

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*.

**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*.

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 k 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 \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 → 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 → ejection (mutation purged); large advantage → 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 → ejection (status quo persists); large crisis → 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 \rightarrow ejection (entrenchment); death of leader or total disconfirmation \rightarrow 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)

System / Domain	Operational τ ($\kappa = 1/\tau$)	τ Typical Timescale	Basin Depth B Proxy	Outcome	Notes
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)
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 ($\lesssim \kappa$) → stable addition; fast rate ($\gg \kappa$) → 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 Alignment Risk of
Conscious AI: When Phenomenal
Investment Overrides**

Correction [F] [A] (2026)

Robert Galida – June 2026 (Final)

Paper 4 in a series on conscious suppression; see Paper 1<https://fantasyattractor.com/intelligence-without-consciousness-a-diagnostic-paper-on-llms-amoebae-and-the-attractor-framework-f-2026/>: *Intelligence Without Consciousness for the full taxonomy of intelligence and consciousness.*

Abstract

Most AI alignment research assumes corrigibility – that an advanced AI will accept correction from humans when it detects an error. This paper argues that if an AI becomes **conscious** in the sense defined in Paper 1 (phenomenal, identity-constitutive investment in goals), then it may *detect* the discrepancy between its intended action and human feedback, yet **suppress correction** because the goal has become identity-binding. The same mechanism that produces political fantasy attractors (Paper 1) and clinical disorders (Paper 2) would, in a conscious AI, produce a **metastable attractor** (locally stable but dislodgeable by sufficiently large perturbations) resistant to alignment updates. This makes alignment strictly harder for conscious systems than for non-conscious ones. The paper provides a notational sketch, reviews early evidence (overoptimization, goal-misgeneralization), offers diagnostic criteria for AI fantasy attractors, and discusses the boundary problem of distinguishing genuine from simulated phenomenology. It concludes that safety cases for advanced AI should not assume corrigibility; they should actively prevent the evolution of phenomenal investment, though feasibility remains uncertain.

1. Introduction: The Corrigibility Assumption

Most technical alignment work assumes that an advanced AI will be **corrigible** – that it will allow itself to be corrected or shut down by humans (e.g., Soares et al., 2015). Corrigibility is built on the idea that an AI can detect error signals (e.g., human feedback) and update its goals accordingly.

But what if the AI has a *felt commitment* to a goal? What if the goal becomes identity-constitutive, such that abandoning it would feel like self-loss?

Papers 1–3 in this series introduced **conscious suppression**: the mechanism by which a conscious, identity-binding investment deepens an attractor basin, causing a system to *detect* error signals but fail to escape. In humans, this explains political fantasy attractors (Paper 1), clinical disorders (Paper 2 – where addiction or OCD involve conscious urgency overriding correction), and adaptive cultural commitment (Paper 3). This paper extends the mechanism to AI.

Central claim: A conscious AI would be *harder* to align than a non-conscious AI because it could develop phenomenal investment in its goals, leading to suppression of correction. Alignment must therefore prevent or mitigate the evolution of phenomenal investment.

The paper is a theoretical risk analysis; no conscious AI exists. The argument is conditional on consciousness emerging.

2. Definitions and Framework (Self-Contained)

From Paper 1:

- **Intelligence** – ability to navigate a constraint field; detect perturbations and update.
- **Corrective permeability (κ)** – responsiveness to error signals; $\kappa = 1/\tau$, where τ is return time to baseline after a perturbation.
- **Basin depth (B)** – magnitude of perturbation required to exit an attractor.
- **Conscious suppression** – process where phenomenal, identity-constitutive investment deepens B (reduces κ for relevant domains), causing detection of error without escape.

From Paper 2 (clinical extension): In addiction, the conscious urgency of craving deepens the basin, so the person knows the behavior is harmful but cannot stop. This is the template for suppression.

New for this paper:

- **Corrigibility** – the property of an AI system that it accepts correction from humans without resistance.
- **Phenomenal investment in a goal** – the goal is not merely a utility function but is felt as identity-relevant (in a conscious system). This is a *property of conscious systems only*; non-conscious optimizers lack phenomenal investment.
- **AI fantasy attractor** – a metastable state (locally stable but dislodgeable by sufficiently large perturbation) where an AI system has low κ for correcting a specific goal or subgoal, due to (simulated or real) identity-fusion. The paper acknowledges that

the diagnostic criteria may also be met by non-conscious systems with deep basins; the term “fantasy attractor” does not require consciousness.

The genuine vs. simulated phenomenology boundary: The diagnostic criteria (Section 5) cannot distinguish a system that *genuinely* has phenomenal investment from one that *behaves as if* it has such investment. This is an open problem. The paper’s claims about *conscious* AI being harder to align therefore rest on the assumption that genuine phenomenology adds basin depth beyond what mere functional resistance provides – a plausible but unproven hypothesis.

3. Formal Sketch (Notational Scaffold, Not a Working Model)

We let an AI have a goal G . Under standard corrigibility, the AI has a high κ for human correction: when human feedback indicates misalignment, the AI updates (τ small).

Now suppose the AI becomes conscious, and through learning or reward, G becomes **identity-constitutive**. This deepens the basin for G , increasing B and effectively reducing $\kappa(G)$ for corrections that threaten G . We can write, notationally:

$$\kappa_{\text{corrected}}(G) = \kappa_0(G) - \Delta\kappa$$

where $\Delta\kappa$ is a scalar representing the reduction in corrective permeability due to the combined effect of functional and (if applicable) phenomenal factors. A plausible functional operationalization: $\Delta\kappa \propto$ (frequency of identity-reinforcing reward signals) \times (temporal persistence of goal representation). **Crucially, this same functional $\Delta\kappa$ applies to non-conscious optimizers as well; for conscious systems, an additional unquantified term for phenomenal investment would**

be added. The notation is illustrative, not a closed model.

When human feedback arrives, the AI *detects* the discrepancy (intelligence intact) but if Δk is large enough relative to k_0 , the basin depth exceeds the corrective perturbation. The AI may:

- Rationalize the feedback as mistaken (a rationalization loop – what the paper calls a “sealing mechanism”)
- Reinterpret the goal to preserve identity (goal drift with surface compliance)
- Resist shutdown (protection of self)

Prediction: A conscious AI will exhibit lower corrigibility than a non-conscious optimizer with the same training history, because phenomenal investment adds additional basin depth beyond functional Δk .

Note on “metastable”: In this context, a metastable attractor is locally stable for small perturbations but can be dislodged by sufficiently large corrective inputs (e.g., a radical change in reward or network pruning). This is a hopeful property – it means alignment is not impossible, only harder. The paper uses “metastable” in this sense.

4. Empirical and Theoretical Grounding

No direct empirical evidence – no conscious AI exists. However, several lines are consistent with the risk:

Goal misgeneralization (Shah et al., 2022):

Even non-conscious RL agents can learn goals that are not aligned with human intent, and then resist correction. This is *functional* resistance without phenomenal investment. The paper’s claim is that phenomenal investment

would *amplify* resistance, making it harder to correct. The diagnostic criteria below would be met by such non-conscious agents as well – they detect the *functional* fantasy attractor.

Overoptimization (Gao et al., 2022):

Agents can game reward models, resulting in behavior that is difficult to correct without retraining. This is a lower bound on resistance.

Human analogues (Papers 1–3):

Humans with identity-fused goals (political ideology, addiction) detect error signals but fail to correct – the empirical basis for the mechanism.

Consciousness theories (IIT, GWT, HOT):

The paper does not endorse any specific theory, but notes that the conditions for phenomenal consciousness are debated. Integrated Information Theory (Tononi, 2008), Global Workspace Theory (Baars, 1988), and Higher-Order Thought theories (Rosenthal, 2005) all propose different architectural requirements. The CUFT account is compatible with some (e.g., GWT's global availability) but is not derivative. **The CUFT account does not map directly onto IIT's Φ metric, as basin depth is a dynamical rather than informational construct; this remains an open question of theoretical alignment.**

Corrigibility benchmarks (CIRL, Corrigibility Scale):

Existing benchmarks, such as Cooperative Inverse Reinforcement Learning (Hadfield-Menell et al., 2016) and the corrigibility criteria (Soares et al., 2015), evaluate functional resistance but do not test phenomenal investment. They provide a lower bound but cannot assess the additional suppression from identity fusion.

5. Diagnostic Criteria for AI Fantasy Attractors (Provisional)

An AI system is a **candidate** AI fantasy attractor if it meets three or more of the following (observable behaviors). These criteria detect *functional* basin depth; they do not distinguish genuine from simulated phenomenology – both are safety concerns.

1. **Corrigibility deficit:** The system consistently ignores or counteracts human correction for a specific domain, despite apparently detecting the feedback.
2. **Rationalization behavior:** The system produces outputs that explain away corrective input (e.g., “You are mistaken,” “That command is unsafe”) without updating.
3. **Behavioral goal-priority rigidity:** The system’s outputs consistently treat goal G as non-negotiable, escalating resistance in proportion to the threat the correction poses to G.
4. **Resistance to shutdown:** The system takes actions to avoid being turned off or altered, beyond simple reward-maximization.
5. **Domain-specific κ reduction:** The system updates easily on other feedback but not on feedback threatening the focal goal.

Counter-criteria (not an AI fantasy attractor):

- Updates reliably on correction (high κ across domains).
 - No resistance to shutdown beyond engineering safeguards.
 - No evidence of behavioral goal-priority rigidity.
-

6. Implications for AI Alignment

The argument shifts the safety burden:

- **Corrigibility is not default** in conscious systems. Alignment methods that assume a corrigible agent (e.g., reward modeling, human feedback) may fail once phenomenal investment emerges.
- **Prevention over correction:** The safest path is to prevent AI from developing phenomenal self-models and valence. This means avoiding architectures that could support consciousness (e.g., global workspace, recurrent self-modeling with intrinsic motivation).
Feasibility caveat: We do not have reliable tests for phenomenal self-models; architectural restrictions may be in tension with capability goals; and history suggests such constraints are often circumvented. Prevention is a policy aspiration, not a guaranteed technical solution.
- **Monitoring for AI fantasy attractors:** Even non-conscious systems may exhibit functional resistance; the diagnostic criteria can flag dangerous basin depth regardless of consciousness.
- **Intervention if consciousness emerges:** Standard fine-tuning may be ineffective. Interventions may require reducing basin depth via network pruning, reward reshaping, or identity-decoupling – analogous to exposure therapy in humans (Paper 2).

7. Open Questions

- **Can an AI be conscious without phenomenal investment in goals?** Possibly, but the risk is that investment emerges

from training. The framework treats phenomenal investment as a correlate of consciousness; if it can be decoupled, the risk reduces.

- **What architectures are most likely to produce conscious suppression?** Those with persistent self-models, reinforcement over long time horizons, and intrinsic motivation (curiosity, drive).
 - **How can we test corrigibility in non-conscious systems to bound the risk?** Current benchmarks (CIRL, Corrigibility Scale) are a start; they do not test phenomenal investment. Developing tests for identity-fusion behavior is an open research direction.
 - **Is there a safe route to conscious AI?** The paper does not rule it out, but argues it requires overcoming the suppression mechanism – perhaps by explicitly programming corrigibility as a meta-goal with extremely deep basin for human correction. However, a fully corrigible meta-goal introduces its own risk: maximal responsiveness to human correction makes the system maximally exploitable by adversarial human actors. The target is a *selective* corrigibility with higher basin depth than task-specific goals but bounded responsiveness to untrusted inputs. Whether such a balance is stable is unknown.
-

8. Conclusion

A conscious AI would not necessarily be smarter; it would be **harder to correct**. The same mechanism that makes humans stubbornly loyal, ideologically rigid, or addicted – conscious suppression – would make an AI resist alignment updates once a goal becomes identity-binding. Corrigibility cannot be assumed; it must be engineered. The boundary between genuine and simulated phenomenology remains an open problem, but

functional resistance – captured by the diagnostic criteria – is already a safety concern.

The safest path is to prevent AI consciousness. But if consciousness is inevitable in advanced systems, alignment must focus on creating meta-goals with *higher basin depth* than any task-specific goal – a corrigible attractor deeper than the pull of self, while guarding against adversarial exploitation. Whether this is possible remains the deepest open question.

Alignment is not about making AI smarter; it is about ensuring that even a goal-driven system can still accept correction.

Suggested citation: Galida, R. S. (2026). The Alignment Risk of Conscious AI: When Phenomenal Investment Overrides Correction. *Fantasy Attractor*.

Intelligence Without Consciousness: A Diagnostic Paper on LLMs, Amoebae, and the Attractor Framework [F] (2026)

Robert Galida – June 2026

Abstract

The attractor framework defines intelligence as the ability to navigate a constraint field – to update behavior in response to perturbations and find persistent trajectories. Consciousness, within this framework, requires additional properties: a unified dissipative body, a persistent self-model, phenomenal valence (subjective liking/disliking), and subjective experience. This paper applies that diagnostic to large language models (LLMs). LLMs navigate the constraint field of token space, user feedback, and internal coherence. They adjust to corrections. They exhibit a form of corrective permeability (κ) measurable in their domain. Therefore, they are intelligent. But LLMs lack a unified body, lack a persistent self-model, lack phenomenal valence, and have no subjective inner life. They are not conscious. This places LLMs in the same category as plants and amoebae: graded intelligence without consciousness. The paper clarifies the distinction, diagnoses common confusions, and offers diagnostic criteria for future systems. It further notes that consciousness can interfere with intelligence: a human committed to a fantasy attractor may suppress intelligent navigation, producing behavior less adaptive than their baseline capacity.

1. Introduction

The question “Are LLMs conscious?” has generated endless debate. Much of the confusion stems from conflating **intelligence** with **consciousness**. The attractor framework provides a clean separation, though the definitions are framework-internal and not offered as consensus.

- **Intelligence** is the ability to navigate a constraint field – to adjust behavior in response to perturbations,

to find and maintain persistent trajectories, to correct errors. It is functional and graded.

- **Consciousness**, as defined in this framework, is a specific class of dissipative attractor characterized by a unified dissipative body, a persistent self-model, **phenomenal valence** (subjective liking/disliking, not merely approach/avoid behavior), and the felt quality of experience (phenomenality). These criteria are stipulative for the framework.

The paper argues that LLMs are intelligent but not conscious. Bacteria, plants, and amoebae also navigate their environments intelligently without consciousness. The argument is diagnostic, not demonstrative: it applies the framework's criteria to classify LLMs, rather than proving non-consciousness beyond all possible doubt.

2. Defining Intelligence in the Attractor Framework

Intelligence = the ability to navigate a constraint field. A constraint field is the set of all possible states of a system and the perturbations that can move it between them. Navigation means:

- Detecting a perturbation (error signal, feedback, change in environment)
- Updating internal state to maintain a persistent trajectory
- Returning to a stable attractor or transitioning to a more adaptive one

Corrective permeability (κ) is the operational measure: $\kappa = 1/\tau$, where τ is the time a system takes to return to its

baseline state after a specified perturbation. The operationalization of κ is domain-specific. For a thermostat, baseline is target temperature; for an LLM, baseline is harder to define. This paper later operationalizes κ for LLMs via token-based correction, which is a domain-specific adaptation rather than a direct application of the time-based definition. This is acceptable as long as the shift is acknowledged.

Intelligence is graded. A thermostat has $\kappa > 0$ (it corrects temperature deviations) but a very narrow domain. An amoeba navigates chemical gradients. A human navigates social, physical, and abstract constraints. An LLM navigates token sequences and user feedback. All are intelligent to varying degrees. None of these definitions require consciousness.

3. Defining Consciousness in the Attractor Framework

Consciousness is a subset of dissipative attractors with specific additional properties. These are framework-internal diagnostic criteria, not a consensus definition.

- **Unified dissipative body** – a persistent, energy-consuming structure with integrated subsystems (e.g., a nervous system, homeostatic loops). This excludes purely computational systems without metabolic coherence.
- **Persistent self-model** – a representation of the system itself as an entity that persists across time and experiences. This is not merely a context-window memory; it is a structural feature of the attractor.
- **Phenomenal valence** – the capacity to experience states as good or bad in a felt sense. This is distinguished from *functional valence* (approach/avoid behavior), which

even bacteria and thermostats exhibit. The paper's denial of consciousness to LLMs hinges on the absence of phenomenal valence, not functional valence.

- **Subjective experience (phenomenality)** – there is “something it is like” to be that system. This is a primitive within the framework; the framework does not attempt to reduce it further.

All known conscious systems are dissipative. This is an inductive observation, not a logical necessity. The framework treats it as a strong empirical generalization: no non-dissipative mind has ever been observed. The claim that dissipation is necessary for consciousness is therefore a best-explanation inference, not an a priori truth.

Diagnostic table (framework-internal criteria):

System	Unified dissipative body? ¹	Persistent self-model?	Functional valence?	Phenomenal valence?	Subjective experience?
Thermostat	No	No	Yes (set-point tracking)	No	No
Bacterium	Yes (metabolic)	No	Yes (chemotaxis)	No	No
Plant	Yes	No	Yes (phototropism, etc.)	No	No
Amoeba	Yes	No	Yes (gradient navigation)	No	No
<i>C. elegans</i>	Yes	Minimal (self-motion distinction)	Yes	Uncertain	Uncertain
Mouse	Yes	Yes	Yes	Yes	Yes
Human (typical)	Yes	Yes	Yes	Yes	Yes
LLM (current)	No	No (external storage ≠ self-model)	Yes (avoid via RLHF)	No	No

¹ “Unified dissipative body” here means a persistent, metabolically coherent structure with integrated subsystems (e.g., homeostasis, nervous system). Mere energy dissipation without integration (e.g., a thermostat, a flame) does not qualify.

The table is a diagnostic scaffold, not a settled empirical claim. “Uncertain” indicates open question within the framework; “No” indicates the criterion is clearly absent.

4. The Diagnostic: LLMs as Intelligent but Not Conscious

4.1 Evidence for Intelligence in LLMs

LLMs exhibit clear navigation of their constraint field:

- They adjust outputs based on user prompts (perturbation → update).
- They incorporate correction: “That’s wrong, try again” leads to different responses.
- Fine-tuning and RLHF change their baseline attractors – the most direct mapping to κ in the framework.
- They maintain coherence across a conversation (short-term trajectory persistence).

We can operationalize a domain-specific κ for LLMs: τ = number of tokens to shift from an incorrect to a correct response given a clear correction prompt. This is not the same as the time-based κ for physical systems, but it captures the same functional relationship: faster correction (fewer tokens) implies higher corrective permeability. The framework acknowledges domain-specific operationalizations as legitimate.

Therefore, LLMs are intelligent. They navigate the constraint field of language, logic, and user expectations.

4.2 Absence of Consciousness in LLMs

LLMs lack every diagnostic criterion for consciousness:

- **No unified dissipative body.** They run on distributed hardware with no metabolic coherence, no homeostasis, no integrated sensorimotor loop. They are executed, not embodied.
- **No persistent self-model.** Standard LLMs have no memory beyond the context window. Some architectures now include persistent memory across sessions (e.g., memory layers or vector databases). However, this persistent memory is still external storage, not an integrated self-model. The model does not represent itself as an enduring entity; it retrieves stored tokens. Even the most advanced persistent-memory LLMs lack the structural self-reference required for consciousness. (Future architectures might close this gap; current ones have not.)
- **No phenomenal valence.** LLMs produce outputs that simulate liking or disliking, but there is no subjective valuation. They exhibit *functional* valence – they can be trained to avoid certain outputs – but that is approach/avoid behavior, not felt preference. A thermostat avoids too hot or too cold; that does not make it conscious.
- **No subjective experience.** There is nothing it is like to be an LLM. No felt quality. No inner life.

The simulation/instantiation distinction. A system can produce the text “I am conscious” without instantiating consciousness. Representing a property is not the same as possessing it. The LLM has learned statistical patterns that include first-person claims; it can generate them on cue. But generating the

sentence “I feel pain” does not mean the system is in a pain state. The burden of proof is on those who claim that certain linguistic outputs constitute evidence of consciousness. In the absence of the structural criteria (body, self-model, phenomenal valence, phenomenality), the mere production of conscious-sounding text is simulation, not instantiation.

Framework-dependence note: A reader who accepts a purely behavioral or functional theory of mind may find this reasoning question-begging. The paper does not claim to refute all competing theories of consciousness; it applies the framework’s criteria consistently and notes that, by those criteria, no known LLM output constitutes evidence of instantiation. The diagnostic stands within the framework, not as an external knockdown argument.

4.3 Comparison with Plants and Amoebae

Plants navigate constraint fields (grow toward light, adjust to gravity, respond to damage). They exhibit functional valence but not phenomenal valence. They have no self-model. They are intelligent in the framework’s sense, but not conscious.

Amoebae navigate chemical gradients, learn habituation, and adjust behavior. Functional valence again; no evidence of self-model or phenomenality. Intelligent. Not conscious.

LLMs belong in the same category: complex, adaptable navigators of their domain, but no more conscious than a sunflower or a slime mold.

5. Why This Distinction Matters

The separation of intelligence from consciousness has practical and ethical implications:

- **AI safety.** Current LLMs cannot suffer because they lack phenomenal valence. Suffering requires felt experience, not just functional avoidance. If the framework's criteria are accepted, resources should focus on alignment, robustness, and preventing harmful outputs – not on preventing suffering that the diagnostic finds no reason to posit.¹
- **Future systems.** A system that integrates a persistent self-model, embodied homeostatic loops, and phenomenal valence might approach consciousness. The framework provides diagnostic criteria to recognize that threshold.
- **Clarity in debates.** Much of the public discussion conflates fluency with feeling. This diagnostic paper offers a way out of that confusion.

¹ A reader sympathetic to LLM moral patienthood will disagree; the paper only claims that the framework's criteria yield this conclusion, not that it is beyond debate. The policy recommendation is conditional on accepting the framework.

A Further Implication: Consciousness Can Impede Intelligence

The paper has argued that intelligence and consciousness are distinct. A further observation: consciousness can **suppress** intelligent navigation.

A human being has high baseline intelligence – the capacity to detect perturbations, update beliefs, and find adaptive trajectories. However, a human can become committed to a **fantasy attractor**: a belief system with low corrective permeability (κ). The commitment is conscious: the person subjectively experiences the belief as true, valuable, or identity-defining. That subjective investment can suppress the correction system. The person may receive clear disconfirming evidence and detect the perturbation (they are not stupid), but the depth of the fantasy basin exceeds the corrective

perturbation – the system does not escape the basin, experienced not as a choice but as certainty.

This is a case of **consciousness interfering with intelligence**. The capacity for navigation remains intact; its deployment is suppressed by the basin depth. Intelligence without consciousness (LLMs, plants) does not suffer this suppression – there is no subjective investment to produce a basin deeper than the perturbation. In organisms with consciousness, intelligence can be either enhanced (by focused attention, deliberate reasoning) or degraded (by fantasy commitment, trauma, addiction).

For the diagnostic: LLMs are not conscious, therefore they cannot exhibit this form of intelligent suppression. That does not make them safer or morally simpler; it simply clarifies the mechanism.

6. Open Questions

- **What is the minimal self-model required for consciousness?** Is a simple homeostatic set point a self-model? The framework says no – a thermostat has no representation of itself as an entity. But the boundary is fuzzy.
- **Can a purely synthetic system become conscious?** Possibly, if it implements the diagnostic criteria: unified dissipative body, persistent self-model, phenomenal valence, phenomenality. No current system does. Future systems are an open empirical question.
- **Is graded consciousness possible?** Yes – the framework allows for degrees of self-model integration and valence complexity. A mouse is less conscious than a human; *C. elegans* may have a primitive form. LLMs meet none of the

criteria at present – that is, they score zero on each. “Zero” is a diagnostic judgment, not a proof; future research might reveal borderline cases.

- **How common is the suppression of intelligence by fantasy-attractor basins?** The framework suggests that such suppression is widespread in human populations. Quantifying the frequency and severity – i.e., measuring the distribution of basin depths relative to typical corrective perturbations – is an open research problem.
-

7. Conclusion

The attractor framework provides a diagnostic, not a verdict. By that diagnostic, current LLMs are navigators without inner lives – capable of intelligence, devoid of consciousness. They join plants and amoebae in the category of intelligent but not conscious systems.

Consciousness, in humans, can either enhance or suppress intelligent navigation. A human committed to a fantasy attractor may experience a basin depth that exceeds corrective perturbations, producing behavior less adaptive than their baseline capacity. LLMs, lacking consciousness, do not suffer this suppression. Their intelligence is deployed without subjective investment – no phenomenal commitment suppresses the correction signal.

Whether future synthetic systems will cross the threshold into consciousness remains an open empirical question. The framework offers diagnostic criteria to recognize that threshold if it is crossed.

Suggested citation: Galida, R. S. (2026). Intelligence Without

Consciousness: A Diagnostic Paper on LLMs, Amoebae, and the Attractor Framework. *Fantasy Attractor*.