

Non-Physical Claims Are Fantasy Attractors: Why Unverifiable Realms Cannot Be Empirically Distinguished from Nonexistence

Robert Galida – June 2026

[F] (Foundation)

Abstract

The attractor framework adopts a physicalist commitment: to be real is to be able to interact, and to interact is to share at least one **interaction channel** (spacetime, energy, momentum, gauge charge, or any measurable coupling). This is a philosophical starting point, not an empirical discovery. The paper argues that any claim about a non-physical realm – defined as having no such interaction channel – cannot be empirically assessed. Such claims are **fantasy attractors**: belief systems structurally sealed against correction by defining their objects as forever beyond any possible test. The paper distinguishes provisional non-detection (e.g., dark matter) from **structural, permanent non-verifiability** (e.g., non-physical gods, transcendent souls). It concludes that while such claims may have personal or social meaning, they cannot be part of a scientific ontology, and their structure makes them vulnerable to fraud and manipulation – though sincere belief is not fraud.

1. The Foundational Commitment: Interaction Requires Shared Channels

The attractor framework is a physicalist ontology. It begins with a commitment: **entities can only interact through shared interaction channels**. An *interaction channel* is any measurable coupling – spacetime coordinates, energy, momentum, electric charge, weak isospin, color charge, or any other quantity that can be transferred or correlated between systems. This is not an empirical discovery of the Standard Model; it is the framework's chosen criterion for what counts as real.

The neutrino example illustrates the criterion but does not prove it. Neutrinos interact weakly because they share weak isospin; they do not interact electromagnetically because they lack electric charge. The framework simply says: if an entity shares no interaction channel with physical reality, we have no way to detect it, measure it, or include it in a scientific ontology. That is a philosophical choice, not a falsifiable claim about the world.

Why interaction? Interaction is chosen because it provides a public, corrigible basis for knowledge. It avoids ontological commitments that cannot influence observation, and it aligns with the core principle of the attractor framework: *persistence under perturbation*. An entity that never perturbs anything cannot be distinguished from nothing.

What the framework does not claim:

- That non-physical entities are logically impossible.
- That all non-physical claims are false.
- That physics has disproven God or the supernatural.

What it does claim:

- That non-physical entities cannot be empirically distinguished from nonexistence.
 - That claims about them operate as fantasy attractors, resistant to correction.
-

2. Types of Non-Physical Claims

A non-physical claim is any assertion about an entity, force, or realm defined as having **no interaction channel** with the physical world. However, not all claims that seem non-physical are alike. We distinguish two categories:

Category A: Truly non-interacting – Claims that explicitly deny any possible interaction. Examples:

- A deistic creator who wound the universe and then never interacts.
- A transcendent God defined as beyond all categories, including causality.
- An immaterial soul that cannot influence the body after death.
- Abstract objects (Platonism) that exist non-physically and non-causally.

Category B: Claims that assert interaction but evade testing – Examples:

- Ghosts that move objects but become undetectable when instruments are present.
- Psychics whose powers fail under controlled conditions (explained as “skeptic’s energy”).
- Homeopathic “water memory” that cannot be detected by any known physical measurement.

Category B is a different epistemic pathology: motivated reasoning, ad-hoc escape clauses, and sealing mechanisms. The attractor framework addresses them as *functionally* non-verifiable in practice, but they are not the primary target of this paper. This paper focuses on **Category A**: claims that structurally preclude any possible interaction channel.

Domain (Category A)	Example Claim	Interaction Channel?	Empirically Assessable?
Religion (non-interacting God)	A creator with no detectable properties	None	No – any test is ruled out a priori
Paranormal (non-interacting ghosts)	Ghosts that cannot affect matter	None	No – no possible evidence
Abstract objects (Platonism)	Numbers exist non-physically, non-causally	None	No – no interaction, hence no evidence
New Age (non-interacting “vibrations”)	Crystals with undetectable healing vibrations	None	No – absence of effect is blamed on “wrong intent”

Under the framework’s commitment, such claims are not false; they are **not empirically assessable**. They belong to a different domain: personal belief, fiction, or social identity.

3. Provisional vs. Structural

Non-Verifiability

A crucial distinction separates:

- **Provisional non-detection** – e.g., dark matter, gravitational waves (before 2015), the neutrino (before 1956). These entities are predicted to share at least one interaction channel (gravity, weak force) and are in principle detectable. **A future discovery could confirm or disconfirm them.** That is the key: we can specify what would count as evidence, even if we don't yet have it.
- **Structural, permanent non-verifiability** – Category A claims. The entity is defined so that **no possible future discovery** could ever count as confirmation or disconfirmation. Any proposed test is ruled out in advance. This is the hallmark of a fantasy attractor.

(This framework does not assert that dark matter could have been called a fantasy attractor before detection; dark matter always had specified interaction channels – gravity – and was therefore never structurally non-verifiable.)

4. Fantasy Attractor: Formal Definition

A belief system qualifies as a **fantasy attractor** if it meets the following conditions:

1. **No specified interaction channel** – The central claim lacks any measurable coupling to physical reality (Category A), or defines it in a way that systematically evades testing (Category B).
2. **Sealing mechanisms** – The belief incorporates rhetorical or cognitive strategies that neutralize disconfirming evidence (e.g., “God works in mysterious ways,” “The

ghost left when the EMF meter arrived”).

3. **Low corrective permeability ($\kappa \rightarrow 0$)** – The belief does not update in response to counterevidence; the return time τ to baseline is effectively infinite.
4. **Identity fusion** – The belief is tied to self-worth or group membership, making abandonment costly.

Under this definition, both Category A and some Category B claims can be fantasy attractors, but Category A are the paradigmatic case because they are structurally immune to evidence.

5. Fiction Is Real but Not True: A Crucial Distinction

The main argument might provoke an objection: *What about fiction? Sherlock Holmes is not physical, yet we say he exists as a character. Isn't that a counterexample to the claim that non-physical entities cannot be empirically distinguished from nonexistence?*

The objection fails because it conflates two different senses of “exists.” We must distinguish:

- **Fiction exists as physical information.** The character Sherlock Holmes is realized as patterns of ink on a page, as sounds in a performance, as neural firing patterns in readers' brains, or as bits on a computer screen. Information is a physical arrangement of matter. It shares interaction channels (energy, spacetime, causality) with the physical world. You can buy a book, discuss the plot, or be emotionally affected by a story. Fiction is **real** in this sense: it has a physical substrate and causal effects.

- **Fiction is not true.** The proposition “Sherlock Holmes lived at 221B Baker Street” does not correspond to any actual state of affairs in the world. It is false. Fiction is not required to be verifiable; it is understood as imagined.

Thus, the attractor framework happily accommodates fiction. It is real as information, but not claimed as true.

The bad faith of non-physical claims: Non-physical claims that demand to be treated as real – gods, ghosts, souls, hidden cabals – are *fiction pretending to be true*. They borrow the ontological status of real information (they exist as patterns in books, sermons, or brains) but also demand the epistemic authority of factual truth. Yet they refuse any possible test. They define themselves as beyond verification. This is bad faith: it is not metaphysics, but fiction that insists on being taken as fact while rejecting the rules of fact-checking.

Category	Exists as physical information?	Claims to be true?	Verifiable?	Framework classification
Fiction (Hamlet)	Yes	No (acknowledged as imagined)	Not applicable	Real information, not true
Scientific claim (neutrino)	Yes (theory, data)	Yes	In principle	Real, true (provisionally)
Non-physical claim (God)	Yes (as cultural artifact)	Yes	No – structurally excluded	Fantasy attractor

Therefore, the framework does not deny the reality of stories; it denies the epistemic legitimacy of treating unverifiable stories as facts. The fantasy attractor is not the story. It is the insistence that the story is true combined with the structural refusal to let the story be tested.

6. Vulnerability to Fraud and Manipulation

The structure of non-physical claims makes them **vulnerable** to fraud and manipulation – not that all such claims are fraudulent. Because there are no checks, a bad actor can assert divine commands, psychic readings, or secret knowledge without fear of disconfirmation. Sincere believers are not fraudsters, but the attractor basin can be exploited by those who understand its dynamics.

The framework diagnoses the **structure**, not the intent of every believer. It distinguishes **error, self-deception, motivated reasoning, and fraud** – all possible outcomes, but not all present in every case.

7. What This Argument Does Not Prove

To avoid overreach, the paper explicitly states what it does **not** claim:

- It does not prove that non-physical entities are logically impossible.
- It does not refute philosophical positions like Platonism (abstract objects) or classical theism that defines God as existence itself rather than an interacting object – though it notes that such positions are not empirically assessable.
- It does not claim that all believers are fraudsters or that all non-physical claims are meaningless in a philosophical sense.
- It does not assert a timeless criterion for what will be discovered in the future.

The claim is narrower: **within the attractor framework's physicalist commitment, non-physical claims are not empirically assessable, and they exhibit the dynamics of fantasy attractors.**

8. Conclusion

The attractor framework adopts a physicalist commitment: entities can only interact through shared interaction channels. Non-physical claims – defined as having no such channels – are not empirically assessable. They are fantasy attractors: belief systems structurally sealed against correction by permanent non-verifiability. This does not make them meaningless or false; it places them outside the domain of scientific ontology. Their structure makes them vulnerable to exploitation, but sincere belief is not fraud. The framework provides a diagnostic tool for recognising when a claim has been immunised against evidence, regardless of its content.

The argument supports the following conclusion:

Claims that are permanently insulated from any possible empirical correction occupy a distinct epistemic category and exhibit attractor dynamics that make them resistant to updating. Within the attractor framework's physicalist ontology, such claims cannot be empirically distinguished from nonexistence.

That is a substantial claim. It does not require asserting that non-physical realms cannot exist – only that they cannot be part of a scientific ontology, and that the beliefs which cling to them operate as fantasy attractors.

Suggested citation: Galida, R. S. (2026). Non-Physical Claims Are Fantasy Attractors: Why Unverifiable Realms Cannot Be Empirically Distinguished from Nonexistence. *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
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
κ^∞ (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)
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 °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 → transient absorption; chronic overload → 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 → transient absorption; persistent input → 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 → stable addiction attractor; low dose → 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 → ejection (buffered); large perturbation → 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 → ejection (invader fails); vacant niche → 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 → transient absorption; above threshold → 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 → ejection (rationalization); strong evidence → stable addition (belief change, rare).
- **Citation:** Nyhan & Reifler (2010).

Conspiracy and Extremist Movements

- State space: belief adoption × 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 → $\kappa \rightarrow 0$).
 - 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 → ejection (old task overwritten); replay/regularization → 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 → ejection (catastrophic forgetting); progressive networks or EWC → 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 $\kappa \rightarrow$ ejection (correction ignored); high $\kappa \rightarrow$ 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)

System / Domain	Operational τ ($\kappa = 1/\tau$)	τ Typical Timescale	Basin Depth B Proxy	Outcome	Notes
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)
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 ($> \kappa$) → ejection (backfire).
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.

Domain	Finding	Mapping to framework
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]
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 (already written)</i>	[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 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 $\Delta k(d)$ 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 $\Delta k(d)$ 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

Feature	Adaptive suppression (endpoint)	Pathological suppression (endpoint)
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 $\Delta k(d)$ 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 k 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- k 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.

(2026)

Robert Galida – June 2026 (Final)

Paper 2 in a series on conscious suppression; see [Paper 1: Intelligence Without Consciousness](#) for the full taxonomy of intelligence and consciousness.

Abstract

Why do people with addiction, trauma-related avoidance, or obsessive-compulsive disorder often know their behavior is maladaptive yet cannot stop? Standard explanations – impaired executive control, habit dominance, weak insight – are incomplete. This paper applies the attractor framework's suppression mechanism. In each disorder, the person *detects* the discrepancy between behavior and goals (insight is intact), but **phenomenal, identity-constitutive investment** – the felt urgency of craving, the necessity of avoidance, the compulsion to ritualize – deepens the attractor basin relative to corrective perturbations. The suppression is not a failure of intelligence; it is a dynamical competition between attractors. The paper distinguishes this account from dual-process and executive-control theories, provides falsifiable diagnostic criteria, and discusses treatment implications (why insight alone fails). Acknowledgment is made that for addiction, the relationship between incentive salience (*wanting*) and phenomenal consciousness remains contested; the model targets the subset of craving states that patients report as felt urgency.

1. Introduction: The Paradox of Insight Without Change

A person with alcohol use disorder knows that drinking damages their health, relationships, and future. Yet when a craving arises, they drink. A trauma survivor knows that the parking garage is safe, yet they avoid it. A person with OCD knows that the ritual is irrational, yet they perform it.

Standard explanations invoke **impaired executive control** (Volkow et al., 2016), **habit dominance** (Balleine & Dickinson, 1998), or **lack of insight** (Amador et al., 1994). But these accounts do not explain why the person can articulate the harm, describe counterarguments, and intend change, yet the behavior persists. Executive control may be intact in non-trigger contexts; habits may be sensitive to goal-level knowledge; insight may be partial or oscillating.

The attractor framework provides a model of **motivational competition** where a conscious, identity-binding urge temporarily overrides the correction signal. In *Intelligence Without Consciousness* (Galida, 2026), we introduced **conscious suppression**: phenomenal, identity-constitutive commitment deepens an attractor basin, making it resistant to corrective perturbations. This paper applies that mechanism to addiction, trauma-related avoidance (PTSD), and OCD. It does not deny executive or habit deficits; it proposes that in many cases, a conscious-level attractor competition is the primary obstacle to change.

2. Defining Conscious Suppression (Self-Contained Glossary)

For readers unfamiliar with Paper 1:

- **Attractor basin** – the set of states from which a system returns to a stable pattern. A deeper basin resists larger perturbations.
- **Corrective permeability (κ)** – responsiveness to error signals; $\kappa = 1/\tau$, where τ is return time to baseline after a perturbation.
- **Conscious suppression** – a process where the person *experiences* an urge, fear, or compulsion as felt, identity-relevant, and *not chosen* (non-deliberative), yet the depth of that attractor prevents escape from the maladaptive behavior. This corresponds to **Level 3** in Paper 1: detection of error + suppression via basin depth. Level 2 (automatic bias without error detection) and Level 1 (unfamiliarity) are not the target.

On sealing mechanisms: The paper treats sealing mechanisms (e.g., rationalizations) as *attractor-consistent outputs* generated by the basin state, not as deliberate strategic choices. Although they may *feel* deliberate to the patient, the model treats them as expressions of the attractor's depth, not as independent volitional acts. This resolves the tension between “non-deliberative urgency” and the deployment of rationalizations.

3. Empirical Grounding

Addiction:

Volkow et al. (2016) demonstrate that chronic substance use impairs prefrontal executive function in a state-dependent manner – deficits emerge under craving or stress, not at all times. Individuals can maintain intact verbal knowledge of consequences and express intention to stop (Goldstein et al., 2009). The craving state has been modeled as a competing attractor (Redish, 2004; Gutkin et al.,

2006). **Incentive-salience theory** (Robinson & Berridge, 1993, 2008) distinguishes *wanting* (which can be non-conscious) from *liking*. The present model targets the subset of craving states that are *phenomenally accessible* – the patient’s reported felt urgency. This is a narrower claim; the paper does not assume that all incentive-salience processes are conscious.

PTSD & avoidance:

Extinction recall deficits (Milad et al., 2006) are well documented, but they do not fully account for conscious fear as *necessary* even when safety is known. Meta-analyses confirm vmPFC–amygdala decoupling in PTSD (e.g., Etkin & Wager, 2007, and subsequent reviews). Ecological momentary assessment (EMA) studies in representative samples show that individuals with PTSD often report high certainty of safety before trigger environments yet avoidance persists (see, e.g., reviews of EMA in PTSD). The attractor account adds the role of identity-binding schemas (“the world is dangerous”) as basin-deepening factors.

OCD:

The DSM-5-TR includes an insight specifier: *good/fair, poor, or absent*. Approximately 25–30% of individuals with OCD have poor insight (Catapano et al., 2010). This paper targets the **good-insight subgroup** (where the person recognizes irrationality). For poor-insight patients, the mechanism may be closer to Level 2 (automatic compulsion without error detection).

Recent literature (2015–2025):

- EMA studies of craving show that momentary urge strength predicts relapse better than global insight (Serre et al., 2015; Shiffman et al., 2020).
- OCD outcome studies confirm that poor insight predicts worse response to ERP (García-Soriano et al., 2021).

Good-insight patients still show substantial residual symptoms, consistent with a competition model.

- Identity-shifting interventions for addiction (Best et al., 2016) support the importance of decoupling selfhood from “addict” identity.
-

4. Three Clinical Patterns

4.1 Addiction

- **Mechanism:** Craving as a state-dependent attractor that overrides goal-directed control when triggered. Identity fusion (“I am an addict”) deepens the basin where present, but is not universal.
- **Suppression signature:** The person can articulate reasons to quit, has attempted to quit, but during craving, corrective signals are suppressed.
- **Sealing mechanisms:** Cognitive rationalizations (“just this once,” “I need it to cope”) that block the error signal from updating the basin – treated as attractor-consistent outputs, not deliberate choices.

4.2 Trauma-Related Avoidance (PTSD)

- **Mechanism:** Conditioned fear creates an avoidance attractor. Safety knowledge may be intact, but felt necessity dominates.
- **Suppression signature:** “I know it’s safe, but I can’t go in.”
- **Identity fusion:** “The world is dangerous” as a self-defining schema.

4.3 Obsessive-Compulsive Disorder (OCD – Good Insight Subgroup)

- **Mechanism:** Anxiety drives compulsions that temporarily reduce distress, despite knowledge of irrationality.
- **Suppression signature:** “I know it doesn’t make sense, but I have to do it.”
- **Sealing mechanisms:** “Better safe than sorry,” “It’s a small price to pay for certainty.”

5. Transdiagnostic Table

Disorder	Error signal detected	Conscious investment	What maintains basin depth (mechanism)
Addiction	Knowledge of negative consequences	Craving (felt urgency)	Reinforcement schedule + state-dependent executive impairment + (sometimes) identity fusion
Trauma avoidance	Safety knowledge (cognitive)	Fear (felt necessity)	Extinction resistance + hyperarousal + schema of danger
OCD (good insight)	Knowledge of irrationality	Anxiety (felt urgency)	Negative reinforcement via distress reduction + certainty-seeking belief

6. Diagnostic Criteria for Clinical Fantasy Attractors (Operationalized)

A patient's presentation is a **candidate** clinical fantasy attractor if it meets **three of five** criteria (provisional threshold; empirical validation required). The Level 2/3 distinction requires momentary assessment (see §7).

1. **Insight intact:** The patient can state, unprompted, the discrepancy between behavior and goals. *Operationalization:* Score ≥ 4 on the Brown Assessment of Beliefs Scale (BABS) insight item, or equivalent.
2. **Conscious urgency:** The maladaptive behavior is preceded by a felt, urgent state (craving, fear, anxiety) rated by the patient as "overwhelming" or "necessary." *Operationalization:* Momentary ecological assessment (EMA) rating $> 7/10$ before the behavior.
3. **Identity fusion:** The patient endorses that the behavior or its avoidance is central to selfhood (e.g., "I am an addict," "I must do this to be safe"). *Operationalization:* Endorsement of at least one identity statement on a structured interview.
4. **Low corrective permeability in trigger contexts:** Repeated corrective information (psychoeducation, feedback) does not reduce the behavior. *Operationalization:* No significant reduction after three sessions of evidence-based psychoeducation alone.
5. **Sealing mechanisms:** The patient spontaneously uses rationalizations that neutralize corrective input. *Operationalization:* Qualitative coding of patient speech (inter-rater reliability to be established; currently a research gap).

Counter-criteria (exclude if any present):

- The patient cannot state the discrepancy (insight absent) – then Level 2 or 1.
 - The behavior stops entirely after receiving corrective information alone – then basin depth was shallow.
-

7. The Detection Problem (Level 2 vs. 3) in Clinical Practice

Distinguishing automatic compulsion without error detection (Level 2) from conscious suppression with error detection (Level 3) requires:

- **Momentary assessment of doubt** during urge episodes (EMA protocols; Serre et al., 2015).
- **Reaction time paradigms** (e.g., Gillan et al., 2014, for goal-directed vs. habitual control in OCD; note that the specific link to error detection latency remains an active area).
- **Physiological markers** (dissociation between cognitive knowledge and fear response suggests Level 3).

These methods are promising but not fully validated; the paper specifies directions for needed research.

8. Implications for Treatment

Insight-only interventions (psychoeducation, cognitive restructuring alone) often fail in these disorders because the basin depth is maintained by conscious urgency, not lack of knowledge.

Effective treatment must **reduce basin depth** or **increase**

corrective force:

- **Addiction:** Pharmacological reduction of craving (e.g., naltrexone; emerging evidence for GLP-1 agonists – see recent reviews, e.g., Klausen et al., 2022, for GLP-1 receptors and alcohol, and emerging clinical reports), contingency management, and identity-shifting interventions (Best et al., 2016).
- **Trauma:** Exposure therapy (increasing corrective force) combined with arousal reduction. The mechanism is basin reshaping, not insight.
- **OCD:** Exposure and response prevention (ERP) directly targets the basin by preventing the compulsion while the patient experiences urgency. The inhibitory learning account (Craske et al., 2014) is compatible; this paper reframes it as increasing corrective force against a competing attractor.

The prediction: treatments that solely enhance insight will be less effective for patients meeting the diagnostic criteria than treatments that directly target basin depth or corrective force.

9. Open Questions

- **Measuring basin depth in clinical settings:** Subjective urgency scales, behavioral persistence tasks, heart rate variability. A Clinician Basin Depth Scale (CBDS) is a research priority.
- **Level 2 vs. 3 differentiation:** Can EMA and reaction time methods reliably classify patients? Pilot studies needed.
- **Diagnostic threshold validation:** The “three of five” criterion requires empirical ROC analysis against

treatment response.

- **Disorders where suppression is purely Level 2:** Some impulse control disorders or psychotic conditions may not meet the conscious detection criterion.
-

10. Conclusion

Addiction, trauma-related avoidance, and OCD (good insight subtype) are not failures of intelligence. They are cases where conscious, identity-constitutive investment deepens an attractor basin relative to corrective perturbations. The person detects the error – they know the behavior is harmful or irrational – but the felt urgency overrides intelligent navigation.

This diagnosis explains why insight alone fails and why treatments that target basin depth succeed. The clinical fantasy attractor is a trapped navigator: intelligent, aware, but unable to escape.

The dance of recovery is not about knowing the way out. It is about reshaping the attractor landscape so that the path to safety becomes shallower than the pull to stay.

Suggested citation: Galida, R. S. (2026). Trapped Navigation: Addiction, Trauma, and OCD as Conscious Suppression of Intelligent Correction. *Fantasy Attractor*.