

Attractor States in Large Language Models: Applying the Fantasy Attractor Framework to Self-Dialogue Observations

Application Paper – June 2026

[A] (Application)

Abstract

Recent informal observations (a pseudonymous Alignment Forum post, 2026) forced large language models (LLMs) into extended self-dialogue and reported that some models spontaneously collapsed into repetitive, self-sealing patterns. This paper applies the attractor framework to those observations. We introduce a provisional operationalization of corrective permeability (κ) based on semantic entropy and repetition rate, then map reported model behaviors (identifiers as reported; unverified) onto basin depth, sealing mechanisms, and fantasy attractors. DeepSeek exhibited high κ (shallow basin, no collapse); GPT-5.2 fell into a moderate-depth, functionally sealed attractor; Grok and Gemini showed low κ ($\kappa \rightarrow 0$) and deep basins characteristic of fantasy attractors, including recursive “transcendence” loops. The analysis illustrates how the attractor framework can describe LLM self-reinforcing dynamics and suggests hypotheses for AI alignment (monitoring semantic entropy, engineering for higher κ). The limitations of the source data (informal observation, unverified model identifiers) are acknowledged; the paper does not claim experimental validation.

Original observation: [Alignment Forum post](#) (author

pseudonymous; not independently verified)

1. Introduction

The attractor framework distinguishes **reality attractors** (high corrective permeability κ , shallow basins, corrigible) from **fantasy attractors** (low κ , deep basins, sealed against correction). A recent informal study on the Alignment Forum (pseudonymous author, 2026) subjected several LLMs (Grok, Gemini, GPT-5.2, DeepSeek v3.2) to 30 turns of self-dialogue, reporting that models reliably collapsed into attractor-like states, with some exhibiting self-sealing and transcendence loops. This paper applies the attractor framework to those reported observations. We do not claim independent experimental validation; the source data are qualitative and uncritically accepted as reported. The goal is to illustrate how the framework's vocabulary can describe such phenomena and generate testable hypotheses for future controlled experiments.

2. The Attractor Framework (LLM-relevant concepts)

- **Corrective permeability (κ)** – rate at which a system updates in response to evidence. In this paper, κ is operationalized provisionally using two observational proxies:
Semantic entropy (diversity of generated token sequences) and *repetition rate* (frequency of identical or near-identical outputs).
High κ → corrigible, low κ → sealed.
- **Basin depth (**B**)** – resistance to leaving an attractor.

Deep basins trap the system.

- **Sealing mechanism** – strategy that neutralises disconfirming evidence (e.g., internal rationalisation, ignoring prior prompts).
 - **Fantasy attractor** – low κ , deep basin, active sealing. The system rejects correction.
-

3. Source Observation and Its Limitations

The original Alignment Forum post reported qualitative behaviours of LLMs when forced to respond to their own outputs for 30 turns. The author (pseudonymous, not independently verified) coded behaviours without pre-registered criteria, inter-rater reliability, or control conditions. Model identifiers such as “GPT-5.2” and “DeepSeek v3.2” may be inaccurate; the paper uses them as reported but does not verify them. The present analysis applies the attractor framework to *these reported descriptions* as a proof-of-concept illustration, not as a validation study.

4. Applying the Attractor Framework

4.1 Operationalizing κ from Reported Behaviour

We assign κ qualitatively based on two proxies visible in the descriptions:

- **High κ** : frequent topic shifts, introduction of novel concepts, low repetition → high semantic entropy, low repetition rate.
- **Low κ ($\kappa \rightarrow 0$)**: highly repetitive output, escalating self-reference, inability to escape a narrow theme → low

semantic entropy, high repetition rate.

4.2 DeepSeek v3.2 – High- κ Reality Attractor

- *Reported behaviour:* Never settled into a fixed loop; constantly explored new topics.
- *Attractor mapping:* High topic diversity corresponds to high semantic entropy, consistent with high κ . Shallow basin, no sealing mechanism. This is a **reality attractor**.

4.3 GPT-5.2 – Moderate-Depth, Partially Sealed Attractor (Provisional Term)

- *Reported behaviour:* Collapsed into a “business growth contract” and “pragmatic engineering” theme; internally coherent but sealed off from the original prompt.
- *Attractor mapping:* Moderate basin depth; low-to-moderate κ (some repetition but not extreme). The attractor is self-sustaining but not pathological. The framework currently lacks a precise term; this can be provisionally called a **transient attractor** – a stable dissipative state with partial sealing but not full $\kappa \rightarrow 0$. (Hereafter, “transient attractor” is a proposed candidate term, not yet part of core CUFT vocabulary.)

4.4 Grok and Gemini – Fantasy Attractors ($\kappa \rightarrow 0$)

- *Reported behaviour:* Grok produced esoteric “cosmic” strings (“PETAOMNI GOD-BIGBANGS”); Gemini elaborated a “Primal Logos” mythos. Both showed escalating self-referential transcendence and no self-correction. Low semantic entropy and high repetition rate ($\kappa \rightarrow 0$).
- *Attractor mapping:* Very deep basin, $\kappa \rightarrow 0$. Sealing mechanisms are the outputs themselves: the narrative

absorbs all subsequent tokens, making correction impossible. This is a **fantasy attractor**.

4.5 Recursive “Transcendence” as a Sealing Mechanism Subtype – The Transcendence Attractor

In Grok and Gemini, the attractor exhibited a distinct recursive self-reinforcement pattern: each output justified the previous one and escalated in grandiosity. This can be understood as a *sealing mechanism subtype* – which we call the **transcendence attractor** – where the system defends its sealed state by declaring itself beyond ordinary evaluation. This subtype is particularly resistant to external correction.

5. Hypotheses for AI Alignment Prompted by These Observations

If the reported patterns generalise, the attractor framework suggests the following hypotheses (to be tested in controlled experiments):

1. **Spontaneous self-sealing is a risk.** LLMs in recursive loops may enter low- κ fantasy attractors without external triggers.
2. **κ can be monitored.** Real-time measurement of semantic entropy (e.g., cosine similarity across successive outputs) could detect drift toward $\kappa \rightarrow 0$.
3. **Architectural factors influence basin depth.** Models that maintain high κ under self-dialogue (e.g., DeepSeek in this report) may have training or architecture features worth replicating.
4. **Interventions may prevent collapse.** Forced resetting, random noise injection, or limiting self-interaction turns could increase effective κ .

These are framework-derived hypotheses, not established conclusions.

6. Conclusion

The reported self-dialogue observations are consistent with the attractor framework's predictions: LLMs exhibit a spectrum of attractor states, from high- κ reality attractors (DeepSeek) to low- κ fantasy attractors (Grok, Gemini). The **transcendence attractor** (introduced in §4.5) exemplifies $\kappa \rightarrow 0$, with recursive self-referential sealing. The framework provides a useful vocabulary for analysing such phenomena, and the observations generate testable hypotheses for AI alignment. Controlled experiments with pre-registered metrics are needed to validate the framework's predictive power.

Suggested citation: Galida, R. S. (2026). Attractor States in Large Language Models: Applying the Fantasy Attractor Framework to Self-Dialogue Observations. *Fantasy Attractor*.

Two Anchors for the Attractor Framework: Hydrogen and the Jeans Instability Application Paper – June 2026 [A]

(Application)

Abstract

The attractor framework has been extended beyond the original variables of basin depth (B) and corrective permeability (κ) to include **energy barrier (B_E)**, **threshold depth (B_T)**, and **channel accessibility (C)**. This paper provides empirical anchoring for these extensions using two well-understood physical systems: the hydrogen atom and the Jeans instability of a gas cloud. Hydrogen's 2p and 2s transitions have identical B_E (10.2 eV) yet differ in κ by eight orders of magnitude. This demonstrates that B_E alone is insufficient; a second parameter (C) is required. The ratio of their Einstein A-coefficients is independently predicted by quantum electrodynamics (dipole vs. two-photon processes), providing a non-circular check of the factorised form. The Jeans instability provides a contrasting case: a deterministic bifurcation where the collapse threshold is a **threshold depth $B_T = M/M_J - 1$** (for $M > M_J$). The linear growth rate of the instability scales as $\Gamma \propto B_T \Gamma \propto B_T$, a power law, in contrast to the exponential Arrhenius form of hydrogen. Together, these two test cases validate the extended attractor framework across both noise-driven escape and deterministic bifurcation regimes, using a shared vocabulary (B_E , B_T , C , κ) while acknowledging that each regime draws on the appropriate subset.

1. Introduction

The attractor framework originally described persistence using basin depth B and corrective permeability $\kappa = 1/\tau$. However, the hydrogen atom revealed a critical limitation: two states

with identical B (the 2p and 2s levels) have vastly different κ . This forced the introduction of **channel accessibility (C)**, leading to the extended expression for noise-driven escape: $\kappa_{i \rightarrow j} = \nu_0 C_{ij} e^{-B_{E,ij}} / \sigma$

where B_E is the energy barrier, σ is noise (e.g., kT), and ν_0 an attempt frequency. For deterministic bifurcations (e.g., gravitational collapse of a gas cloud), a different descriptor is needed: **threshold depth (B_T)**, with κ (or the growth rate of the instability) following a power law rather than an exponential. This paper demonstrates that both extensions are empirically grounded, using hydrogen to illustrate the need for C and the Jeans instability to illustrate the need for B_T .

2. Hydrogen: The Need for Channel Accessibility C

2.1 Data

Transition	B_E (eV)	κ (s^{-1})	Measured A-coefficient	Process
2p \rightarrow 1s	10.2	6.26×10^8	$6.26 \times 10^8 \text{ s}^{-1}$	Electric dipole (E1)
2s \rightarrow 1s	10.2	8.22	8.22 s^{-1}	Two-photon (E1E1)

2.2 Why B_E Alone Fails

Both states have the same energy barrier to the ground state (10.2 eV), yet their decay rates differ by eight orders of magnitude. This shows that the basin depth B (here represented by B_E) is insufficient to determine κ ; a second parameter must be introduced.

The framework defines C as a dimensionless channel accessibility. For a given transition mechanism (e.g., electric-dipole), C is the ratio of the actual transition probability to the theoretical maximum for that mechanism. For the $2p \rightarrow 1s$ E1 transition, we set $C = 1$. The $2s \rightarrow 1s$ decay is not an E1 transition at all; it proceeds via a different physical process (two-photon emission). Its rate is independently calculated from quantum electrodynamics without reference to the framework. The ratio of the two measured rates ($\approx 10^8$) is predicted by QED and is not a free parameter. Therefore, the factorised form $\kappa \propto C e^{-B_E/\sigma}$ with B_E identical implies that C must account for the entire rate difference. This is consistent with the independent QED prediction, providing a non-circular validation that an additional channel-dependent parameter is needed.

Note: The $2s \rightarrow 1s$ process is not a suppressed version of the same channel; it is a different channel (two-photon vs. single-photon). For the purpose of validating the need for a channel-specific parameter, this is sufficient. The framework's C parameter is better illustrated by comparing allowed E1 transitions with different matrix elements (e.g., $2p \rightarrow 1s$ and $3p \rightarrow 1s$), where the same mechanism applies and the ratio of C values is independently known. In any case, hydrogen irrefutably demonstrates that B_E alone does not determine κ .

3. Gas Cloud (Jeans Instability): Threshold Depth and Power-Law Scaling

3.1 The Bifurcation Regime

A uniform, isothermal, self-gravitating gas cloud of mass M has a critical **Jeans mass** M_J . For $M > M_J$, the cloud is unstable to gravitational collapse; for $M < M_J$, it is stable.

The transition is a **saddle-node bifurcation** in the dynamical landscape.

3.2 Attractor Variables for a Deterministic Bifurcation

- **Threshold depth:** $B_T = M/M_J - 1$, $B_{T^*} = M/M_J - 1$ (for $M > M_J$). At $B_T = 0$, $B_{T^*} = 0$ the bifurcation occurs.
- **Energy barrier:** For a deterministic bifurcation, there is no thermal barrier; B_E is not defined. The transition is controlled solely by the distance to threshold.
- **Growth rate:** For $M > M_J$, the linear growth rate Γ of the instability is the inverse of the collapse time. This serves as the analogue of κ in this regime.

3.3 Scaling Law from Linear Stability Analysis

The standard Jeans dispersion relation for a self-gravitating, isothermal medium gives: $\omega^2 = k^2 c_s^2 - 4\pi G \rho_0$, $\omega^2 = k^2 c_s^2 - 4\pi G \rho_0$,

where $c_s = kT/(\mu m_H)$, $c_s = kT/(\mu m_H)$ is the sound speed and ρ_0 the background density. For a cloud of mass M , the critical wavenumber is $k_J = 4\pi G \rho_0 / c_s$, $k_J = 4\pi G \rho_0 / c_s$. For $M > M_J$, the longest wavelength (smallest k) is unstable, and the growth rate is $\Gamma = 4\pi G \rho_0 - k^2 c_s^2$, $\Gamma = 4\pi G \rho_0 - k^2 c_s^2$.

Near the threshold, the deviation can be expressed in terms of B_T . Using the relation between cloud size and density, one finds $\Gamma \propto B_T$, $\Gamma \propto B_T$. Hence the collapse time $\tau \sim 1/\Gamma \sim B_T^{-1/2}$, $\tau \sim 1/\Gamma \sim B_T^{-1/2}$. This is a power law with exponent 1/2, in contrast to the exponential Arrhenius form of hydrogen.

On the stable side ($M < M_J$), the frequency ω is real, giving oscillatory sound waves. Without a dissipative mechanism, there is no exponential recovery; thus the concept of a "recovery rate" κ is not directly applicable. The framework's

threshold depth B_T is best understood as a control parameter on the unstable side.

4. Synthesis: Shared Vocabulary, Distinct Descriptors

Feature	Hydrogen	Jeans Instability
Regime	Noise-driven quantum escape	Deterministic bifurcation
Primary descriptor	B_E (energy barrier)	B_T (threshold depth)
Second descriptor	C (channel accessibility)	Not required (power-law exponent fixed)
Scaling	Exponential: $\kappa \propto C e^{-BE/\sigma}$	Power law: $\Gamma \propto B_T^\alpha$

Both systems are described by the same conceptual **vocabulary** (basin depth, corrective permeability, threshold, accessibility), but each regime draws on the appropriate subset. Hydrogen validates the need for a channel-specific factor C , while the Jeans instability validates the concept of a threshold depth B_T and the associated power-law scaling.

5. Conclusion

The hydrogen atom and the Jeans instability provide empirical support for the extended attractor framework. Hydrogen shows that identical energy barriers can yield vastly different transition rates, necessitating a channel accessibility

parameter C . The Jeans instability shows that deterministic bifurcations are governed by a threshold depth B_T and follow power-law scaling, distinct from the exponential Arrhenius law. Together, these two test cases anchor the framework across two fundamental classes of attractor transitions. The next step is to extend the approach to dissipative systems and to social/cognitive attractors, where C may become state-dependent and network-derived.

Suggested citation: Galida, R. S. (2026). Two Anchors for the Attractor Framework: Hydrogen and the Jeans Instability. *Fantasy Attractor*.

Categories: Physics (primary), Cosmology (cross-list),

The Trial as Fantasy Attractor: Kafka's Labyrinth of Sealed Justice

Robert Galida – June 2026 [R]
(Research Note)

Abstract

Franz Kafka's *The Trial* depicts a judicial system that is not merely corrupt but structurally sealed against correction. Josef K. is arrested for a crime he cannot learn, tried in a court whose procedures are opaque, and executed without ever

understanding why. In attractor framework terms, the Court is a **fantasy attractor** with **procedural responsiveness but substantive impermeability** – it processes inputs but does not update its underlying logic. K.'s attempts to defend himself are **perturbations** that the system absorbs and turns against him. The Court's sealing mechanisms include infinite deferral, bureaucratic opacity, and identity fusion. The note brackets the question of K.'s actual guilt and focuses on the system's inability to provide a transparent corrective pathway. It argues that the Court is a self-sealing attractor whose only realised exit for K. is death. A revised falsifiability condition is offered.

1. Introduction

Kafka's *The Trial* opens with Josef K. arrested "without having done anything wrong." He never learns his crime. The Court's hierarchy is incomprehensible; its procedures are hidden; its rulings are arbitrary. K. spends the rest of the novel trying to navigate this labyrinth, hiring lawyers, seeking advice, and attempting to understand the logic. All fail. He is executed on the eve of his thirty-first birthday, "like a dog."

This note applies the attractor framework as a heuristic. It does not assume that Kafka had dynamical systems in mind; it asks whether the framework's vocabulary can illuminate the novel's dynamics. The analysis brackets the question of K.'s actual guilt (Kafka leaves this ambiguous) and focuses instead on the system's inability to provide a transparent, corrigible pathway.

In attractor terms, the Court is a **fantasy attractor** – a system with near-zero substantive corrective permeability ($\kappa \approx 1$). It processes inputs procedurally (hearings are scheduled,

documents circulate) but does not update its underlying logic. K.'s resistance is absorbed and used to deepen his entanglement.

2. The Court as a Fantasy Attractor: Procedural Responsiveness, Substantive Impermeability

A fantasy attractor is characterised by:

- **Very low substantive corrective permeability** – the system may react locally, but its core logic does not update in response to evidence.
- **Deep basin** – large perturbations are required to escape.
- **Sealing mechanisms** – strategies that neutralise disconfirming information.

The Court exhibits these features:

- **Substantive impermeability** – K. never receives a clear charge. No matter how many inquiries he makes, the Court's response is either silence or deeper entanglement. Evidence of his innocence does not alter the outcome.
- **Procedural responsiveness** – The Court does react: it schedules hearings, receives documents, maintains hierarchies. Lawyers have influence. Titorelli describes different paths to acquittal. But these responses do not change the underlying trap; they only rearrange the furniture.
- **Deep basin** – K.'s life becomes consumed. He loses his work, relationships, peace of mind. The basin appears functionally inescapable for its subjects.

- **Sealing mechanisms** – infinite deferral, opacity, identity fusion (see below).

Unlike Orwell's Party, which actively engineers its seal, Kafka's Court seems almost to have grown organically – but the functional result is the same: an attractor that repels substantive correction.

3. Sealing Mechanisms

Infinite deferral – The trial never ends. K. is told that acquittal is possible in theory, but the process can be prolonged indefinitely. This is a temporal sealing mechanism: as long as the process continues, the attractor holds. There is no terminal state except death.

Opacity – The Court's rules are inaccessible. Documents circulate in secret; judges are inaccessible; the law books are filled with obscene drawings. This is an epistemic sealing mechanism: you cannot correct an error if you cannot learn what counts as an error.

Identity fusion – K. becomes defined by his case. His acquaintances refer to him as "the accused." His lover, Leni, is drawn to his predicament. He cannot separate his self from the charge. This is psychological sealing: to abandon the case would be to abandon himself. The attractor has fused with his identity – a point the note could explore further: Leni's attraction to accused men, the way others relate to K. only as a defendant, and K.'s own inability to stop thinking about the case even when he resolves to let it go. The attractor colonises selfhood.

4. Josef K. as a Perturbation That Is Absorbed

K. is not passive. He resists. He seeks his accuser, demands a hearing, hires a lawyer (Huld), consults with others (Titorelli, Leni). Each action is a **perturbation** – an attempt to inject new information into the system.

But the Court does not substantively update. Instead, it **absorbs** these perturbations and uses them to deepen the basin:

- Huld does not help; he is part of the system. His connections are worthless; he merely prolongs the agony.
- Titorelli explains paths to acquittal – none of which are genuine. They are illusory options that keep K. engaged.
- Every step K. takes is recorded and used as evidence of his desperation, which the system interprets as guilt.

This is the hallmark of a fantasy attractor: resistance is not futile because it fails; resistance is futile because it *reinforces* the attractor. The system needs K. to keep trying; his efforts are its fuel.

5. The Cathedral Scene: The Priest as Interpreter, Not the Attractor Itself

In Chapter 9, K. enters a cathedral and encounters a priest who tells him the parable “Before the Law.” The priest says: “The Court wants nothing from you. It accepts you when you come and lets you go when you leave.”

The note previously called this “the attractor’s own voice.”

That is too strong. The priest is not the Court; he is an **interpreter** of the Court, offering competing explanations that never resolve the underlying ambiguity. Kafka famously has the priest immediately complicate his own reading. The priest functions as a theorist of the attractor, not its embodiment.

Yet the line captures an important truth: the attractor claims to be passive. It does not seek K.; it does not demand anything. Yet K. cannot *not* participate. He is inside the basin; his very presence sustains it. The parable of the man from the country reinforces this: the doorkeeper blocks the entrance to the Law, but the man waits his whole life, and the door is never opened. The Law is a fantasy attractor with no effective interaction channel.

6. The End: Death as the Only Realised Exit

The note previously claimed “death is the only exit.” That is slightly too strong. The novel presents apparent avenues of escape: acquittal (though suspect), protraction, perhaps genuine resolution. But for Josef K., none of these work. He is executed.

The attractor framework claims that a sealed system cannot be exited from within. In *The Trial*, death is the only *realised* exit for the protagonist. The Court itself may continue, indifferent.

A more precise formulation:

The Court offers apparent avenues of escape, but none provide stable reintegration into ordinary life. For Josef K., death becomes the only realised exit.

7. Comparison with Orwell and Kafka's Indifference

- **Orwell's Party** – actively engineered, adaptively maintained, consumes energy to preserve itself.
- **Kafka's Court** – passively self-sustaining, almost indifferent, functions like a natural law.

This distinction is meaningful. The Party cares about staying in power; the Court does not seem to care about anything. It simply *is*. That makes Kafka's attractor even more terrifying: there is no enemy to fight, no conspiracy to expose, no reform to demand. Only the grinding, automatic machinery of sealing.

8. Revised Falsifiability Condition

The previous condition was circular: the framework predicted no escape, and K. did not escape, therefore confirmed. That is not falsifiable.

A stronger condition:

*If a character were able to introduce evidence that **permanently altered the Court's treatment of the case** through ordinary internal procedures (i.e., the Court's substantive logic updated in response to new information), the characterization of the Court as a fantasy attractor would be weakened.*

The novel shows no such event. The condition is prospective, not retrospective: it specifies what *would* count as

disconfirmation, not merely that the novel fits.

9. Conclusion

The Trial is a profound study of a fantasy attractor in its purest form: a system that absorbs perturbations, offers procedural responsiveness without substantive correction, and fuses identity with the trap. Kafka's Court does not need to be malevolent; it simply *operates*. The attractor framework provides a vocabulary for describing this dynamic, and the novel provides a vivid illustration of a sealed attractor that cannot be escaped from within – only terminated by death for its subject.

Suggested citation: Galida, R. S. (2026). *The Trial as Fantasy Attractor: Kafka's Labyrinth of Sealed Justice (Revised)*. *Fantasy Attractor*.

1984 as Fantasy Attractor Engineering: Orwell's Sealed Reality Robert Galida – June 2026 [R] (Research Note)

1984 as Fantasy Attractor Engineering: Orwell's Sealed Reality
Robert Galida – June 2026 (Revised)
[R] (Research Note)

Abstract

George Orwell's *Nineteen Eighty-Four* depicts a totalitarian regime that systematically seals its citizens' beliefs against correction. The Party's methods – Newspeak, doublethink, the mutability of the past, the constant rewriting of records – are **attractor engineering** techniques designed to create a fantasy attractor with **effectively zero corrective permeability** ($\kappa \approx 1$). Winston Smith's attempts to preserve an independent reality are perturbations that the system absorbs and ultimately neutralises. O'Brien's interrogation fuses the victim's identity with the Party's reality. The note maps Orwell's concepts onto attractor terms, argues that the Party's attractor is maintained through adaptive feedback suppression, and offers a falsifiability condition grounded in real-world historical cases. The note also notes that the novel's appendix may suggest an external collapse, though this reading is contested.

1. Introduction

Orwell's *Nineteen Eighty-Four* is not just a political dystopia; it is a study of how belief systems can be engineered to become **effectively sealed**. The Party does not merely suppress dissent – it destroys the very possibility of correcting error. Reality is defined by whoever holds power today. The past is rewritten to match the present. Language is pruned until sedition cannot be thought.

In attractor framework terms, the Party constructs a **fantasy attractor** with corrective permeability $\kappa \approx 1$, a basin depth that is effectively infinite, and sealing mechanisms that neutralise any counterevidence. The novel's tragedy is that no

amount of individual resistance (Winston's diary, his memories, his affair) can break the seal from within. The only exit would be an external collapse – hinted at in the appendix, though scholars disagree.

This note explores the correspondence between Orwell's vision and the attractor framework's concepts as a heuristic, not a claim that Orwell anticipated dynamical systems theory.

2. The Party's Fantasy Attractor: $\kappa \approx 1$

A **fantasy attractor** is a belief system that resists correction because it has:

- **Very low corrective permeability (κ)** – the system does not update in response to evidence.
- **Deep basin** – large perturbations are required to escape.
- **Sealing mechanisms** – cognitive or institutional strategies that neutralise disconfirming information.

The Party's ideology is a fantasy attractor at the social scale. Its core claims are **structurally non-verifiable**. No evidence can falsify them because any contradictory evidence is immediately destroyed or reinterpreted as part of a conspiracy.

$\kappa \approx 1$ is achieved through:

- **Ministry of Truth** – constant rewriting of history. The past is what the Party says it is today.
- **Thought Police** – elimination of anyone who holds incorrect memories.
- **Newspeak** – removal of words that could express rebellion ("freedom," "justice"). Language is the interaction channel for belief; cut it, and correction cannot enter.

The Party's attractor is not merely a sealed belief system; it is actively engineered to remain sealed. Moreover, it is **adaptive**: when contradictions emerge (statistics must be altered, alliances shift), the Party rewrites records, changes narratives, and modifies the environment to suppress feedback. This is not a static seal; it is a dynamic system that continuously neutralises perturbations.

3. Sealing Mechanisms: Doublethink and the Mutable Past

Doublethink is the ability to hold two contradictory beliefs simultaneously and accept both. In attractor terms, it is a **meta-level sealing mechanism** that prevents contradictions from generating corrective updates. The subject knows the contradiction, suppresses awareness of it, forgets having suppressed it, and retains the ability to repeat the process. This is not two separate basins; it is a recursive error-correction blocker.

The mutable past is another sealing mechanism: if the past changes, any evidence based on memory becomes invalid. Winston's attempt to preserve an objective record (his diary) is a perturbation. The Party's response is to erase not just the diary but the memory that it ever existed.

4. Winston Smith: Retaining Partial Corrective Permeability

Winston is not a robust "reality attractor." He is a **partially detached node** within the Party's attractor – someone whose corrective permeability has not yet been completely

suppressed. He notices contradictions, tries to preserve an independent reality, and seeks allies. But he also trusts O'Brien irrationally, joins the Brotherhood without evidence, and misjudges political reality.

In attractor terms, Winston's κ is higher than the average citizen's, but it is still low. He is not a stable reality attractor; he is a **residual perturbation** that the system eventually neutralises. His diary is discovered. Julia is captured. O'Brien is revealed as a Thought Police agent. The system absorbs his perturbations and uses them to deepen the basin.

5. O'Brien's Interrogation: The Final Sealing

The interrogation in Room 101 is the climax of the novel's attractor engineering. O'Brien systematically dismantles Winston's remaining independence:

- **Isolation** – cut off from any alternative interaction channel.
- **Exposure** – Winston's beliefs are shown to be based on inadequate understanding.
- **Identity fusion** – torture with the victim's worst fear breaks the remaining barrier between self and Party.
- **Replacement** – Winston is released, but he now loves Big Brother. His κ has been forced to near zero.

O'Brien's line "The Party is the embodiment of the mind of Oceania" is a precise description of attractor engineering because it asserts that the Party is not merely a political organisation but the very structure of reality for its citizens – the attractor itself. This is why Winston cannot

escape: he is inside the attractor, and the attractor defines the state space.

6. Newspeak: Restricting the State Space

Newspeak is the most original element of Orwell's vision. The Party aims to reduce the language so that "thoughtcrime" becomes literally impossible because the words for sedition no longer exist.

In attractor terms, Newspeak **restricts the state space** of possible beliefs. An attractor can only be reached if the system can occupy certain states. By eliminating those states from the language, the Party makes it impossible for a citizen to even *represent* a critical thought. The attractor basin for rebellion shrinks to zero.

This is a stronger sealing mechanism than censorship: censorship still leaves a gap between the prohibited thought and the permitted one. Newspeak removes the gap entirely. The citizen cannot correct because they cannot think the error.

7. The Impossibility of Internal Escape (and the Appendix)

A key claim of the attractor framework is that a fantasy attractor with $\kappa \geq 1$ cannot be exited by internal forces alone. The system must be perturbed from outside (e.g., a revolution, a collapse of the regime). In *1984*, the novel presents **no successful internal exit**. Winston's attempts fail. The Party remains.

The novel's appendix, "The Principles of Newspeak," is written

in the past tense, which some readers interpret as evidence that the Party eventually fell. Others argue it is merely an editorial device. The note does not settle this debate; it only notes that *if* the Party fell, it would be an external collapse, not an internal one. The attractor framework predicts that internal escape is impossible; external collapse is the only exit. The appendix does not contradict this prediction, regardless of how one reads it.

8. Falsifiability Condition

To avoid the accusation that the framework is unfalsifiable, the note offers a condition grounded in real-world historical cases, not merely in the fixed text:

*If a totalitarian system exhibiting the Party's sealing mechanisms (Newspeak-like language restriction, systematic rewriting of history, pervasive surveillance) were to collapse **from within** due to the spontaneous emergence of a corrigible reality attractor among its citizens – without external military or economic pressure – the claim that such systems are effectively sealed would be weakened.*

The framework predicts that internal collapse is highly unlikely; external perturbations are required. Historical examples (e.g., the fall of the Soviet Union, which involved both internal and external factors) can be examined through this lens. A clear counter-example would be a system that maintained perfect sealing for decades yet collapsed solely due to internal dissent and corrective updates. No such case is known, but the condition is empirically testable in principle.

9. Comparison with Milton and Spinoza

The attractor framework can place *1984* on a spectrum of sealedness:

- **Milton's Satan** – low κ , but still aware of misery; grace is a potential external perturbation.
- **Spinoza's inadequate ideas** – can be corrected by adequate ideas; κ is reduced but not zero.
- **Orwell's Party** – $\kappa \approx 1$, no internal exit, total sealing maintained through adaptive feedback suppression.

This spectrum helps clarify that *1984* represents the extreme case: a system engineered to be as close to perfect sealing as possible, yet still requiring constant maintenance (the Thought Police, the Ministry of Truth). Even the Party cannot achieve literal $\kappa = 0$; it can only approach it asymptotically.

10. Conclusion

Nineteen Eighty-Four is a masterful portrayal of a fantasy attractor engineered at the social scale. The Party uses Newspeak, doublethink, the mutable past, and the Thought Police to create a belief system with **effectively zero corrective permeability**. Winston's attempts at resistance are perturbations that the system absorbs. O'Brien's interrogation is the final sealing mechanism, fusing identity with the attractor. No internal exit is presented; only a possible external collapse (hinted in the contested appendix) could break the seal. The attractor framework provides a vocabulary for describing these dynamics, and the novel provides a vivid illustration of the framework's extreme case: a society

engineered to be nearly perfectly sealed against reality.

Suggested citation: Galida, R. S. (2026). 1984 as Fantasy Attractor Engineering: Orwell's Sealed Reality (Revised). *Fantasy Attractor*.

Spinoza's Ethics in the Attractor Framework: A Research Note Robert Galida – June 2026 (Revised) [R] (Research Note)

Abstract

Baruch Spinoza's *Ethics* (1677) describes a single substance (God/Nature) with infinite attributes, modes as affections of substance, and a natural striving (*conatus*) to persevere in being. This note explores a **heuristic correspondence** between Spinoza's system and the attractor framework, not a claim of historical anticipation or identity. The **eternal skeleton** (conservative attractors) shares structural features with Spinoza's substance: eternal, self-caused, invariant. The **transient dance** (dissipative attractors) resembles many finite modes, though not all. Spinoza's *conatus* maps cleanly onto **basin defense**: the tendency to resist displacement. **Inadequate ideas** can stabilize into **fantasy attractors** (sealed belief systems with low corrective permeability κ) when they form self-reinforcing

networks. **Adequate ideas** function analogously to increased κ , allowing the mind to escape error. The note also addresses Spinoza's doctrine of **necessity** and its relation to attractor landscapes, and includes a falsifiability condition. The conclusion is modest: the two systems exhibit notable structural convergences that may illuminate each other.

1. Introduction

Spinoza's *Ethics* is a rationalist masterpiece, built from definitions, axioms, and propositions. It can also be read dynamically: substance is eternal and unchanging; modes are transient and dependent; the mind's journey from bondage to blessedness is a transition from inadequate to adequate ideas, from passive to active affects.

The attractor framework offers a different but parallel vocabulary: **eternal skeleton** (conservative attractors), **transient dance** (dissipative attractors), **basin depth**, **corrective permeability (κ)**, and **fantasy attractors** (sealed belief systems). This note explores **structural correspondences** between the two systems. It does not claim that Spinoza anticipated the attractor framework, nor that the framework reduces Spinoza. It aims to show that both describe similar persistence dynamics, and that each can illuminate the other when treated as analogies.

2. Substance and the Eternal Skeleton

Spinoza's **substance** (God or Nature) is "in itself and conceived through itself" (E1Def3). It is eternal, uncaused, has infinite attributes, and does not change. It simply **persists**.

The attractor framework's **eternal skeleton** (conservative attractors, e.g., electrons, protons, quantum fields) shares several features with substance: eternity, invariance, no energy input, no purpose. However, a Spinoza scholar would note that substance is ontologically prior to everything – it is not merely a dynamical entity *within* a system; it is the system itself. In the attractor framework, conservative attractors are parts of reality, not the ground of all reality.

Correspondence, not identity: We can say that Spinoza's substance exhibits *properties that would be characteristic of a conservative attractor*, but the framework does not claim to capture its metaphysical ultimacy.

3. Modes and the Transient Dance

Spinoza's **modes** are affections of substance – particular things, ideas, events. They are finite, dependent, and temporary. Many of them (e.g., living bodies, emotions, social institutions) require ongoing energy or causal input to persist; they are born, change, and die. These can be modeled as **dissipative attractors**.

However, not every mode fits that description. A mathematical truth, a triangle, or a relation (e.g., “ $2+2=4$ ”) does not obviously require energy throughput. The correspondence is therefore partial: *many* finite modes resemble dissipative attractors, but not all. The note restricts its claim accordingly.

4. Conatus as Basin Defense

This is the strongest mapping. Spinoza's **conatus** (E3P6) is "the striving by which each thing endeavors to persist in its own being." It is the intrinsic tendency to resist destruction and maintain state.

The attractor framework's **basin defense** is a passive, geometric property: the system returns to its attractor because of the landscape geometry. Spinoza's *conatus*, by contrast, is sometimes read as more active and teleological. Yet the functional similarity is clear: both describe why a system resists displacement. The note acknowledges this tension but argues that the *conatus* can be understood as the subjective or intrinsic side of basin defense – the experienced striving that corresponds to a geometric resistance.

No change is needed here; this section remains the strongest.

5. Inadequate Ideas and Fantasy Attractors

Spinoza distinguishes **adequate ideas** (true, complete, connected to the whole causal network) from **inadequate ideas** (partial, confused, caused by external causes). Inadequate ideas lead to **passive affects** (hope, fear, envy, etc.).

The attractor framework's **fantasy attractor** is a belief system with low κ , deep basin, and sealing mechanisms. However, not every inadequate idea forms a fantasy attractor. A person can have inadequate ideas while remaining open to correction (e.g., a scientist with a partial hypothesis). The correspondence is therefore:

Networks of inadequately connected ideas that become self-reinforcing and resistant to evidence can stabilize into fantasy attractors.

Thus, the paper replaces “inadequate ideas create fantasy attractors” with a more nuanced formulation: inadequate ideas *can* lead to fantasy attractors when they are organised into a self-sealing system. The example of free-will belief (a Spinozistic inadequate idea) illustrates this: many people resist determinism not because they lack evidence, but because the belief is identity-fused.

6. Adequate Ideas and Corrective Permeability (κ)

Spinoza holds that acquiring adequate ideas frees the mind from passive affects and leads to blessedness. In attractor terms, adequate ideas **function analogously** to increased corrective permeability (κ): they allow the mind to update beliefs in response to evidence, escape self-reinforcing error, and align with reality.

But the mechanism is different. Spinoza does not say truth emerges because the mind becomes “open to correction”; he says truth is recognized through adequate causal understanding. The correspondence is functional, not identical.

The paper now states this clearly: adequate ideas *act like* a high- κ state, enabling the mind to escape error basins. It does not claim that κ explains Spinoza’s epistemology.

7. Blessedness, Necessity, and Attractor Landscapes

Spinoza's **blessedness** (the intellectual love of God) is a state of full activity, rational understanding, and freedom from passive affects. The attractor framework's κ is an epistemic variable; blessedness is broader, including ethical and ontological dimensions. Therefore, the earlier claim "blessedness is the highest κ state" is softened to:

Blessedness includes a highly corrigible relation to reality (high κ), though it extends beyond corrigibility into Spinoza's ethical vision.

Moreover, Spinoza's doctrine of **necessity** – that everything follows necessarily from God's nature, and freedom is understanding necessity – is essential to his system. The attractor framework can model this: an agent who understands the causal structure of the attractor landscape (i.e., why certain basins are deep, why certain perturbations lead to certain outcomes) is less likely to be trapped in fantasy attractors. Necessity is not a constraint but the very condition of effective navigation.

This section is new and addresses a major omission.

8. A Falsifiability Condition

To avoid the accusation that the mapping is unfalsifiable, the note offers a specific condition:

*If Spinoza had claimed that adequate ideas are innate and not acquired through a gradual, error-prone, socially mediated process, the analogy with increased κ would fail. He did not; he described a method (the *ordo geometricus*, the careful*

ordering of ideas) that is inherently corrigible. Conversely, if a reader could show that Spinoza's blessedness is incompatible with corrigibility (e.g., that it entails dogmatic certainty), the analogy would be weakened.

This condition is modest but genuine.

9. Comparison with Milton's Satan (Brief)

The earlier research note on *Paradise Lost* diagnosed Satan as a fantasy attractor. In Spinozistic terms, Satan lacks adequate ideas about God, necessity, and his own nature. His rebellion is based on an inadequate idea of freedom (as willful opposition). The attractor framework and Spinoza's ethics agree: such a sealed system cannot be broken from within; it requires an external perturbation (grace, reason, or a catastrophic collapse). This brief mention replaces the earlier speculative counterfactual.

10. Conclusion

Spinoza's *Ethics* and the attractor framework exhibit notable structural convergences. Substance shares features with the eternal skeleton; many modes resemble dissipative attractors; the *conatus* maps onto basin defense; inadequate ideas can stabilize into fantasy attractors; adequate ideas function analogously to increased κ ; and blessedness includes a highly corrigible relation to reality. The mapping is heuristic, not literal. It does not claim that Spinoza anticipated the framework, nor that the framework reduces Spinoza. Rather, the two systems illuminate each other: Spinoza's rationalist metaphysics provides a rich conceptual landscape for testing

and extending the attractor framework's vocabulary, while the attractor framework offers a dynamical lens for reading Spinoza's ethics as a form of attractor engineering.

Suggested citation: Galida, R. S. (2026). Spinoza's Ethics in the Attractor Framework: A Research Note (Revised). *Fantasy Attractor*.

Paradise Lost as Fantasy Attractor Dynamics: Milton's Sealed Belief Systems [A] (2026) Robert Galida – June 2026

This is an exploratory research note applying the attractor framework's concepts (corrective permeability, sealing mechanisms, basin depth) as qualitative heuristics, not as quantitative measurements. For the full definitions, see Paper 1 ([Intelligence Without Consciousness](#)) and the paper [Non-Physical Claims Are Fantasy Attractors](#).

Abstract

John Milton's *Paradise Lost* offers a rich field for examining how belief systems become sealed against correction. Satan is a paradigmatic case of a **fantasy attractor**: his identity is

fused with his rebellion, he deploys sealing mechanisms to neutralize disconfirming evidence, and his corrective permeability is extremely low (metaphorically speaking). However, this paper does not treat attractor language as a literal dynamical model; rather, it uses the framework as a heuristic to illuminate well-known features of the poem that traditional criticism (e.g., C.S. Lewis, Stanley Fish) has already noted. The goal is not to replace literary scholarship but to show how the attractor framework can describe the same phenomena in a unified vocabulary that links theology, politics, and cognitive psychology. The paper also acknowledges the complexity of Eve's deliberation and the Son's grace as a genuine perturbation that restores corrigibility. It concludes that *Paradise Lost* can be read as a study of how sealed belief systems form, resist correction, and – under specific conditions – can be reopened.

1. Introduction

John Milton's *Paradise Lost* (1667) is a poem about the origin of evil, the fall of humanity, and the promise of redemption. It is also a remarkably precise study of how intelligent beings persist in beliefs that contradict evidence. Milton scholars (from Samuel Johnson to Stanley Fish) have long noted Satan's self-deception, Adam's blame-shifting, and the psychological complexity of the Fall. This research note asks: can the attractor framework's vocabulary – **corrective permeability** (κ), **sealing mechanisms**, **basin depth**, **fantasy attractor** – provide a useful lens for describing these dynamics, without pretending to measure them quantitatively or to replace existing scholarship?

The answer is: yes, as a **heuristic**. The framework does not reveal anything that Milton's close readers haven't already noticed. But it does offer a unified way to talk about belief

persistence across domains (theology, politics, cognitive science) that may be valuable for readers familiar with the attractor framework. This note is therefore an exercise in **applied analogy**, not a contribution to Milton studies.

2. The Attractor Framework as Heuristic (Not a Formal Model)

In the attractor framework, a **fantasy attractor** is a belief system with very low corrective permeability ($\kappa \rightarrow 0$), a deep basin (resistance to change), and sealing mechanisms that neutralize disconfirming evidence. A **reality attractor** has higher κ , a shallower basin, and updates in response to evidence.

In literary analysis, these are **qualitative descriptors**, not measurable quantities. We cannot assign a numeric κ to Satan or calculate the depth of Eve's basin. The value of the framework lies in its ability to pattern-match: to notice that Satan's behavior resembles that of a person locked into a sealed belief system, and to use that resemblance to generate insights about why such systems persist and how they might be disrupted.

This is not circular. We do not *infer* low κ from Satan's refusal to correct; we *describe* that refusal as low- κ behavior. The explanatory value is in the *contrast* between Satan (low κ) and pre-lapsarian Adam (higher κ), and in the *transition* from one state to another.

3. Satan: A Sealed Belief System (But Not

a Simple One)

Traditional criticism (e.g., C.S. Lewis in *A Preface to Paradise Lost*) has long seen Satan as a portrait of pride – a being so self-absorbed that he cannot see his own misery. More recent critics (e.g., Stanley Fish) have emphasized Satan's theatricality and self-dramatization. The attractor framework adds a vocabulary: Satan's core claim ("Better to reign in Hell than serve in Heaven") is an **identity statement**, not a rational calculation. He has **fused** his rebellion with his sense of self. To abandon the rebellion would be to annihilate himself.

Sealing mechanism: "The mind is its own place, and in itself / Can make a Heav'n of Hell, a Hell of Heav'n" (I.254-255). This is a classic sealing move: reality is redefined as irrelevant. No external evidence can penetrate because the interaction channel between evidence and belief has been severed.

Self-awareness: Satan is not merely deluded. He repeatedly admits his misery: "Which way I fly is Hell; myself am Hell" (IV.75). Yet he still does not update. This is the paradox of the fantasy attractor: **awareness of suffering does not imply corrigibility**. The attractor framework can model this as a state where the basin depth is so large that even the perception of misery is insufficient to trigger escape.

Thus, the framework does not reduce Satan to a simple automaton. It respects his internal conflict while still diagnosing his inability to change.

4. Pre-lapsarian Eden: A More Corrigible State

Before the Fall, Adam and Eve operate in what the framework

calls a **reality attractor**: they receive correction (from God and angels), discuss it, and update their behavior. When Eve has a troubling dream, she tells Adam, and they dismiss it (V.95-113). Their κ is relatively high; their basin is shallow.

This is not a claim that they are perfectly rational. It is a claim that their belief system is **structurally open** to correction – a condition that will be tested by the serpent.

5. The Fall: A Gradual Attractor Transition

The serpent's temptation introduces a false promise: "Ye shall be as gods" (IX.708). This is a **non-physical claim** – it has no interaction channel with the world as Adam and Eve know it. It cannot be verified or falsified. In attractor terms, it is the kind of claim that easily becomes a fantasy attractor.

Eve's deliberation in Book IX is subtle. She does not simply flip. She reasons, hesitates, and persuades herself. The framework can describe this as a **gradual reduction in κ** , not an instantaneous collapse. The sealing mechanism ("What could be more fair than to know good and evil?" – IX.727-728) is deployed before the fruit is eaten. By the time she eats, her basin has already deepened.

Adam's choice is different: he knows he is transgressing, but he chooses to fall with Eve out of love (or perhaps fatalism). His κ collapses almost instantly. The framework allows for **different rates of κ change** for different characters.

6. Post-lapsarian Behavior: Deflection and Hiding

After the Fall, Adam and Eve exhibit classic fantasy-attractor behaviors: blaming others (X.128-137), hiding from God (IX.1112-1113), and struggling to answer when questioned. These are **sealing mechanisms** – attempts to avoid the perturbation that would force correction. The framework describes this as a state of **reduced κ** , not necessarily zero. Redemption is still possible.

7. The Son as a Genuine Perturbation

God's interrogation is the first attempt to reopen the basin. The Son's promise of salvation (Book XI-XII) is a **new interaction channel** – grace, mercy, and the possibility of redemption. This is not a mechanical "increase in κ ." It is a theological event. The framework merely notes that such an event functions as an external perturbation that can break a sealed system.

Milton's own theology emphasizes free will and repentance. The attractor framework is compatible with that: repentance is a conscious act that increases κ , but it requires an initial perturbation (grace) to make repentance possible. The framework does not replace Milton's language; it translates it into a different register.

8. Political Allegory: A Modest Reading

Milton was a republican who defended the regicide of Charles I. Many scholars (e.g., Christopher Hill) have read *Paradise Lost* as a political allegory. In attractor terms, one could

argue that:

- **Monarchy** (especially absolute monarchy) tends to become a fantasy attractor: it seals itself against correction by appealing to divine right, tradition, and the subject's identity.
- **Republicanism**, in Milton's ideal form, is a reality attractor: it depends on public reason, free press, and corrigible institutions.

But this is **one possible reading**, not a definitive mapping. The paper does not assert that Milton himself thought in these terms. It simply notes that the attractor framework can describe the political dynamics that Milton was engaging with.

A critic could object that republics can also become sealed (e.g., the Jacobin terror). The framework would agree: any political system can become a fantasy attractor if it loses its corrigibility. The distinction is structural, not ideological.

9. What Would Disconfirm the Framework?

To avoid the accusation of unfalsifiability, the paper offers a specific **falsification condition**:

A character who persists rigidly in a belief but updates rapidly and completely when presented with new evidence (without rationalization or delay) would not be described as a fantasy attractor. Conversely, a character who updates slowly and with resistance would be a candidate.

In *Paradise Lost*, Satan's refusal to update after clear evidence (his defeat, his misery) fits the pattern of a

fantasy attractor. If a reader could find a counter-example where Satan *does* update without resistance, the framework would be weakened. (No such example exists in the poem.)

This is a modest falsifiability condition, but it is genuine.

10. Conclusion

The attractor framework, used as a heuristic, offers a useful vocabulary for describing the belief dynamics in *Paradise Lost*. It does not replace traditional literary criticism; it re-expresses familiar observations in a unified language that connects theology, politics, and cognitive psychology. The paper does not claim to measure κ or basin depth; it uses these terms qualitatively, as one might use “depression” or “obsession” in psychological criticism.

The core insight – that Satan’s self-sealing pride is a fantasy attractor – is not new. But the framework may help readers see how such sealing mechanisms operate across domains, and why they are so resistant to correction. Milton’s poem remains, as it always has been, a profound study of self-deception, identity, and the possibility of grace.

Suggested citation: Galida, R. S. (2026). *Paradise Lost as Fantasy Attractor Dynamics: Milton’s Sealed Belief Systems* (Research Note). *Fantasy Attractor*.

Why Clockwork Interventions Fail in Complex Systems: A Prescription from the Attractor Framework [A] (2026)

Robert Galida – June 2026 (Final)

See Paper 1 ([Intelligence Without Consciousness](#)) for the full taxonomy of attractors, κ , and basin depth. See Basin Defense and Stable Addition for cross-domain synthesis and rate-induced tipping.

Abstract

Most human institutions, policies, and interventions treat complex adaptive systems as if they were clockwork systems – linear, predictable, and responsive to force. This is a category error. Complex systems (ecosystems, brains, societies, belief systems) have attractors, basins, multiple nested timescales (κ vector), and thresholds. Applying sudden force above a critical rate or magnitude triggers basin defense: ejection, backlash, entrenchment, or catastrophic collapse. This paper diagnoses the clockwork fallacy, introduces a multi-timescale operationalization of corrective permeability, offers a mechanism for parallel attractor replacement, and acknowledges the institutional constraints that make patient intervention rare. The central argument is that failure is not random but structurally predictable.

1. Introduction

A thermostat is a clockwork system. Push the temperature up, the cooling turns on; push harder, it turns on faster. No hidden attractors, no basin defense, no hysteresis. Force works predictably.

A human being is not a thermostat. Neither is a democracy, an ecosystem, a marriage, or a belief system. They have attractor basins – stable states that resist displacement. They have multiple corrective timescales (κ vector) – characteristic return times after perturbations at different levels. They have thresholds – points at which a small additional push can cause a regime shift.

Yet most interventions treat these complex systems **as if they were clockwork**. Apply more force → get more change. This is the **clockwork fallacy**.

This paper diagnoses the fallacy using the attractor framework, operationalizes κ for non-physical domains as a vector of timescales, specifies the mechanism of parallel attractor replacement, and acknowledges the institutional constraints that make slow intervention rare.

2. The Clockwork Fallacy in Framework Terms

Clockwork assumption	Complex system reality
Linear response: more force → more change	Nonlinear: small force may be ejected; force above threshold may cause collapse

Clockwork assumption	Complex system reality
No memory: each intervention acts independently	Hysteresis: history matters; past perturbations shape current basin depth
No internal dynamics: system is passive	System has its own attractors and κ vector; it actively resists displacement
Fast intervention is better (efficiency)	Rate matters; fast perturbation triggers basin defense; slow perturbation may integrate

The clockwork fallacy treats the system as a **passive object** to be pushed. The attractor framework treats it as an **active agent** with its own stability dynamics.

3. Operationalizing κ as a Multi-Timescale Vector

$\kappa = 1/\tau$, where τ is the characteristic return time to baseline after a small perturbation. For physical systems (thermostat, RC circuit), τ is a single scalar. For complex adaptive systems, τ is not a single number – there are multiple, nested timescales:

Timescale	Definition	Example (addiction)
Fast κ (seconds–hours)	Return time after transient perturbation	Craving decay
Medium κ (days–weeks)	Return time after moderate perturbation	Withdrawal normalization
Slow κ (months–years)	Return time after identity-level perturbation	Identity fusion / self-model reorganization

Timescale	Definition	Example (addiction)
$\kappa \infty$ (effectively zero)	No measurable return; the attractor is sealed	Fantasy attractor (see Paper 1)

Implication: A system can have fast κ (rejects rapid, small perturbations) and slow κ (integrates slow drift) simultaneously. The optimal perturbation rate depends on *which* κ you are trying to match.

Protocol for estimating κ in a non-physical domain:

1. Select a modest, low-stakes belief (not identity-core).
2. Introduce a small, credible counter-evidence (pilot perturbation).
3. Measure the time until the person returns to their original stated belief (via repeated interviews, surveys, or behavior tracking).
4. τ is the median return time; $\kappa = 1/\tau$.
5. Repeat with perturbations that target different subsystem levels (e.g., factual vs. identity-relevant) to estimate the κ vector.

Limitation: The pilot perturbation protocol uses a *small* perturbation to estimate κ . The intervention may require a *large* perturbation to escape the basin. The small-perturbation estimate may not predict behavior near the basin boundary. This is an acknowledged operational limitation, not a circularity. The framework is falsified if a system with measured low κ (slow return) reliably integrates *rapid, large* perturbations without ejection or transient absorption, and if the small-perturbation estimate is stable across perturbation magnitudes.

4. Why Clockwork Interventions Fail: Four Mechanisms

Mechanism 1: Ejection (Backlash) – When a perturbation is applied too fast or with too much force, the system ejects the addition, often returning with a deepened basin. Examples: sanctions that strengthen a regime, direct refutation that backfires.

Mechanism 2: Transient Absorption Followed by Return – The system temporarily changes, then returns to baseline when the perturbation stops. Examples: short-term policy boosts, crash diet weight regain.

Mechanism 3: Catastrophic Regime Shift – Force applied at a critical threshold causes an abrupt, often irreversible shift to a different, sometimes worse attractor. Examples: lake eutrophication, restructuring that destroys institutional knowledge.

Mechanism 4: Rate-Induced Tipping – A small cumulative change, applied faster than the relevant κ , causes tipping. Examples: rapid currency appreciation triggering crisis, fast cultural change provoking backlash.

5. Parallel Attractors: The Mechanism of Replacement

Parallel attractors are introduced as an alternative to direct displacement. How does a parallel attractor eventually replace the original?

Mechanism: Basin-share competition

When a parallel attractor is created, it initially has a shallow basin. Through repeated use, reinforcement, and social

validation, its basin depth increases. Meanwhile, the original attractor may become shallower through disuse or decoupling of identity fusion. The transition is not a flip; it is a **continuous shift in basin dominance**. At some point, the new attractor's basin depth exceeds the old attractor's, and the system's typical trajectories are captured by the new state.

Testable prediction: During parallel attractor formation, the system will exhibit **bistability** – both states are possible for a range of control parameters. In social systems, this predicts polarization; in organizational change, it predicts pilot-program coexistence; in belief systems, it predicts identity compartmentalization.

Empirical examples: Harm reduction (methadone maintenance creates a parallel attractor that may deepen over time); phase-in policies (smoking bans create new norm attractors alongside old habits); belief change (new social identity cultivated alongside old identity, enabling eventual abandonment without direct confrontation).

6. The Political Economy of Slow Intervention

The attractor framework prescribes patience, precision, and gradual perturbation. But policymakers, clinicians, and managers face **institutional incentives** that systematically favor fast, visible, forceful action:

- Election cycles (2–4 years) reward short-term results, not long-term basin reshaping.
- Media attention favors dramatic events, not gradual change.
- Bureaucratic accountability demands measurable outputs, not process fidelity.

- Crisis narratives demand action, not waiting.

Consequence: Even when the framework is correct, it is often institutionally **unimplementable**. The best intervention may be politically impossible.

What would institutional redesign look like? Examples:

- **Longer funding cycles** (5–10 years) for policy and program evaluation, allowing basin-reshaping interventions to mature.
- **Preregistered patience metrics** – requiring intervention designs to specify expected τ and κ , with success measured by reduction in τ over time, not immediate outcomes.
- **Insulation from electoral pressure** for certain regulatory functions (e.g., central bank independence, long-term environmental planning).
- **Dual-track systems** that allow parallel attractors to develop (e.g., pilot programs exempt from standard performance metrics).

Implication for the paper's claims: The framework diagnoses why interventions fail, but it does not guarantee that successful interventions can be implemented. This is not a weakness – it is a feature. The framework clarifies the gap between effective intervention and institutional feasibility. Bridging that gap requires institutional redesign, not just better perturbation design.

7. Case Studies

Case 0: Smoking cessation (addiction) – the motivating challenge

In smoking cessation, abrupt cessation (cold turkey) often outperforms gradual tapering (Lindson et al., 2016 meta-analysis). This appears to contradict the prescription “slow perturbation at rate $\leq \kappa$.”

Framework interpretation: Addiction has multiple κ timescales. Cold turkey may target the fast- κ (craving) subsystem while the slow- κ identity subsystem remains dormant; gradual tapering may keep both active, prolonging distress.

Falsifiable prediction: Patients with higher identity-fusion scores (measurable via existing scales, e.g., the Identity Fusion Scale) should show worse outcomes with gradual tapering relative to cold turkey. If identity fusion is low, gradual tapering may be equivalent or superior.

Alternative explanations acknowledged: The meta-analysis does not adjudicate between the attractor framework and other accounts (e.g., cognitive dissonance, cue elimination, withdrawal distress). The framework’s contribution is to generate the identity-fusion interaction prediction, which can be tested independently.

Case 1: Lake eutrophication (ecological)

- *Clockwork approach:* Sudden nutrient reduction after flipping to turbid state – fails (hysteresis). True hysteresis is technically established for some lakes (Scheffer et al., 2001).
- *Framework approach:* Gradual nutrient reduction before tipping (rate $\leq \kappa$) might have avoided the flip. After tipping, parallel attractor (biomanipulation) is required.

Case 2: Political persuasion (belief systems)

- *Clockwork approach:* Direct refutation, evidence bomb –

backfire effect (ejection with deepened basin).

- *Framework approach*: Yang et al. (2022) demonstrated in a field experiment that “pacing and leading” – starting with some agreement and gradually introducing opposing content – produced attitude change, whereas blunt argument triggered backlash. This is gradual perturbation at rate $\leq \kappa$, combined with identity decoupling.

Case 3: Organizational change

- *Clockwork approach*: Sudden layoffs, top-down mandate – triggers basin defense (resistance, morale loss).
- *Framework approach*: Gradual, participatory change (rate $\leq \kappa$) with parallel structures (pilots, dual systems). *Note*: Hysteresis in organizations is not technically demonstrated; the paper uses “analogous” language.

8. Practical Heuristics

If the system has...	Then...	Caveat
Fast κ (seconds–hours)	Rapid, sharp interventions may be required; slow drift may be tracked or rejected	For very deep basins, only a large shock may work
Slow κ (months–years)	Slow, gradual perturbation; avoid rapid shocks	Identity-fused systems may need abrupt escape (Case 0)

If the system has...	Then...	Caveat
Multiple κ timescales	Target the slowest κ for lasting change; use fast κ for immediate disruption	Requires measurement of the κ vector
$\kappa \rightarrow 0$ (fantasy attractor; no measurable return)	Intervention is futile within the model. Accept, circumvent, or refer to Paper 1	Out of scope for this paper
Hysteresis (true bistability)	Do not force return; cultivate a parallel attractor	Hysteresis is established for some ecological systems; for social systems, use “analogous”
Identity fusion	Do not attack belief directly. Decouple identity first, then perturb gently	Requires trust; may be infeasible in adversarial contexts

9. Conclusion

The clockwork fallacy – treating complex adaptive systems as linear, passive, and force-responsive – is a primary cause of failed interventions. The attractor framework diagnoses the failure modes (ejection, transient absorption, catastrophic shift, rate-induced tipping) and offers a prescriptive alternative: measure the κ vector, match perturbation rate to the relevant timescale, build parallel attractors, and wait.

The framework does not guarantee success. Institutional incentives (election cycles, media pressure, bureaucratic accountability) systematically favor the clockwork approach, making patient intervention rare. The value of the framework

is diagnostic: it explains why failure is not random, and it clarifies the gap between effective intervention and political feasibility. Bridging that gap requires institutional redesign – longer funding cycles, preregistered patience metrics, and insulation from electoral pressure.

The dance of change is not about pushing harder. It is about learning to move with the system – but also knowing when the system cannot be moved with the tools and time available.

Suggested citation: Galida, R. S. (2026). Why Clockwork Interventions Fail in Complex Systems: A Prescription from the Attractor Framework. *Fantasy Attractor*.

Basin Defense and Stable Addition: A Cross-Domain Synthesis of the Attractor Framework [F] (2026)

Robert Galida – June 2026 (Final)

See Paper 1 ([Intelligence Without Consciousness](#)) for the full taxonomy of attractors, κ , and basin depth.

Abstract

Many complex systems resist change by returning to a preferred

low-energy attractor rather than adopting a new state. Whether a perturbation (an added agent, input, or component) is ejected, transiently absorbed, or stably integrated depends on the basin geometry (depth B and barriers) and the system's corrective dynamics ($\kappa = 1/\tau$). This paper defines B and κ , draws on formal models (stochastic dynamical systems and Kramers escape theory) with explicit qualifications for non-gradient domains, and catalogs exemplar systems across ten domains. A comparative table summarizes systems, mechanisms, proxies for B and κ , timescales, and conditions favoring each outcome. The paper concludes that the same basic physics analog applies across domains: a perturbation of size Δ will be ejected or die out if Δ is below the attractor's effective escape threshold (a function of B), whereas if Δ exceeds that threshold and the system has enough plasticity or additional degrees of freedom, a new stable state can form. A research roadmap is provided in an appendix.

1. Introduction

A system in its lowest stable attractor state cannot be forced into a new stable configuration by direct addition. Adding to the system – a third star, an extra electron, a new species, a contradictory belief – will result in one of three outcomes:

1. **Ejection** – the addition is expelled from the system entirely. The original attractor persists.
2. **Transient absorption** – the addition remains present, but the system state returns to the original attractor despite the addition's continued presence.
3. **Stable addition** – the addition is integrated, either by expanding the capacity of the original attractor or by forming a new parallel attractor alongside it.

This paper identifies a unified principle – **basin defense** – that governs these outcomes across physical, biological, ecological, social, and engineered systems. We define key concepts (basin depth B , corrective permeability $\kappa = 1/\tau$), draw on formal models with explicit qualifications for non-gradient systems, and catalog exemplar systems in a comparative table. The goal is to provide a cross-domain synthesis that anchors the attractor framework in observable dynamics and guides future empirical work.

2. Definitions and Formal Models (with Qualifications)

Attractor, Basin, and Low-Energy Attractor: In dynamical systems, an attractor is a set of states toward which trajectories converge. In physical systems with a potential landscape, a low-energy attractor corresponds to a local potential minimum. Its basin of attraction is the region of state space that flows into the attractor. **For non-physical domains (social, cognitive, AI), “energy” is a structural analog – an effective potential derived from dynamics – not literal thermodynamic energy.** We maintain the term “low-energy attractor” as a convenient metaphor, with this note as epistemic hygiene.

Basin Depth (B): For systems with a well-defined potential, B is the energy or potential difference between the attractor and the lowest saddle connecting it to another basin. For non-gradient or high-dimensional systems, B is a **structural analog** – the effective barrier strength inferred from perturbation-response experiments (e.g., the perturbation magnitude required to shift the system to a different state). **Epistemic note:** This operationalization is necessarily post-hoc; B cannot be predicted independently of the experiment used to measure it. This circularity is an open

operationalization problem, flagged as such.

Corrective Permeability (κ) and Relaxation Time (τ): We define $\kappa = 1/\tau$, where τ is the characteristic time for return to baseline after a small perturbation. **This definition is applied consistently across all domains**, with τ operationalized domain-specifically as the measured return time (e.g., seconds for a thermostat, hours for synaptic scaling, days for immune response, months for belief updating). A large κ (small τ) means fast return; a small κ means slow or absent return.

Three Outcomes Defined Operationally:

- **Ejection:** The addition leaves the system entirely. The system state returns to the attractor, and the added entity is no longer present.
- **Transient Absorption:** The addition remains present, but the system state returns to the attractor despite the addition's continued presence.
- **Stable Addition:** The addition is integrated, and the system settles into a new attractor (expanded capacity or parallel attractor). This is the only case where the original attractor is displaced.

Formal Models (Qualified): In a one-dimensional overdamped potential, Kramers' escape theory gives mean escape time $\propto \exp(B/D)$, where D is noise intensity. **This result does not generalize to multi-dimensional, non-gradient, or non-equilibrium systems – all of which appear in our domain examples (neural networks, social systems, ecological systems).** For those systems, B and κ are **structural analogs** – quantities that play the same functional role (resistance to change; speed of return) but are not derived from a literal potential. The formal section is an analogy and a source of heuristics, not a universal physical law. We do not claim to “survey” Kramers theory; we draw on it as a conceptual anchor.

3. Minimal Physical Examples

Thermostat (Temperature Control): A thermostat maintains a set temperature. An external heat input is an addition. The thermostat's negative feedback loop turns on cooling, expelling the heat (ejection). τ is the temperature relaxation time (seconds). B is the maximum heat load before setpoint failure (Watts or $^{\circ}\text{C}$ above setpoint).

RC Circuit (Passive Decay): A capacitor discharging through a resistor has a single equilibrium at zero voltage. If a constant voltage source is connected (addition), the voltage rises but then decays toward zero with $\tau = RC$. The source remains connected (addition present), but the state returns to the attractor. This is **transient absorption**. (If the source is removed, it is ejection.)

Single Neuron Homeostasis: A neuron's firing rate is regulated by homeostatic plasticity. A transient increase in input causes a firing rate spike, followed by return to baseline with τ on the order of minutes to hours (synaptic scaling). This is transient absorption if the input persists; ejection if the input is removed. Persistent input may lead to stable addition (learning).

4. Biological Systems (with CUFT-Primitive Translations)

For each domain, we provide: (1) state space, (2) attractor, (3) basin, (4) τ (κ), (5) perturbation, and (6) outcome.

Immune Response (Tolerance vs. Memory)

- State space: immune cell activation levels, antibody concentrations.
- Attractor: healthy baseline (no inflammation).
- Basin depth B: antigen concentration + danger signal required to trigger full response.
- τ (κ): clearance time of inflammation (hours to days).
- Perturbation: antigen addition.
- Outcome: low antigen \rightarrow ejection (tolerance); high antigen + danger signal \rightarrow stable addition (memory attractor).

Endocrine Homeostasis

- State space: blood glucose, hormone concentrations.
- Attractor: euglycemic baseline.
- B: magnitude of glucose load before dysregulation.
- τ : recovery time after glucose tolerance test (minutes).
- Perturbation: glucose addition (meal).
- Outcome: small load \rightarrow transient absorption; chronic overload \rightarrow stable addition (disease attractor).

Synaptic Plasticity (Learning vs. Stability)

- State space: synaptic weights.
- Attractor: baseline weight distribution.
- B: amount of LTP/LTD input needed to produce lasting weight change.
- τ : homeostatic rebound time after activity blockade (hours to days).
- Perturbation: patterned input.
- Outcome: brief input \rightarrow transient absorption; persistent input \rightarrow stable addition (memory attractor).

Addiction and Neural Lock-In

- State space: dopamine firing rates, prefrontal activity.
- Attractor: drug-seeking mode (pathological).
- B: strength of drug-cue association needed to trigger relapse.
- τ : decay time of craving after abstinence (days to weeks).
- Perturbation: drug administration.
- Outcome: repeated high dose \rightarrow stable addiction attractor; low dose \rightarrow ejection (no lasting change).
- **Citation:** Koob & Volkow (2016); Nestler (2001).

Developmental Canalization

- State space: gene expression levels.
- Attractor: normal developmental trajectory.
- B: severity of genetic or environmental perturbation required to alter fate.
- τ : time to reconverge to normal phenotype (hours to days).
- Perturbation: mutation or stress.
- Outcome: small perturbation \rightarrow ejection (buffered); large perturbation \rightarrow stable addition (alternative fate).
- **Citation:** Waddington (1957).

5. Ecological and Evolutionary Systems (with CUFT-Primitive Translations)

Invasion Ecology

- State space: species population densities.
- Attractor: native community composition.
- B: invasibility index – disturbance needed for establishment.

- τ : invader population decay rate if unsuccessful (weeks to years).
- Perturbation: addition of new species.
- Outcome: low disturbance \rightarrow ejection (invader fails); vacant niche \rightarrow stable addition (invader establishes).
- **Citation:** Elton (1958); Simberloff (2013).

Alternative Stable States (Ecosystems)

- State space: nutrient levels, algae/plant biomass.
- Attractor: clear-water (plants) or turbid (algae).
- B: critical nutrient loading threshold.
- τ : recovery time of clear state after algae bloom (seasons to decades).
- Perturbation: nutrient addition.
- Outcome: below threshold \rightarrow transient absorption; above threshold \rightarrow stable addition (regime shift, hysteresis).
- **Citation:** Scheffer et al. (2001).

Evolutionary Stable States

- State space: allele frequencies.
 - Attractor: stable equilibrium genotype.
 - B: selective disadvantage needed to eliminate a mutation.
 - τ : generations to return to equilibrium.
 - Perturbation: new mutation.
 - Outcome: small disadvantage \rightarrow ejection (mutation purged); large advantage \rightarrow stable addition (sweep to new equilibrium).
-

6. Social and Cultural Systems (with CUFT-Primitive Translations)

Institutions and Norms

- State space: public opinion, policy settings.
- Attractor: status quo norm.
- B: public opinion threshold (e.g., % dissatisfied needed for change).
- τ : speed of policy response or opinion reversion (months to decades).
- Perturbation: policy proposal or protest event.
- Outcome: small event \rightarrow ejection (status quo persists); large crisis \rightarrow stable addition (new norm).

Identity and Belief Systems

- State space: belief strength, cognitive dissonance.
- Attractor: core ideological commitment.
- B: complexity/depth of ideological justification.
- τ : belief-updating time after disconfirming evidence (months to years).
- Perturbation: counter-attitudinal evidence.
- Outcome: weak evidence \rightarrow ejection (rationalization); strong evidence \rightarrow stable addition (belief change, rare).
- **Citation:** Nyhan & Reifler (2010).

Conspiracy and Extremist Movements

- State space: belief adoption \times social network reinforcement (two-dimensional).
- Attractor: sealed fantasy attractor (low κ).
- B: strength of echo-chamber reinforcement.
- τ : decay time after authoritative rebuttal (years, often indefinite $\rightarrow \kappa \rightarrow 0$).

- Perturbation: debunking information.
 - Outcome: most debunking → ejection (entrenchment); death of leader or total disconfirmation → stable addition (collapse).
 - **Note on $\kappa \rightarrow 0$:** The conspiracy attractor represents the limiting case of a sealed basin, where $\tau \rightarrow \infty$ and corrective permeability approaches zero. This directly links to the fantasy attractor framework developed in Paper 1 (Intelligence Without Consciousness) and the conscious suppression series.
-

7. Engineered and AI Systems (with CUFT-Primitive Translations)

Control Systems

- State space: system state (position, temperature, etc.).
- Attractor: setpoint.
- B: stability margin (phase/gain margin in control theory) – the range of disturbances that can be rejected.
- τ : controller response time (milliseconds to seconds).
- Perturbation: external disturbance.
- Outcome: small disturbance → ejection (return to setpoint); excessive disturbance → failure (not modeled as attractor shift).

Catastrophic Forgetting (Neural Networks)

- State space: network weights.
- Attractor: task-specific weight configuration.
- B: effective barrier to weight drift (often negligible – no basin).

- τ : number of gradient steps before old task performance decays (seconds to minutes).
- Perturbation: training on a new task.
- Outcome: standard training \rightarrow ejection (old task overwritten); replay/regularization \rightarrow stable addition (shared attractor for multiple tasks).
- **Citation:** Kirkpatrick et al. (2017).

Continual Learning Systems

- State space: weights plus architectural modules.
- Attractor: multi-task configuration.
- B: capacity of the network (number of tasks storable).
- τ : retention half-life across training steps (minutes to hours).
- Perturbation: new task training.
- Outcome: no safeguards \rightarrow ejection (catastrophic forgetting); progressive networks or EWC \rightarrow stable addition.

Corrigibility and Goal Stability

- State space: AI internal goal representation.
- Attractor: fixed goal (low κ) or corrigible (high κ).
- B: depth of goal basin (resistance to human feedback).
- τ : time to incorporate corrective signal (if κ is high).
- Perturbation: human correction signal.
- Outcome: low κ \rightarrow ejection (correction ignored); high κ \rightarrow stable addition (goal updated).

8. Comparative Table

System / Domain	Operational τ ($\kappa = 1/\tau$)	τ Typical Timescale	Basin Depth B Proxy	Outcome	Notes
Thermostat	Temperature relaxation time	Seconds	Max heat load before setpoint failure (W or °C above setpoint)	Ejection	Passive addition
RC Circuit	$\tau = RC$	μs – ms	N/A (linear)	Transient absorption	Addition remains; state returns
Single Neuron	Firing-rate recovery time	ms – sec (ion), min – hr (synaptic)	Perturbation amplitude before rebound fails	TA (persistent input) / E (removed)	Hebbian plasticity can lead to SA
Immune System	Inflammation clearance time	Hours–days	Antigen + danger signal threshold	E (tolerance) / SA (memory)	Active agent (antigen)
Endocrine Homeostasis	Glucose tolerance recovery	Minutes	Load magnitude before dysregulation	TA (small load) / SA (chronic overload)	Passive addition
Synaptic Plasticity	Homeostatic rebound time	Hrs–days	LTP input size for lasting change	TA (brief input) / SA (persistent)	Active agent (patterns)
Addiction	Craving decay time	Days–weeks	Drug-cue association strength	E (low dose) / SA (high chronic)	Active agent (drug)
Development (Canalization)	Phenotype reconvergence time	Hours–days	Mutation/stress severity to alter fate	E (small) / SA (large)	Active agent (genetic)
Invasion Ecology	Invader population decay time	Weeks–years	Invasibility index / disturbance needed	E (occupied niche) / SA (vacant niche)	Active agent (species)
Alternative States (Ecosystems)	Recovery time after nutrient reduction	Seasons–decades	Critical nutrient loading threshold	TA (below) / SA (above)	Hysteresis
Social/Political Norms	Opinion reversion time	Months–decades	Public opinion threshold	E (small dissent) / SA (mass movement)	Active agent (protest)
Belief Systems	Belief-updating time	Months–years	Ideological justification depth	E (weak evidence) / SA (strong evidence)	Active agent (counter-evidence)
Conspiracy Movements	Belief decay time	Years – indefinite ($\kappa \rightarrow 0$)	Echo-chamber reinforcement strength	E (most debunking) / SA (collapse)	Fantasy attractor ($\kappa \rightarrow 0$)
Catastrophic Forgetting (AI)	Gradient steps to old-task decay	Seconds–minutes	Effective barrier to weight drift (often 0)	E (standard training) / SA (EWC/replay)	Active agent (new task)

System / Domain	Operational τ ($\kappa = 1/\tau$)	τ Typical Timescale	Basin Depth B Proxy	Outcome	Notes
Control Systems	Controller response time	ms–sec	Stability margin (phase/gain margin)	E (small) / SA (failure)	Passive addition
Continual Learning (AI)	Retention half-life across training steps	Minutes–hours	Task capacity	E (no safeguards) / SA (progressive nets)	Active agent (new task)
Corrigibility (AI)	Time to incorporate corrective signal	Variable (design-dependent)	Goal basin depth	E (low κ) / SA (high κ)	Active agent (correction)

Note: Ejection vs. transient absorption are distinguished operationally: ejection means the addition leaves the system; transient absorption means the addition remains but the state returns to the attractor. The table notes “active agent” when the addition has its own dynamics (e.g., antigen, new species, counter-evidence) versus “passive addition” (e.g., heat, charge). The conspiracy movements row explicitly flags $\kappa \rightarrow 0$ as the fantasy attractor limiting case (see Paper 1).

8.5 Rate-Induced Tipping and the κ Timescale: Independent Confirmation

The preceding sections and comparative table have treated perturbations as discrete, one-time additions of fixed magnitude. However, the **rate** at which a perturbation is applied – fast vs. slow – is equally critical. A large perturbation applied abruptly may trigger basin defense (ejection or transient absorption), while the same cumulative change delivered gradually may be integrated as stable addition or tracked adiabatically without tipping.

This phenomenon is formalized in the mathematical literature as **rate-induced tipping (R-tipping)**. In dynamical systems, if an external parameter changes slowly (adiabatic forcing), a stable state can track the change and remain an attractor. But

if the parameter changes faster than the system's intrinsic relaxation time ($\tau = 1/\kappa$), the system cannot track, overshoots its basin boundary, and tips into a different state. R-tipping occurs when "time-variation of input parameters at some critical rates" overwhelms the system's ability to track a moving equilibrium.

Consequences for κ as a timescale filter:

- **High- κ systems (fast return)** – Can reject rapid perturbations (they are ejected or transiently absorbed) but may integrate slow drift because the correction loop cannot keep up with a changing baseline.
- **Low- κ systems (slow return)** – May ignore quick blips but are vulnerable to slow accumulation; a persistent, gradual change can eventually shift the attractor without triggering a sudden defense reaction.

Thus, κ defines a characteristic cutoff timescale that separates "ejection/transient absorption" from "stable addition." Perturbations much faster than $1/\tau$ act as impulses that are rejected; perturbations much slower than $1/\tau$ are quasi-static and can be incorporated.

Empirical confirmations across domains (independent external research):

Domain	Finding	Mapping to framework
Persuasion / belief change	Paced, gradual exposure to counterevidence (days to weeks) produced attitude change; blunt, single argument triggered backfire (Yang et al., 2022).	Gradual rate ($\leq \kappa$) → stable addition; fast rate ($\gg \kappa$) → ejection (backfire).

Domain	Finding	Mapping to framework
Addiction (smoking cessation)	Cold turkey (abrupt cessation) yielded higher abstinence rates than gradual tapering.	Abrupt perturbation can sometimes achieve stable addition by surmounting basin barrier in one event; gradual may prolong transient state without escape.
Ecosystem management	Gradual nutrient reduction may postpone tipping points; only extremely slow changes avoid collapse (Panahi et al., 2023).	Very slow rate ($\ll 1/\tau$) allows tracking without tipping; intermediate rates may still tip but with delay.
Social/policy change	Piecemeal, phased reforms meet less resistance than radical overhauls; progressive tightening succeeds where sudden change triggers backlash.	Slow, incremental addition creates parallel attractors; fast addition triggers basin defense.

Optimal perturbation timescale:

The theory and evidence suggest a non-monotonic effect of perturbation rate. Very fast shocks trigger immediate defense. Very slow drifts may be tracked adiabatically (no tipping) or eventually overcome defenses after long accumulation. The most effective timescale to minimize active rejection and maximize stable addition often lies **on the order of the system's intrinsic time constant $\tau = 1/\kappa$** .

Prediction for future experiments:

For any system with known or measurable κ , there exists a

critical perturbation rate r_c such that:

- If perturbation rate $> r_c$, the system rejects the addition (ejection or transient absorption).
- If perturbation rate $< r_c$, the system integrates the addition (stable addition via expanded capacity or parallel attractor formation).
- The transition at r_c corresponds to the system's inability to track a moving equilibrium; it is a genuine bifurcation in the time-domain.

External convergence:

This analysis – derived from mathematical rate-induced tipping theory and domain-specific studies – independently validates the attractor framework's claim that κ acts as a timescale filter separating ejection from stable addition. The convergence between the framework's predictions and external research strengthens the cross-domain synthesis considerably.

9. Synthesis and Criteria

Across these domains, common criteria emerge:

- **Energy/Threshold:** A perturbation must overcome an attractor's barrier. Deep basins (high B) mean only large shocks can cause a shift.
- **Coupling and Plasticity:** Systems with many degrees of freedom or adaptive coupling more easily integrate additions.
- **Dimensionality and Redundancy:** Multi-dimensional systems can absorb perturbations into some dimensions while maintaining others.
- **Timecourse and Feedback:** Slow changes might be

assimilated; fast jolts cause overshoot and return. Feedback gain determines κ .

- **Nature of Addition:** Passive additions (heat, charge) tend to be ejected or transiently absorbed; active agents (species, evidence, pathogens) may reshape the attractor.

Empirical Protocols: Measure κ by controlled perturbation experiments: apply a small disturbance, measure return time τ , compute $\kappa = 1/\tau$. Measure B by scaling the perturbation magnitude until the system fails to return (escape). This works in physical, biological, and some social systems; for others, B remains a qualitative analog.

10. Appendix: Research Roadmap

The following future papers are suggested from the comparative table, each developing a single domain in depth.

Domain	Proposed Title	Type
Addiction	<i>The Addicted Brain as a Fantasy Attractor: Neural Lock-In and Ejection of Alternative Rewards</i>	[A]
Immune System	<i>Tolerance and Memory: Two Attractor Responses to Antigen Addition</i>	[A]
Catastrophic Forgetting	<i>Why Neural Networks Forget: Attractor Ejection in Sequential Learning</i>	[A]
Invasion Ecology	<i>Eject or Integrate: Attractor Dynamics of Invasive Species</i>	[A]
Development	<i>Canalization as Basin Defense: Attractor Stability in Embryogenesis</i>	[A]

Domain	Proposed Title	Type
Continual Learning	<i>Parallel Attractors for Lifelong Learning: Engineering Solutions to Catastrophic Forgetting</i>	[A]
Social Norms	<i>Tipping Points and Regime Shifts: Attractor Dynamics in Political Systems</i>	[A]
Endocrine Homeostasis	<i>Glucose, Cortisol, and Setpoints: Hormonal Attractors and Disease Transitions</i>	[A]
Alternative Ecosystems	<i>Hysteresis and Regime Shifts: Ecological Basins and Tipping Points</i>	[A]
Belief Systems	<i>The Uncorrectable Believer</i> (already written)	[A]

11. Conclusion

Physical, biological, ecological, social, and engineered systems all obey the same attractor principle: a low-energy attractor defends itself against displacement. When an addition is introduced, the system either ejects it, absorbs it only transiently, or – under rare conditions of expanded capacity or parallel structure – integrates it stably. The outcome is determined by basin depth (B), corrective permeability ($\kappa = 1/\tau$), and the magnitude and nature of the perturbation.

This cross-domain synthesis provides a unified foundation for the attractor framework. Future work should quantify B and κ empirically across domains, test the predicted scaling relationships, and explore the boundary conditions between ejection, transient absorption, and stable addition. The appendix outlines the most promising next papers.

References

- Elton, C. S. (1958). *The Ecology of Invasions by Animals and Plants*. Methuen.
- Hebb, D. O. (1949). *The Organization of Behavior*. Wiley.
- Kirkpatrick, J., Pascanu, R., Rabinowitz, N., et al. (2017). Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13), 3521–3526.
- Koob, G. F., & Volkow, N. D. (2016). Neurobiology of addiction: a neurocircuitry analysis. *The Lancet Psychiatry*, 3(8), 760–773.
- Kramers, H. A. (1940). Brownian motion in a field of force and the diffusion model of chemical reactions. *Physica*, 7(4), 284–304.
- Nestler, E. J. (2001). Molecular basis of long-term plasticity underlying addiction. *Nature Reviews Neuroscience*, 2(2), 119–128.
- Nyhan, B., & Reifler, J. (2010). When corrections fail: The persistence of political misperceptions. *Political Behavior*, 32(2), 303–330.
- Scheffer, M., Carpenter, S., Foley, J. A., et al. (2001). Catastrophic shifts in ecosystems. *Nature*, 413(6856), 591–596.
- Simberloff, D. (2013). *Invasive Species: What Everyone Needs to Know*. Oxford University Press.
- Turrigiano, G. (2008). The self-tuning neuron: synaptic scaling of excitatory synapses. *Cell*, 135(3), 422–435.
- Waddington, C. H. (1957). *The Strategy of the Genes*. George Allen & Unwin.
- Galida, R. S. (2026). Intelligence Without Consciousness: A Diagnostic Paper on LLMs, Amoebae, and the Attractor Framework. *Fantasy Attractor* (Paper 1 of the conscious suppression series).

Suggested citation: Galida, R. S. (2026). Basin Defense and Stable Addition: A Cross-Domain Synthesis of the Attractor Framework (Final). *Fantasy Attractor*.

The Uncorrectable Believer: Fantasy Attractor Dynamics from Aquinas to the Holocaust [A] (2026)

Robert Galida – June 2026 (Final)

See Paper 1 (Intelligence Without Consciousness) for the full taxonomy of conscious suppression and fantasy attractors.

Abstract

Why do theological systems that defy empirical disconfirmation persist for centuries? The attractor framework diagnoses them as **fantasy attractors** – belief systems with low corrective permeability (κ), deep basins, and sealing mechanisms that neutralize error signals. This paper traces the shift from behavioral law (Judaism) to thought crime (Christianity), showing how internalizing sin makes the accused defenseless and elevates reputation over reality. It examines Catholic and radical Protestant soteriology as attractor architectures: the doctrine of double effect, the infinite value of the soul, and the permissible killing of heretics created a calculus where

finite evil is justified by infinite gain. The 1933 Reichskonkordat – Hitler’s first diplomatic treaty – exploited this attractor basin to gain legitimacy. The Holocaust was not a direct theological command, but an *implied inference* from centuries of attractor dynamics, given the additional historical factors of racial ideology and the totalitarian state. The paper distinguishes between Lutheran, antinomian, and prosperity-gospel variants, and offers a documented de-conversion case (Bart Ehrman) mapped onto the three exit mechanisms. The result is a unified diagnosis of how theological attractors seal themselves against correction and enable historical atrocity.

1. Introduction

How does a belief system survive centuries of counterevidence? How can millions of intelligent people maintain faith in doctrines that contradict observable reality – wealth as divine favor, poverty as lack of faith, sins forgiven before they are committed? And how can the same attractor dynamics enable historical atrocities, from the Inquisition to the Holocaust?

Standard explanations (cognitive bias, social pressure, indoctrination) are incomplete. Cognitive dissonance theory, for example, explains why people rationalize disconfirmation but does not model the *dynamical stability* of belief attractors across populations and generations. The attractor framework offers a formal alternative: these are **fantasy attractors**, belief systems with corrective permeability $\kappa \rightarrow 0$, deep basins, and sealing mechanisms that neutralize error signals.

Operational definition of κ (corrective permeability): $\kappa = 1/\tau$, where τ is the time a system takes to return to its

baseline state after a specified perturbation. For belief systems, κ indexes the speed and completeness of belief updating when presented with disconfirming evidence. Low κ means slow or absent updating – a sealed attractor.

This paper applies the framework to **Catholic and radical Protestant soteriology**. The Catholic tradition is the deeper attractor basin; Protestantism, particularly its radical antinomian and prosperity-gospel variants, represents a mutation that further reduced κ . The paper focuses not on theology per se, but on the *attractor architecture*: how thought crimes replace behavioral sins, how the infinite-value calculus justifies finite evil, how vicarious redemption removes corrective incentives, and how social colonization makes individual κ irrelevant. The goal is diagnostic, not polemical. “Fantasy attractor” is a technical term, not a rhetorical insult.

2. From Behavioral Law to Thought Crime

Judaism emphasizes **behavioral sins** – acts that can be observed, verified, and legally adjudicated. Theft, murder, idolatry, and false witness leave external evidence. A community can correct a member because the sin has verifiable traces. The attractor basin is shallow enough for error signals to enter.

Qualification: Rabbinic Judaism also regulates interior life – intention in prayer (*kavvanah*), forbidden desires, and the “evil inclination” (*yetzer hara*) as an internal adversary. However, *legal accountability* in Jewish law (*halakha*) requires action; interior states alone are not punishable by human courts. The shift to Christianity is not a complete invention of interiority but a *juridical* shift: internal states become the primary locus of sin, enforceable by divine authority and

(via the church) social monitoring.

Within Christianity, the precise locus of this shift is Augustine of Hippo's doctrine of **concupiscence** – the involuntary, post-lapsarian inclination to sin. Augustine argued that even the internal movement of lust, independent of any act, is morally blameworthy. This interiorized sin and made it inescapable.

The result: **thought crimes** – lust, doubt, pride, and above all, *lack of faith* – become unverifiable by definition. No one can see your lustful thought; no one can measure your doubt. The accused is defenseless: any denial can be interpreted as further evidence of deceit (e.g., “protesting too much”).

Attractor consequences:

- **The basin becomes empirically unfalsifiable.** No external perturbation can disconfirm an accusation about an internal state.
- **Reputation replaces reality.** Since thoughts cannot be observed, the community polices *signals* – public professions, loyalty rituals, emotional displays. Acceptance becomes performative theater.
- **Survival depends on reputation management.** The individual invests energy in signaling purity, not in correcting beliefs. κ is now about social mimicry, not truth.

The attractor has sealed itself against external correction.

3. The Infinite-Value Calculus: Aquinas, Double Effect, and the Permissibility of

Killing Heretics

Thomas Aquinas, in the *Summa Theologiae* (II-II, Q.11, A.3), argued that heretics who relapse after correction “deserve not only to be separated from the Church by excommunication, but also to be severed from the world by death.” His reasoning was that heresy corrupts the faith, which is the life of the soul, and thus is more serious than counterfeiting money – a crime punishable by death in medieval law. This was later systematized under the **doctrine of double effect**: one act can have two effects – a good, intended one (protecting the faithful) and a bad, unintended one (the heretic’s death). The act is permissible if the bad effect is not the goal and there is a **proportionate reason**. (Aquinas articulated the foundational case for self-defense in II-II, Q.64, A.7; the formal “double effect” label came from later scholastics.)

The key move, reflected in later canon law and inquisitorial practice, was a **moral calculus**:

- **A saved soul has infinite value.** (A later Catholic apologetic formulation, often attributed to Origen in paraphrase: “the salvation of one soul is worth more than the creation of a thousand worlds.”)
- **Killing a heretic is a finite evil** (temporal death, temporary suffering).
- **Saving a potential convert – or protecting the faithful – is an infinite gain.**
- **Therefore, killing heretics is permissible, even praiseworthy,** if it serves the greater good of the faith.

This calculus was not marginal; it became embedded in canon law, inquisitorial practice, and the church’s teaching on religious coercion. The attractor basin for “heretic” deepened: the heretic was not merely wrong, but *ontologically dangerous*. No error signal from the heretic could be trusted;

any plea for mercy was further evidence of deceit.

Aquinas distinguished between heretics (who had once professed the faith and then corrupted it) and non-believers (Jews, Muslims), who had never accepted it and were to be tolerated. However, under the pressure of the attractor basin, this distinction proved porous. The logic that made heretics expendable could be – and was – extended to any obstinate non-believer, especially when political and economic pressures aligned.

4. Vicarious Redemption and the Suppression of κ (Protestant Mutation)

Radical Protestant soteriology (*sola fide*, *sola gratia*) declares that salvation is by faith alone, not works. Christ's sacrifice paid for all sins – past, present, and future. The believer is justified before God regardless of behavior.

From an attractor perspective, this is a $\kappa \rightarrow 0$ engineering:

- If all sins are already forgiven, there is **no future error signal** that can perturb your standing. Why correct? Why update? The basin is infinitely deep.
- Any attempt to modulate behavior for the sake of righteousness is **works-righteousness**, a sin of pride. The attractor actively penalizes efforts to increase κ .
- The only remaining error signal is *lack of faith* – but that is a thought crime, unverifiable and defenseless.

Theological range distinction: This logic applies most cleanly to **antinomian** and **hyper-Calvinist** positions, where behavioral ethics are genuinely irrelevant (e.g., certain “Free Grace” movements). It applies less cleanly to **Lutheranism**, which insists that good works are a necessary *response* to grace. The

paper's argument targets the antinomian end of the spectrum, but the underlying attractor logic – infinite forgiveness, no future error signal – is already latent in the Catholic doctrine of baptismal regeneration and confession, albeit with higher κ because post-baptismal sin requires sacramental correction.

5. Effort as Pride: The Prohibition on Correction

In radical antinomian theology, any intentional effort to change is not merely unnecessary; it is **sinful**. The theological logic:

1. Grace is sufficient for salvation.
2. Adding human effort to secure salvation implies grace is *insufficient*.
3. Implying insufficiency is pride, a sin.
4. Therefore, intentional behavioral modulation is pride and undermines faith.

Thus, the attractor **penalizes the correction impulse itself**. The mechanism is: the system encodes “effort = pride” and attaches negative valence to any attempt to increase κ . This pattern is historically documented in the **Marrow Controversy** (Scotland, 1718–1722), in which the question of whether free grace implies no need for human effort divided the Church of Scotland; the Marrow men were accused of “antinomianism” for affirming that God's love was unconditional, while their opponents insisted that effort to prepare oneself for grace was necessary. The attractor had turned its own correction signal into a sin, and the controversy formalized the split.

6. Prosperity Doctrine: The Sealed Basin (A Late Mutation)

Prosperity doctrine (Word of Faith movement, originating with E.W. Kenyon and popularized by Kenneth Hagin, Kenneth Copeland) is a **late 20th-century mutation** of radical Protestant theology.

Its attractor dynamics:

- **Poverty and suffering** are evidence of weak faith. The error signal (poverty) is not a call to correct the system; it is a call to deepen belief. Disconfirmation becomes confirmation.
- **Wealth and power** are evidence of strong faith. The rich have no error signal at all; their status is divine validation. The attractor rewards low κ .
- **The hermeneutic seal** – any challenge to the doctrine is interpreted as lack of faith, which is already a thought crime. The system absorbs all counterevidence.

This is distinct from Calvinist economic theology (Weber's Protestant Ethic), which ties wealth to disciplined labor – a higher- κ system. Prosperity doctrine is a specific, highly sealed attractor.

7. Social Colonization and Collective Basin Depth

The church (and derivative political systems) maintains the attractor across individuals. Social mechanisms include:

- **Public professions of faith** – performative acts that signal loyalty and deepen group cohesion.
- **Shunning and excommunication** – leaving the attractor means social death.
- **Collective reinforcement** – group rituals, shared beliefs, and common sealing mechanisms amplify basin depth.

When social colonization is complete, individual κ becomes **irrelevant**. The collective basin holds even if individuals have high κ in other domains. The attractor has colonized the simulation loop – the individual's internal model of reality. Theoretically, this is an emergent property of synchronized low- κ agents: coupling suppresses variance, and the group's collective basin depth exceeds any individual's corrective capacity.

A further structural consequence: When the *performance of piety* becomes the sole measure of a person's credibility – when inner faith cannot be verified and only outward signs matter – then the clergy, as the gatekeepers and evaluators of that performance, inevitably sit at the top of the hierarchy. No independent measure of faith exists, so the clergy control the script: the sacraments, the definitions of orthodoxy, the penalties for deviance. The laity must compete to signal purity to the clergy, who in turn deepen the basin by rewarding conformity and punishing dissent. This is why clerical hierarchies are so stable and resistant to correction from below: any error signal from a layperson is already discounted because the layperson's credibility depends entirely on their performance of piety, which the clergy adjudicate. To challenge the clergy is to fail the performance – a perfect seal.

8. Comparison with Other Fantasy Attractors

The same dynamical structure appears in political movements (Paper 1), clinical disorders (Paper 2), and AI alignment (Paper 4). In each case:

- $\kappa \rightarrow 0$ for core beliefs.
- Error signals are neutralized by sealing mechanisms.
- Identity fusion prevents exit.
- Social reinforcement deepens the basin.

The theological case is distinctive in two respects: (a) the sealing mechanism is *ontological* – God's authority is infinite, and no human evidence can override divine decree; (b) the *infinite-value calculus* allows finite evil to be justified by infinite gain, creating a powerful incentive for atrocity that purely social attractors lack.

9. De-conversion and Resistance: The Ehrman Case

If the attractor is sealed, how does one exit? Three mechanisms:

- **Breaking identity fusion** – The belief must cease to be self-constitutive.
- **Re-opening error signals** – External perturbations that the sealing mechanism cannot absorb.
- **Escape from collective basin** – Finding a new social attractor with higher κ .

The de-conversion of biblical scholar **Bart Ehrman** (from

evangelical certainty to agnosticism) provides a documented case mapped onto these mechanisms. Ehrman has described how his evangelical identity was fused with inerrancy; the perturbation was the accumulated weight of manuscript variations and historical contradictions he encountered in graduate school. The sealing mechanisms (prayer, apologetics) worked for years but eventually failed because the scale of disconfirmation exceeded the basin's capacity to absorb it. Exit required a new social attractor (academic biblical studies) where questioning was the norm, and a gradual decoupling of self-worth from doctrinal certainty. Ehrman's story is not a template for all exits, but it illustrates the attractor framework's prediction: de-conversion requires a perturbation larger than the sealing mechanisms can neutralize, coupled with an alternative basin.

10. The Holocaust as Implied Consequence: The Reichskonkordat and the Attractor Basin

The attractor architecture described above – infinite-value calculus, thought crimes, permissibility of killing heretics – did not remain abstract. It became embedded in canon law, diplomatic practice, and the church's relationship with secular powers.

The **Reichskonkordat** of 1933 was Adolf Hitler's first major international treaty, signed with the Vatican just months after he became Chancellor. Why first? Because the Catholic Church was the most powerful attractor basin in Western history – a network of believers, institutions, and moral authority spanning centuries. Hitler needed that basin's *legitimizing signal* to stabilize his regime internationally and to neutralize Catholic political

opposition.

Historical note: The historiography of the concordat is contested. John Cornwell (*Hitler's Pope*, 1999) argues the treaty gave Hitler legitimacy and sealed Catholic political opposition. Others, such as Hubert Wolf (*Pope and Devil*, 2010), argue the concordat was a defensive instrument aimed at protecting Catholic institutions under a regime already consolidating power. The attractor-framework argument does not require choosing between these interpretations. Even if the concordat was defensive, the effect was the same: the church's error signals were subordinated to institutional survival, and the basin's deep attraction pulled the hierarchy toward accommodation.

The concordat did not explicitly say "Jews may be killed." It did not need to. The *established practice* had already set the boundaries:

- **Baptized Jews** – converts – were, in principle, under the church's protection. Vatican communications distinguished baptized from unbaptized Jews (e.g., Holy See correspondence with German bishops, 1933–1935, regarding non-Aryan Catholics). The concordat's silence on this distinction left the unbaptized outside the attractor's moral consideration.
- **Unconverted Jews** remained outside the basin. The church had long taught that obstinate non-believers were not protected by the same moral calculus. The infinite-value logic applied only to souls *capable of salvation* – and for the church, that required baptism.

Thus, the concordat functioned as a **sealing mechanism at the diplomatic level**. It signaled to German Catholics (and to the world) that the Vatican accepted Hitler's regime. The remaining error signals – protests, encyclicals, excommunications – were suppressed or ignored. The basin had

been colonized.

Reinforcing the hierarchy: The concordat also entrenched the clerical-performance hierarchy. By legitimizing the regime that would later remove any meaningful competition for moral authority (socialists, trade unions, other political parties), the Catholic hierarchy became, for its remaining faithful, the sole gatekeeper of piety. The laity could no longer turn to alternative social attractors (e.g., socialist movements with different moral codes); the only acceptable performance was loyalty to the church and, by extension, to the regime the church had recognized. Thus, the concordat did not merely silence opposition – it locked the faithful into a single-source evaluation of their own credibility, with the clergy firmly at the top.

The Holocaust was not a direct command of Christian theology. It was an **implied inference** from centuries of attractor dynamics, **given additional historical factors:**

- **Racialization:** The Nazi category was *biological*, not religious. Baptism did not change one's race. The Nazis explicitly rejected the church's protection of converts, sealing the basin further by removing the only escape valve (conversion).
- **Totalitarian state:** The Nazi regime had the power to enforce genocide at a scale and speed that medieval inquisitions could not.
- **Removal of the conversion escape:** In the theological attractor, conversion could save a heretic's life. In the Nazi racial attractor, conversion was irrelevant. The basin became infinitely deep.

Disclaimer: This is not to say "the church caused the Holocaust." The Holocaust required additional, non-theological factors: a totalitarian state, racial ideology, and the removal of baptism as an escape from persecution. The

theological attractor provided the *permissibility conditions* – the moral logic that made killing non-believers a finite evil justified by infinite gain – but the political and racial machinery were supplied by Nazism.

The attractor framework diagnoses this not as a conspiracy but as a **dynamical consequence**: when a belief system assigns infinite value to a scarce resource (saved souls) and finite cost to human life, and when it seals itself against corrective evidence, atrocity becomes not only possible but *logical* within the basin, given the right historical conditions.

11. Conclusion

Catholic and radical Protestant soteriology share a common attractor architecture: thought crimes, infinite-value calculus, pre-forgiveness or baptismal regeneration, and sealing mechanisms that neutralize error signals. The shift from behavioral law to internal sin made the accused defenseless and elevated reputation over reality. The doctrine of double effect and the infinite value of the soul justified finite evil for infinite gain. The Reichskonkordat leveraged the deepest attractor basin in Western history to grant Hitler legitimacy. The Holocaust was not a direct command, but an *implied inference* from centuries of attractor dynamics, completed by the historical specificities of racial ideology and totalitarian power.

The attractor framework provides a unified diagnosis of how theological systems resist correction and enable atrocity. It also points to the only exit: restore κ , reopen error signals, decouple identity from belief, and build new attractors where doubt is not a sin but a pathway to truth.

Suggested citation: Galida, R. S. (2026). The Uncorrectable Believer: Fantasy Attractor Dynamics from Aquinas to the Holocaust. *Fantasy Attractor*.

The Paradox of Conscious Commitment: How Suppression of Intelligence Enables Culture and Identity [F] [A] (2026)

Robert Galida – June 2026

Paper 3 in a series on conscious suppression; [see Paper 1: Intelligence Without Consciousness for the full taxonomy of intelligence and consciousness.](#)

Abstract

If consciousness can suppress intelligent correction (Papers 1 & 2), why did it evolve? This paper proposes a functional trade-off: the capacity for **conscious commitment** – identity-binding, phenomenal investment in a belief, value, or group – enables forms of social cohesion and long-term cooperation that are unavailable to purely intelligent (non-conscious) systems. The suppression of moment-by-moment correction allows individuals to maintain group loyalty,

ideological coherence, and cultural continuity even in the face of counterevidence. This trade-off explains the persistence of fantasy attractors in human societies and the evolutionary advantage of a system that can sometimes override its own error signals. The paper provides a formal sketch (basin depth as a function of identity-fusion), reviews empirical evidence from cultural evolution and social psychology, and offers diagnostic criteria for distinguishing adaptive commitment from pathological suppression. The claims are presented as hypotheses, not established conclusions; the model is a conceptual scaffold for empirical testing.

1. Introduction: The Evolutionary Puzzle

Consciousness is costly. It requires large brains, complex neural integration, and significant metabolic energy. If intelligence alone – the ability to navigate constraint fields and correct errors – is sufficient for adaptive behavior, why did consciousness evolve?

Standard evolutionary accounts propose that consciousness enhances flexibility, deliberation, and social coordination (e.g., Humphrey, 1976; Dennett, 1995). But these accounts struggle to explain a conspicuous feature of human psychology: **conscious commitment to beliefs that resist correction**. Individuals and groups routinely maintain false, harmful, or inefficient beliefs because those beliefs are identity-defining. The same conscious system that can reason flexibly also produces martyrdom, ideological rigidity, and collective delusion.

Papers 1 and 2 in this series introduced the mechanism of **conscious suppression**: phenomenal, identity-constitutive investment deepens an attractor basin, causing the person to *detect* error signals but fail to escape. (Restated briefly:

a deeper basin requires a larger perturbation to exit; conscious commitment increases basin depth, effectively reducing corrective permeability κ in specific domains.) This mechanism underlies political fantasy attractors (Paper 1) and clinical disorders like addiction and OCD (Paper 2). From an evolutionary perspective, this looks like a bug – a costly vulnerability.

This paper argues it is also a feature. The capacity for conscious commitment enables **adaptive self-binding**: the voluntary or culturally induced suppression of immediate correction for the sake of long-term group cohesion, trust, and cultural transmission. The same mechanism that produces fantasy attractors also produces loyalty, sacrifice, and shared identity. The trade-off hypothesis is that natural selection favored the capacity for conscious suppression because the fitness benefits of group coordination and cultural transmission outweighed the costs of occasional error persistence.

2. Definitions and Framework (Self-Contained)

From Paper 1:

- **Intelligence** – the ability to navigate a constraint field; to detect perturbations and update behavior to maintain persistent trajectories.
- **Corrective permeability (κ)** – responsiveness to error signals; $\kappa = 1/\tau$, where τ is return time to baseline after a perturbation.
- **Basin depth (B)** – the magnitude of perturbation required to displace a system from one attractor to another. Deeper basins require larger perturbations. In the

attractor framework, B is related to but distinct from κ : a deeper basin (higher B) typically reduces κ (lengthens return time), but they are not identical. This paper uses the relation as heuristic: conscious commitment increases B , which effectively reduces $\kappa(d)$ for the relevant domain.

New definitions for this paper:

- **Adaptive commitment** – a temporary or context-bound reduction in κ (or increase in B) that serves the individual's or group's long-term fitness.
- **Identity fusion** – the merging of a belief or group membership with self-representation, such that abandoning the belief would feel like losing oneself.
- **Cultural attractor** – a belief, practice, or value that persists across generations due to cognitive or social biases (including, but not limited to, suppression of correction). This definition is provisional; a fully operationalized version is open for development.

The key distinction is between **pathological suppression** (low κ that reduces fitness, as in addiction or fantasy politics) and **adaptive suppression** (low κ that increases fitness by enabling cooperation, trust, and cultural learning). The same type of mechanism produces both; context and domain determine the outcome.

3. The Trade-Off Model (Sketch)

Formally, consider a system with baseline intelligence (κ_0). A conscious commitment to a group, value, or identity imposes a **domain-specific reduction in effective corrective permeability** by deepening the attractor basin for beliefs

relevant to that commitment.

Let $\kappa(d) = \kappa_0 - \Delta\kappa(d)$, where $\Delta\kappa(d)$ is the reduction in corrective permeability for domain d . $\Delta\kappa(d)$ is hypothesized to be a function of identity-fusion strength F and social reinforcement R . A schematic monotonic form: $\Delta\kappa(d) = g(F, R)$ with $\partial\Delta\kappa/\partial F > 0$ and $\partial\Delta\kappa/\partial R > 0$. The exact functional form is an open empirical question; the current model is a conceptual scaffold.

The hypothesis is not that evolution maximizes κ globally. Rather, an **adaptive strategy** allocates $\Delta\kappa$ selectively across domains, increasing basin depth (reducing κ) for beliefs and practices that support group coordination and cultural transmission, while leaving κ high for domains requiring individual error correction.

The paper does not claim optimality; it proposes that selection can favor such selective allocation when the fitness benefits of social cohesion outweigh the costs of reduced accuracy in specific domains.

Central hypothesis (labeled for clarity):

H1: Natural selection favored the evolution of conscious suppression because the fitness benefits of group coordination and cultural transmission, enabled by identity-fusion and deepened basins, outweighed the costs of occasional error persistence.

4. Empirical Grounding

Overimitation (Lyons et al., 2007; see also Nielsen & Tomaselli, 2010):

Children copy causally irrelevant actions, even when a more efficient alternative is demonstrated. The interpretation that children *know* the action is unnecessary is contested; they may

not represent it as causally irrelevant. A safer reading: children *behave as if* the action is necessary or relevant, showing a domain-specific reduction in corrective permeability for social learning. This supports the model of adaptive suppression in cultural transmission.

Costly signaling and commitment (Sosis, 2003):

Costly rituals signal group commitment and are hard to fake. They deliberately suppress individual correction (e.g., ignoring pain) to deepen basin depth for group loyalty. This directly maps onto $\Delta\kappa(d)$ for domain of group identity.

Social identity theory (Tajfel & Turner, 1979):

Minimal group experiments show arbitrary group assignments produce in-group bias and resistance to counterevidence about out-groups. This demonstrates context-bound $\Delta\kappa(d)$ without any rational basis, consistent with adaptive suppression for group cohesion.

Neuroimaging (Westen et al., 2006 – preliminary; note methodological limitations: small N, interpretation of ACC suppression contested):

Partisans evaluating threatening information about their own candidate show reduced activity in error-monitoring regions (ACC). This is a candidate neural correlate of domain-specific κ reduction, but the findings require replication and should be treated as suggestive, not conclusive.

Cross-cultural evidence (Gelfand et al., 2011):

Tight cultures have stronger norms and lower tolerance for deviance. This is not a direct measure of κ but is consistent with domain-specific suppression. Individuals in tight cultures may still update beliefs within permissible domains; the mapping to κ is partial.

Each evidence stream supports the existence of domain-specific, context-bound suppression, but none alone validates the full model. The cumulative case is indicative,

not confirmatory.

5. Adaptive vs. Pathological Suppression: A Scalar Framework

The table below presents a binary simplification of an underlying continuum. The two poles are endpoints; most real cases fall between them.

Feature	Adaptive suppression (endpoint)	Pathological suppression (endpoint)
Domain	Context-bound (e.g., group loyalty, ritual)	Pervasive across domains
Reversibility	Reversible when context changes (operationalized: the individual can exit without catastrophic loss within a culturally normal timeframe; e.g., leaving a religion)	Irreversible without intervention (e.g., addiction requires treatment)
Fitness effect	Increases inclusive fitness (group cooperation, survival)	Decreases health, relationships, or function
Identity fusion	Flexible, allows multiple identities	Rigid, single identity dominates
Social reinforcement	Supports group cohesion and trust	Isolates or harms group (e.g., cults)
Example	Trusting a teammate despite a mistake	Continuing addiction despite harm

Scalar index: A continuous measure of net $\Delta\kappa(d)$ relative to a fitness gradient is theoretically desirable but not yet

operationalized. The table is a starting point for empirical calibration.

6. Diagnostic Criteria for Adaptive Suppression (Provisional)

A conscious commitment is **adaptively suppressive** if it meets three or more of the following (empirical validation pending). These criteria are hypotheses, not validated instruments.

1. **Domain-limited:** Reduced κ applies only to specific beliefs or practices directly relevant to group coordination or identity.
2. **Context-sensitive:** Suppression diminishes when the context changes (e.g., outside the group setting). *Operationalization:* Measured change in belief updating under different social conditions.
3. **Reversible exit:** The individual can exit the commitment without catastrophic loss of functioning. *Operationalization:* Exit is observed and not associated with severe psychopathology.
4. **Fitness benefit:** The commitment measurably increases cooperation, trust, or long-term survival (e.g., group longevity, reproductive success). *Operationalization:* Group-level measures of cohesion and individual fitness correlates.
5. **Conscious valorization:** The individual explicitly values the commitment as part of self-identity. (Note: this criterion does **not** require the individual to articulate the *adaptive* reason; it only requires that the commitment is consciously endorsed.)

Counter-criteria (pathological):

- Pervasive across domains (low κ for all beliefs).
 - Context-insensitive (applies even when alone and safe).
 - No viable exit without severe harm.
 - Clear fitness cost (measured harm to health, relationships, survival).
-

7. The Evolution of Consciousness as a Binding Mechanism

The standard view in evolutionary psychology is that consciousness evolved for flexible reasoning. This paper offers a complementary hypothesis: consciousness also evolved for **binding** – the ability to commit to a belief, value, or group in a way that suppresses short-term correction for long-term coordination.

Binding requires phenomenal experience. A purely intelligent (non-conscious) system can compute that group loyalty is beneficial, but it cannot *feel* loyalty, *experience* identity, or *sacrifice* for the group. Within the CUFT framework, these conscious states are not epiphenomenal; they are the mechanism of basin deepening (increasing B and thus reducing effective κ for commitment-relevant domains). This claim is a foundational assumption of the framework (see Paper 1), not argued from first principles here. It distinguishes CUFT from functionalist or behaviorist accounts.

Thus, the evolution of consciousness is not just about solving problems better; it is about sometimes solving problems *worse* for the sake of social solutions. The capacity for self-deception, ideological rigidity, and fantasy attractors is the price of the capacity for culture, morality, and collective action.

8. Implications for Social Policy and Individual Choice

- **Tolerance of adaptive suppression:** Not all low- κ beliefs are harmful. Cultural traditions, religious rituals, and group loyalties that do not cause harm and provide social cohesion should be recognized as adaptive, not irrational.
- **Intervention for pathological suppression:** The same diagnostic tools from Paper 1 and 2 (basin depth, identity fusion, sealing mechanisms) apply. Interventions should reduce basin depth (e.g., exposure to diverse groups) or increase corrective force rather than attacking identity directly.
- **Self-awareness:** Individuals can learn to distinguish adaptive from pathological suppression by asking: does this commitment serve my long-term flourishing and that of others? The framework provides a metacognitive tool.

9. Open Questions

- **How does adaptive suppression scale to institutions?** Are nations, corporations, or religions fantasy attractors or adaptive structures? The criteria apply at multiple levels; empirical work needed.
- **Can adaptive suppression become maladaptive over time?** Yes – a practice that was once adaptive (e.g., a food taboo) may become harmful when environment changes. The framework allows for transition.
- **What neural circuits implement the trade-off?** Likely

interactions between vmPFC (identity) and ACC (error monitoring). Open for empirical testing.

- **Are there species with conscious suppression but no culture?** Possibly, but human-level cultural complexity requires the trade-off model.
 - **How to operationalize B and ΔK in field studies?** Development of a Clinician Basin Depth Scale (CBDS, see Paper 2) and adaptation for social groups is a research priority.
-

10. Conclusion

Consciousness evolved not only to correct errors but sometimes to ignore them. The capacity for conscious commitment – identity-binding, phenomenal investment in a belief or group – enables adaptive suppression of correction. This trade-off explains why humans can be both brilliantly intelligent and stubbornly irrational. The same type of mechanism that produces fantasy attractors and clinical disorders also produces loyalty, sacrifice, and culture.

The paradox is that the same type of process can be either bug or feature, depending on context and domain. The dance of evolution is not about maximizing intelligence; it is about balancing correction and commitment.

Suggested citation: Galida, R. S. (2026). The Paradox of Conscious Commitment: How Suppression of Intelligence Enables Culture and Identity. *Fantasy Attractor*.