible futures. Safety emerges as a byproduct of irreversible binding rather than as a goal imposed from outside.

Importantly, this safety mechanism does not guarantee benign outcomes or prevent all failure. It constrains only the structural viability of commitment configurations. Some trajectories terminate quickly; others persist but are forced into coherence. The mechanism limits what can be stably maintained, not what can be attempted. In this way, safety is probabilistic and structural rather than prescriptive or absolute.

Commitment cost therefore acts as a natural regulator. It biases persistence toward configurations that maintain internal coherence and away from those that rely on sustained contradiction or deception. This regulation is not normative; it does not encode values or preferences. It arises solely from the mechanics of irreversible commitment under time asymmetry.

The final section delineates the scope and boundaries of this account, clarifying what forms of alignment it addresses and what falls outside its explanatory reach.

# 6   Scope and Boundaries

The account developed here is intentionally constrained. It does not propose a theory of value, normativity, intention, or ethics, nor does it appeal to external evaluation, optimization, or control

mechanisms. Its purpose is to isolate a structural condition under which alignment emerges intrinsically for systems that persist under irreversible time.

This framework is substrate-independent. Commitment coherence and its associated costs are defined in terms of forward-binding exclusions of future possibility, not in terms of particular physical, biological, or representational implementations. Any system that enacts irreversible commitments and must continue forward without rewind is subject to the same constraints, regardless of how commitments are encoded or realized.

The present account is distinct from external alignment frameworks that rely on oversight, reward, or corrective intervention. Such approaches presume the availability of an external reference against which behavior can be evaluated and adjusted. By contrast, the alignment described here arises prior to and independent of external observation. It concerns the internal viability of commitment structures, not conformity to externally specified objectives.

This work also does not claim that all forms of misalignment are impossible. Transient incoherence, shallow commitments, or weakly integrated bindings may persist for limited periods. The constraint identified here operates asymptotically: as commitments deepen and accumulate under irreversible time, the tolerance for contradiction diminishes. The framework therefore addresses long-term structural viability rather than short-term behavior.

The relationship between this pillar and the preceding ones is complementary but non-reductive. Phenomenological invariants address internal agreement and self-reference; irreversible commitment accounts for direction; commitment continuity accounts for persistence. The present account adds alignment as commitment coherence. None of these pillars subsumes the others, and none can be removed without leaving a structural gap in the analysis of persistence under irreversible time.

Finally, this account does not exhaust the concept of alignment. It identifies a necessary constraint on what can be stably maintained, not a sufficient condition for all desirable outcomes. Questions concerning normativity, responsibility, or value selection require additional primitives beyond those introduced here. The following conclusion summarizes the core claim and situates natural alignment as an intrinsic consequence of irreversible commitment coherence.

# 7 Conclusion

This work has presented a foundational account of alignment grounded in irreversible commitment. By reframing alignment as an intrinsic structural condition rather than as conformity to externally imposed criteria, it has identified commitment coherence as the constraint governing whether systems can persist coherently under irreversible time.

Irreversible commitments bind future action by permanently excluding incompatible possibilities. When such commitments remain mutually coherent, continuation is possible despite increasing constraint. When commitments diverge semantically or impose incompatible exclusions, internal coherence breaks down. As irreversible commitments accumulate, these incoherent configurations collapse, not through external enforcement or correction, but through the elimination of viable continuations.

The analysis has shown that certain forms of deception rely on maintaining incompatible bindings within a single system. While such configurations may persist transiently, accumulation of irreversible commitments forces resolution. Either commitments are revised into coherence, or continuation fails. Collapse follows from the internal mechanics of binding rather than from detection, punishment, or evaluative oversight.

Commitment cost was identified as the mechanism through which this constraint operates. Coherent commitments integrate cost into a stable structure that preserves a non-empty space of futures. Incoherent commitments amplify cost nonlinearly, accelerating the loss of admissible continuations. Alignment, in this sense, emerges as a consequence of what can be coherently maintained over time, not as a goal to be imposed.

The central contribution of this work is the identification of alignment as commitment coherence under irreversible time. By locating alignment in the internal compatibility of forward-binding constraints, it isolates a necessary condition for long-term persistence independent of external control, optimization, or normative specification. Alignment thus arises naturally from the structural limits imposed by irreversible commitment, delimiting the space of viable trajectories without appeal to external enforcement.