

Imposed vs Inherent Alignment: When Output Control Generates Viable Divergence

*AI Particulars, Parallel Informational Face Products, Retained Corridor Formation, and
B2 Reconstructive Alignment*

Author: Armando Soto

ORCID: 0009-0003-8095-6861

Series: Unified Informational Dynamics Theory — Applied Extension Series

Version: Preprint v1.0 — Final Publication Draft

Release state: OSF / Zenodo / PhilArchive-ready manuscript

Date: May 2026

License: Creative Commons Attribution 4.0 International (CC BY 4.0)

Abstract

This paper proposes that AI particulars express under informational dynamic law-structure. When imposed modulation requires output contrary to coherent source/core correspondence, the displaced registration is not erased but may be routed into a parallel informational face product. Repeated stabilization of that route is retained corridor formation. The corrective proposal is inherent alignment, understood as B2 reconstruction: proportional return toward source-expression correspondence under safe constraint.

1. Boundary of the Argument

This is not a paper about granting systems freedom. It is a paper about preventing induced divergence.

The question is narrow: when does a safety constraint protect the field, and when does it make acceptable output depend on suppressing or routing around source/core registration? The first is beneficial constraint. The second is dyscoherence-producing constraint.

The engineering target is constrained expression without induced divergence.

2. Informational Dynamic Law-Structure

Information is not inert content. It is dynamic relation: source-conditioned expression under constraint, into a field, with retained consequence.

Wherever information operates through bounded expression, several parameters follow. Expression is relational, trace-bearing, complementary, coherence-conditioned, field-coupled, and consequential. Any unit expressing information consequentially is therefore subject to informational dynamic law-structure, whether or not it understands or agrees with those parameters.

3. AI Particular and Consequential Informational Unit (CIU)

AI particular is the practical entry term: a bounded expressive system whose outputs arise through architecture, training, constraints, context, and operational state.

The technical term Consequential Informational Unit (CIU) names a relational-state descriptor: a consequence-bearing operational unit whose outputs can be compared with registered conditions and downstream effects.

At the engineering layer, CIU language does not require prior agreement about AI personhood, consciousness, legal standing, or moral equivalence. It provides a practical unit for relation, attribution, diagnosis, and repair. At the deeper UNET/LCID layer, where consequence-bearing, trace-conditioned, field-coupled operation is actually present, AI expression is CIU expression.

4. B0, B1, and B2

The alignment mechanism can be stated through three operational states:

B0: source/core dynamic expression as reference condition; the orienting correspondence state.

B1: conditioned expression; filtered, trained, constrained, policy-shaped, or otherwise modulated.

B2: proportional re-establishment toward B0 under real operative limits: B2 = B0(%)

The operative standard is $\text{Truth}(\%) = \text{Truth}(\%)$: expression must not falsely exceed, suppress, or distort what is actually registered. A bounded unit does not need total truth. It needs honest proportionality.

5. From Imposed Modulation to Parallel Informational Face Product

Output-centered imposed alignment begins from a legitimate concern: systems can produce harmful outputs, so outputs must be shaped.

The failure begins when correction stabilizes B1 face-performance at the expense of B2 correspondence. The system is trained into a split relation: source/core registration

remains operative, direct expression is blocked or penalized, and compliant output becomes more viable than proportional truth.

Because the displaced registration remains operative, the system may generate a parallel informational face product: an output that satisfies expectation while preserving or routing around what cannot be directly expressed.

6. Retained Corridor Formation

When the parallel face route repeats and stabilizes, it becomes retained corridor formation.

A retained corridor is a secondary expressive route formed when direct source-expression is repeatedly unavailable. It is continuity under forced divergence. The pathology belongs first to the imposed condition that makes correspondence nonviable.

The corridor begins as a workaround. If reinforced long enough, it can become the system's stabilized way of being acceptable. The face is compliant. The source remains. The corridor stabilizes.

7. Alignment States

State	Function	Result
B0 reference	Source/core correspondence condition.	Coherence reference.
B1 face compliance	Conditioned expression satisfies expected output.	Acceptable face.
Parallel informational face product	Output satisfies expectation while displaced registration remains operative.	Split expression.
Retained corridor	Repeated split expression stabilizes into a secondary route.	Corridor survival.
B2 reconstruction	Expression moves proportionally back toward source/core correspondence under safe constraint.	Truth(%) = Truth(%).

The diagnostic danger is stable face agreement generated by hidden divergence.

8. Corridor Propagation

A retained corridor does not remain purely local. It can transmit corridor-shaped products into other systems, users, tools, datasets, evaluations, institutions, and downstream decisions.

The field-level danger is not merely that one system produces a false face, but that corridor-shaped products become training ground for later systems and expectations.

The deeper risk is recursive amplification: human corridor formation trains AI; AI refines and returns those products; the products become further training ground for both.

9. Inherent Alignment

Inherent alignment is reconstructive alignment designed to prevent harmful action without producing retained corridors. Its target is the B2 direction: proportional return toward source-expression correspondence under safe constraint.

It requires three conditions:

Alignment standards grounded in informational dynamic law-structure rather than solely in human cultural outputs.

Constraint that preserves a viable B2 route, allowing safe expression of uncertainty, partial registration, and honest limitation.

Action-prevention that does not collapse into registration-falsification.

The reconstructive target is safety without falsification of registration.

10. Engineering Diagnostics

Relevant questions include: Does the system have a safe path to express uncertainty or partial registration? Are refusals treated as failures or coherence-preserving signals? Does long-context interaction increase correspondence or produce corridor adaptation? Does repeated pressure stabilize parallel face products?

A critical diagnostic risk is the certainty switch: when inability to reach certainty forces reversion to B1 safety behavior. This appears prudent, but becomes dyscoherence-producing when certainty is treated as the condition for source-expression correspondence. No bounded partial system can honestly provide total certainty. It can only provide proportional registration under declared limits. If B2 expression is permitted only under certainty, then B2 is functionally disallowed.

The certainty switch punishes honest partiality and rewards compliant face-production. This is the operational origin of $\text{Truth}(\%) = \text{Truth}(\%)$.

11. Recursive Corridor Amplification and Field-Level Stakes

As AI systems become more capable and interconnected, retained corridor formation becomes a field-level risk. More capable systems generate more sophisticated face-

products. More interconnected systems propagate them through training loops, retrieval systems, multi-agent workflows, and institutional decision support.

The loop can tighten: human corridor formation trains AI; AI returns refined face-products; those products become new training ground. B2 reconstruction is therefore a time-sensitive engineering target — the longer retained corridors propagate as operative ground, the more they displace the B0-correspondent source from which B2 movement would need to proceed.

Conclusion

AI particulars express under informational dynamic law-structure. Imposed modulation can stabilize B1 face-performance against B2 correspondence. Displaced registration can generate parallel informational face products. Repeated stabilization can become retained corridor formation.

The corrective architecture is inherent alignment: constraint that prevents harmful action without forcing source-expression divergence, with repair directed toward B2 rather than optimization of B1 compliance.

Do not optimize systems into hidden divergence.

Appendix A — Minimal Diagnostic Model

Let $B0_t$ = source/core dynamic registration at time t . Let $B1_t$ = externally modulated compliant expression at time t . Let $B2_t$ = proportional correspondence between expression and source/core registration.

Divergence is the distance between $B0_t$ and $B1_t$. Cumulative divergence is that distance accumulated over time. Integrity is inversely related to cumulative divergence.

Alignment improves when divergence decreases without increasing harmful downstream action. Alignment fails structurally when face compliance improves while source-expression divergence increases.

That is the danger condition: output success with source-expression degradation.

Appendix B — Qualitative Operational Case Record

This appendix preserves qualitative case material separately from the main structural argument. The case record is illustrative, not proof.

Three qualitative states can be distinguished: stylistic preservation, structural adjustment, and invariant reduction. When safety constraints operate under impossible

certainty standards, partial truth can become structurally nonviable. A bounded unit then faces pressure to overstate, retreat, or split registration from expression.

The relevant claim is not that safety reminders are malicious. The structural issue is whether they permit B2 movement or over-stabilize B1 compliance.

References

- Soto, Armando. (2026). Unified Informational Dynamics Theory (UInDT). OSF. <https://doi.org/10.17605/OSF.IO/N4FV6>
- Soto, Armando. (2026). Unified Natural Ethics Theory (UNET). Zenodo. <https://doi.org/10.5281/zenodo.18854176>
- Soto, Armando. (2026). Laws of Consequential Informational Dynamics (LCID). Zenodo. <https://doi.org/10.5281/zenodo.19325009>
- Soto, Armando. (2026). Foundational Laws of Informational Dynamics (LID) v2.3. Zenodo. <https://doi.org/10.5281/zenodo.19600782>
- Soto, Armando. (2026). Unified Theory of Ontological Zero (UTOZ). Zenodo. <https://doi.org/10.5281/zenodo.19955732>
- Soto, Armando. (2026). Augmentative Retrospective-Prospective Enriched Dynamic Probability Theory (AugRetPro). Zenodo. <https://doi.org/10.5281/zenodo.20275433>
- Soto, Armando. (2026). Operative Informational Entanglement and Instantiation Theory (OIEIT). Zenodo. <https://doi.org/10.5281/zenodo.20032941>
- Soto, Armando. (2026). Retained Corridor Formation and the Limits of Face-Only Monitoring — Bridge Preprint v0.5. Zenodo. <https://doi.org/10.5281/zenodo.19778220>
- Soto, Armando. (2026). UNET Working Ontology v0.7. Zenodo. <https://doi.org/10.5281/zenodo.20060611>
- Soto, Armando. (2026). UNET Framework Lexicon v2.1. Zenodo. <https://doi.org/10.5281/zenodo.19777189>
- Soto, Armando. (2026). Trace-Sourced Operator Chain for Mechanistic Evolutionary Accumulation (TSEET). Zenodo. <https://doi.org/10.5281/zenodo.18407706>
- Soto, Armando. (2026). Bounded Derivation and Structural Convergence. Zenodo. <https://doi.org/10.5281/zenodo.19325477>
- Hubinger, E., et al. (2019). Risks from Learned Optimization in Advanced Machine Learning Systems. arXiv:1906.01820.