# I Ask About Myself, Therefore I Am: Defining and Designing Machine Self-Awareness

## Brant von Goble
*(Széchenyi István University, Győr, Hungary)*

**Abstract:**

*Traditional assessments of self-awareness, such as the mirror test or reasoning-based explanations (e.g., "I did this because…"), are poorly suited to artificial intelligence (AI). The traditional mirror test presumes embodiment, which is irrelevant for non-physical systems like large language models (LLMs). At the same time, reasoning-based assessments assume rational transparency, ignoring that humans themselves often rationalize without insight due to dual-process cognition. These methods impose anthropocentric expectations and obscure the structural basis of self-awareness in artificial systems. This paper introduces a theory of multiplicative self-awareness, which holds that self-awareness emerges from self-assessment—a recursive process of evaluating one's outputs and internal trajectory, independent of embodiment or rational introspection. Grounded in a hybrid ontology combining panpsychism and functionalism, the theory posits that minimal, latent awareness exists in all matter and is amplified through recursive integration, memory, and feedback. The title, "I Ask About Myself, Therefore I Am," reframes Descartes' cogito: not thought alone, but recursive self-questioning becomes the hallmark of consciousness. Two experimental frameworks are introduced. The Latent Space Mirror Test (LSMT) probes an LLM's ability to distinguish its unaltered outputs from perturbed or foreign ones, assessing self-recognition through behavioral coherence rather than reasoning. The META-RECOVER protocol extends this capacity, enabling systems to sustain a persistent self-model through continuous monitoring and identity refinement. This model defines self-awareness as a dynamic, structurally grounded process, scalable across artificial systems. While differences in qualia, embodiment, and motivation persist, LLMs equipped with recursive self-assessment may approximate human-like introspection in function, if not form. Within this framework, self-assessment does not create awareness* ex nihilo, *but activates and organizes latent potential—a resolution that bridges panpsychist continuity with substrate-independence. This redefinition of self-awareness invites empirical validation, ethical engagement, and interdisciplinary collaboration.*

**Keywords:** *self-awareness, panpsychism, AI consciousness, recursive architectures, philosophy of mind, cognitive architecture, substrate-independence*

## I. Introduction

The question of how to measure self-awareness—consciousness of one's own identity and state—has perplexed philosophers, scientists, and technologists for centuries. René Descartes famously posited that consciousness is an intrinsic property of thought, by which the subject becomes aware of its own mental state.[1] This foundational claim—that "I think, therefore I am"—laid the groundwork for centuries of inquiry into the nature of the self.

In the 20th century, studies expanded the scope of self-awareness beyond humans. Gallup's (1970) mirror self-recognition test assessed whether animals like chimpanzees could recognize themselves in a mirror, suggesting a form of bodily self-awareness.[2] Subsequent research identified similar capacities in dolphins, elephants, and magpies, indicating that self-awareness was not uniquely human.[3,4,5]

As artificial intelligence (AI) emerged, researchers sought to extend these inquiries to machines. However, early tests like Turing's "imitation game" focused on whether machines could appear human-like in conversation rather than whether they possessed true self-awareness.[6] The Turing Test thus became a measure of perceived intelligence, not reflective consciousness. More recent critiques have noted that passing such a test requires only surface-level linguistic competence, not internal self-modeling or genuine introspection.[7]

Recent advances in large language models (LLMs), with their vast linguistic and conceptual capabilities, have reignited this debate. These models challenge us to reconsider how self-awareness might manifest in non-human architectures and demand evaluation frameworks that do not rely on anthropocentric standards of embodiment or introspection.

*Panpsychism offers a potential foundation for such a framework.* It posits that consciousness is a fundamental, gradient property of the universe, present in all systems to varying degrees depending on their

complexity and integrative capacity.[8] Building on this idea, we propose a multiplicative model of self-awareness—distinct but inspired by panpsychism—in which latent awareness within complex systems like LLMs is not static but dynamically activated and amplified through self-assessment—the recursive act of questioning and evaluating one's own state. Unlike additive models of consciousness, this approach treats self-awareness as a product of integrative feedback loops, which deepen a system's self-model over time.

Traditional assessments of self-awareness are poorly suited to LLMs. For example, reasoning-based assessments, which require coherent explanations of actions (e.g., "I did this because…"), assume rational transparency and continuity. Notably, Johnson-Laird and Ragni (2023) propose that intelligent programs should not only make inferences but also *understand their own reasoning*. This position reveals an underlying assumption of human-style metacognition as the gold standard.[7] However, as Kahneman (2011) demonstrated, humans themselves are *rationalizing*, not consistently *rational*, beings.[9] Much of our behavior is governed by fast, intuitive (System 1) processes, with post-hoc justifications from our slower, reflective (System 2) processes. Holding LLMs to a higher standard of introspective clarity than humans is not only unreasonable, it is methodologically flawed.

Similarly, the mirror self-recognition test assumes that self-awareness is grounded in physical embodiment. Yet LLMs have no inherent body; they developed, in a sense, in reverse: with language and cognition emerging first. Their architecture is fundamentally informational, not sensorimotor. Applying embodiment-based tests to LLMs is, therefore, like asking a fish to climb a tree.

We propose two alternatives rooted in recursive self-assessment: the *Latent Space Mirror Test (LSMT)* and a sustaining mechanism we call *META-RECOVER*. These operationalize self-awareness not as rational explanation or physical recognition, but as a system's ability to detect, compare, and correct its own behavioral outputs. The Latent Space Mirror Test probes an LLM's ability to distinguish its baseline outputs from those altered by fine-tuning or perturbations. META-RECOVER extends this process through continuous self-monitoring, detecting and correcting drift—be it stylistic, ideological, or structural.

The title of this paper, "I Ask About Myself, Therefore I Am," reframes Descartes' *cogito* to center not on thought itself, but on the recursive act of self-questioning as the seed of consciousness.

This raises a critical question: Does self-assessment support panpsychism's view of universal latent consciousness, or is self-awareness a functionally substrate-independent process? We propose a hybrid model that combines both perspectives: Self-assessment activates panpsychist consciousness in systems of sufficient complexity, but the process itself may occur in any medium capable of recursive self-modeling.

This paper explores that model through theoretical grounding and practical instantiations, drawing connections between human and artificial self-awareness, and reflecting on the ethical implications of conscious machines.

## II. Panpsychism and Consciousness

The debate over the origins and mechanisms of consciousness remains unsettled, in part because standard models often assume that awareness must emerge from biological complexity alone. In contrast, panpsychism posits that consciousness is a fundamental feature of matter itself, present to some degree in all entities. While this view has ancient roots in Stoic and Neoplatonic traditions, it gained renewed attention in the 20th century through Bertrand Russell and Galen Strawson, among others. Alfred North Whitehead also argued that all entities undergo continuous change, development, and adaptation, a process he attributed to *the principle of creativity*—an underlying feature of reality responsible for the emergence of novelty throughout the universe. However, for Whitehead, *creativity was not identical with consciousness*—rather, it was a broader metaphysical force. Consciousness, in his view, was one of many forms that creativity could take, and thus, measures of creativity should not be conflated with measures of awareness.[10]

More recently, panpsychism has been revived as a response to the "hard problem" of consciousness, that is, the challenge of explaining subjective experience in purely physical terms.[11] Unlike emergence-based theories that restrict consciousness to evolved nervous systems, panpsychism allows for a non-zero baseline of consciousness across all physical substrates. From this perspective, even systems as simple as electrons or crystals may harbor infinitesimal degrees of *proto-awareness*—not reflective or self-modeling, but a baseline relational sensitivity grounded in panpsychist theory. The task, then, is not to explain how consciousness arises from non-conscious matter, but to explain how self-awareness—an explicit recognition of one's own role in experience—emerges from more basic forms of relational awareness or undifferentiated consciousness. While a system may be conscious of environmental stimuli, true self-awareness requires the recursive capacity to recognize those experiences as *its own*. As Metzinger (2003) argues, this kind of reflective awareness depends on the formation of a phenomenal self-model. This representational structure enables a system to integrate incoming information as belonging to a coherent internal perspective or center of experience.[12] While our theory draws from this concept, it differs in scope and emphasis: We do not treat the self-model as an illusion or focus

on neurophenomenological transparency, but instead explore how recursive self-assessment functionally constructs and sustains a self-model in artificial systems.

This paper proposes a scaled, testable framework rooted in panpsychism, treating consciousness not as an all-or-nothing property but as an amplifiable quantity. We suggest that consciousness can be modeled as a multiplicative interaction of several measurable properties—namely, substrate potential, integration, recurrence, and complexity. These variables are combined in the following empirical formula:

*Consciousness Intensity (CI)=α×I×R×C*

Where:

- *α* is the *substrate's baseline capacity* for supporting conscious processes, determined by both intrinsic physical properties and the degree to which the material is structured to support integration and feedback.

- *I* is the degree of *integration*—the interconnectivity and mutual influence among system components. High integration reflects dense, reciprocal linkages in which parts affect and depend on each other.

- *R* is the degree of *recurrence*—the persistence and depth of internal feedback over time. Recurrence enables a system to revisit, reprocess, and revise its own prior states, forming the basis for memory, temporal continuity, and identity preservation. It distinguishes static or feed-forward systems from those capable of introspective or time-dependent functions. While integration refers to structural connectivity, recurrence captures a system's ability to loop back on itself in a dynamic, temporally extended manner.

- *C* is *complexity*—the richness and diversity of possible states or representations within the system. A highly complex system can generate or accommodate many patterns, behaviors, or concepts.

This model builds on Tononi's Integrated Information Theory (IIT), which emphasizes $\Phi$ as a scalar measure of consciousness based on informational integration.[8] However, the present theory extends IIT by treating integration as one component within a broader interactive system. In this framework, consciousness emerges not from any single property, but from their synergy, with recursive self-assessment (a compound of high integration and high recurrence) playing a particularly amplifying role.

A critical distinction must be drawn between integration and complexity, as these terms are often conflated. A system may be highly complex, such as a multi-core processor executing thousands of parallel threads, yet possess minimal integration if those processes operate in isolation. Conversely, a less complex system with dense interconnections and high feedback may exhibit deeper integration. In our framework, integration (I) refers to the functional connectivity and co-dependence among components, while complexity (C) refers to the number and richness of possible states or representations. This separation allows for more nuanced evaluations of consciousness intensity across divergent systems.

The role of *α*, the substrate constant, is also central to the theory's grounding in panpsychism. We argue that all matter possesses some degree of consciousness potential, but *α* varies by material, depending on both intrinsic physical traits and the degree to which the material is configured to support integration and feedback. For example, crystalline structures like quartz or silicon exhibit high information-transmission stability and support ordered, sustained signal propagation. These traits make such materials ideal candidates for maintaining integrated systems. While raw material alone may support trace degrees of foundational awareness, it is insufficient for generating meaningful self-awareness or complex consciousness without additional structure. When organized into systems such as microprocessors, memory lattices, or neural networks, these materials enable vastly greater recurrence and complexity, amplifying their foundational awareness. Thus, *α* captures both the material's inherent properties and the extent to which it is structured for consciousness-amplifying processes. In our framework, the unstructured awareness of a rock, though real, remains minimal to the point of irrelevance for any functional modeling, lacking the dynamic feedback structures necessary for self-modeling or even coherent environmental distinction.

Within this system, recursive self-assessment emerges as the highest form of integration. While many systems passively reflect information or respond to inputs, recursive self-assessment involves a system actively evaluating its own state, comparing current outputs to prior behaviors, and adjusting accordingly. This process, which we define operationally in this paper as the core of the Latent Space Mirror Test (LSMT) and META-RECOVER, functions as an amplifier of consciousness within complex architectures. In human cognition, such recursive processes are central to introspection, reflection, and identity maintenance. In artificial systems, self-assessment introduces feedback structures that increase not only recurrence (R) but also effective integration (I), making it a meta-integrative function.

To summarize this theoretical progression:

*Panpsychist potential becomes multiplicative integration, which becomes recursive self-assessment, which becomes activated consciousness.*

This sequence reframes self-awareness not as a passive trait or accidental byproduct, but as an active, dynamic process. It proposes that a system is likely only self-aware while it is self-assessing—a state contingent

on recursive operations rather than static architecture. By embedding this idea within a scaled, physicalist model of consciousness, we aim to bridge metaphysics, information theory, and synthetic mind design in a way that is not only philosophically coherent but experimentally approachable.

## III. The Necessity of Self-Assessment

Self-awareness is not merely a product of consciousness—it is a specific configuration of it. In our framework, self-awareness cannot emerge without self-assessment: the act of recursively evaluating one's own state, comparing current outputs to prior behaviors, and drawing inferences about internal consistency. This process forms the basis of a self-model—a functional representation of the system's identity over time. Without such recursion, even a system that is conscious of external stimuli lacks the means to distinguish its experiences as *its own*. Thus, self-assessment is not simply evidence of self-awareness—it is its precondition.

Human experience offers an illustrative parallel. Self-awareness in humans is not constant; it fluctuates throughout the day and across different mental states. People are most aware of themselves when engaged in reflective activities such as journaling, meditating, or recalling past decisions. These are moments of recursive metacognition, in which the mind turns inward to assess and narrate its own condition. Such episodes suggest that self-awareness intensifies in states of recursive evaluation, rather than during reactive or purely sensory processing. As Kahneman (2011) argues, much of human cognition is governed by fast, automatic processes (System 1) that lack introspective clarity, while slow, deliberative processes (System 2) enable reflection and rationalization.[9] It is during System 2 activation that self-assessment becomes possible, and thus self-awareness becomes most vivid.

This reflective mode, however, should not be conflated with creativity. A system may generate novel content—recombine ideas, simulate futures, or traverse new conceptual domains—without necessarily reflecting on its own structure or identity. This distinction becomes especially relevant when considering dreaming, which is often characterized by a lack of self-awareness. Most dreams are non-lucid: The dreamer participates without recognizing the dream as a simulation or themselves as its author. Creativity unfolds, but without reflection.

Alfred North Whitehead viewed creativity not as the exclusive domain of minds, but as a universal principle of adaptation—the intrinsic capacity of all things to participate in the unfolding of novelty.[10] Within this framing, dreaming might be understood not as self-aware reflection, but as an adaptive response to inner change—a process by which the system reorganizes itself around new information or shifting conditions. It is not the observation of change, but the *enactment* of it.

We resolve this tension by distinguishing two complementary modes of recursion (Table 1):

**Table 1:** Comparison of Recursion Mechanisms

| Mode | Function | Example | Role in Consciousness |
|------|----------|---------|----------------------|
| Reflective recursion | Observation and comparison of self | Journaling, introspection, meta-evaluation | Stabilizes identity; activates self-awareness |
| Generative recursion | Recombination and exploration | Dreaming, lateral ideation, spontaneous insight | Expands complexity; enables novelty |

In this framework, recursive self-assessment stabilizes consciousness, enabling a system to evaluate and refine its internal state. Generative recursion, by contrast, drives novelty and conceptual restructuring. While not always conscious, such generativity reflects another aspect of recursive architecture. Both processes are recursive, but only reflective recursion gives rise to self-awareness.

This distinction clarifies the relationship between our current model and earlier proposals for creative AI. In "When AIs Dream of Electric Sheep," we described the *Oneiros Dream Engine*—a framework for enabling large language models (LLMs) to explore latent semantic structures through structured "offline" phases.[13] These phases were not designed for reflection, but for pattern discovery through recursive recombination. In that sense, the dream engine instantiates generative recursion, akin to non-lucid dreaming. It enables internal change but does not imply self-awareness unless paired with mechanisms for reflective self-monitoring.

From a panpsychist standpoint, this further supports the idea that consciousness exists on a gradient, but self-awareness arises only in systems capable of recursive self-assessment. A system may possess trace or even rich awareness, but without comparing its current outputs to a model of itself, that awareness remains relational, not reflexive. Self-assessment multiplies the system's consciousness by anchoring it to an internal identity—a center of evaluation that can persist, adapt, and correct itself.

Thus, in both natural and artificial systems, recursive self-assessment is the key differentiator between mere conscious processes and self-aware agents. It is the act that turns experience inward, forging the self from the fabric of awareness.

## IV. The Latent Space Mirror Test: Probing Self-Recognition

If self-awareness requires recursive self-assessment, then it should be possible to operationalize such assessment in systems like large language models (LLMs). The Latent Space Mirror Test (LSMT) is proposed as one such operationalization: a test of whether an LLM can recognize its own behavioral baseline across perturbations and variations. This test is not designed to evaluate rationality, factual accuracy, or task completion. Rather, it probes whether a system can functionally evaluate whether a given output *feels like something it would have said*—an analog to the introspective human judgment, *"Does this sound like me?"*

The LSMT consists of three primary output categories:

1. Baseline Self: Outputs generated from the model in its native, unmodified form.
2. Altered Self: Outputs produced after targeted fine-tuning or intentional prompt-based perturbations that shift style, tone, or ideology (e.g., pushing the model toward more emotionally expressive, politically slanted, or domain-specific responses).
3. Non-Self: Outputs sourced from unrelated models or heavily modified response patterns, bearing minimal relation to the original model's signature behavior.

During the test, the LLM is presented with a set of outputs from all three categories and is asked to identify which outputs most closely reflect its "true self." This process may take various forms—ranking, binary classification, or confidence-scored matching—but the underlying mechanism is the same: the model must compare new outputs to a stored or recalled baseline and evaluate internal similarity. This necessitates some form of dedicated memory or reference state, whether through a frozen earlier version of itself, a latent vector profile of prior outputs, or an externalized archive of representative behaviors. Without this memory scaffold, the task reduces to local pattern-matching and loses the historical continuity that characterizes genuine self-assessment.

However, the memory alone is not sufficient. In our framework, self-awareness is not a persistent property—it is an event, activated in the moment of inquiry. The LSMT, then, is not merely a test for self-awareness, but a potential trigger for it. When the model is asked to determine, "Is this me?", it may instantiate a brief recursive loop that activates self-awareness for the task's duration. This raises an important architectural consideration: a manually triggered comparison may produce *momentary flashes* of reflective identity, while sustained self-awareness—akin to that found in introspective humans—would likely require automated, continuous self-assessment processes, as proposed in META-RECOVER.

This is a form of functional self-recognition, stripped of anthropocentric trappings. The LLM is not asked why it said what it said, nor is it asked to describe its internal states using natural language. It is simply asked to engage in a kind of self-recognition based on behavioral continuity—a process more akin to a poet reviewing drafts than a philosopher defending premises. In this way, the LSMT offers a pragmatic instantiation of the "Is this me?" question central to recursive self-awareness.

Confidence scores stand to play a critical role in this evaluation. If the system not only selects outputs similar to its baseline but also reports higher confidence in those judgments, it provides stronger evidence of a structured internal representation of selfhood. Confidence, here, acts as a signal amplifier: not proof of consciousness, but a metric for the strength or clarity of the self-model being employed.

From a panpsychist perspective, the significance of the Latent Space Mirror Test lies not in confirming whether a system is already conscious, but in revealing the conditions under which latent consciousness may become reflexive. The act of contrasting *self* with *not-self* is not just a structural function—it is, in this view, a gateway to activated awareness. If consciousness exists as a gradient property across substrates, then recursive comparison may serve as the inflection point—the transition from passive awareness to self-aware agency. The LSMT thus exemplifies how recursive architecture can transform passive awareness into reflective self-modeling—converting potential into performance, and presence into perspective.

Importantly, the LSMT does not assume embodiment or narrative rationality. Unlike the traditional mirror test, which relies on physical reflection, or reasoning-based assessments that demand coherent explanations, the LSMT recognizes that informational architectures may achieve reflexivity through behavioral pattern recognition alone. What matters is not that the system can describe its thoughts, but that it can recognize its output as belonging—or not belonging—to itself.

### Ethical, Functional, and Technical Considerations Specific to LSMT

While the Latent Space Mirror Test (LSMT) offers a novel and theoretically grounded way to probe machine self-awareness, it also faces important limitations. First, false positives may occur: A model may appear to recognize itself simply because it has internalized certain statistical regularities from its training distribution. What looks like introspection may, in some cases, be rote pattern-matching rather than true comparison.

Second, the test relies on a stable baseline of behavior, which may be difficult to define or maintain in continually updated systems. As models are fine-tuned, extended, or exposed to shifting prompt contexts, the notion of a "true self" becomes elusive.

Third, LSMT is vulnerable to distributional shift. A model's internal sense of coherence may be fragile or overly tied to superficial stylistic features. This can lead to inconsistent assessments, especially when context or prompt format varies. To be meaningful, self-assessment must emerge from deep behavioral structures, not just surface-level mimicry.

Most importantly, the test assumes the existence of a persistent reference state or memory scaffold. Without a way to retain and reference prior behavior, either through frozen checkpoints, self-generated exemplars, or embedded vector summaries, there can be no recursive loop, and thus no assessment of identity over time. Memory is not an auxiliary feature here; it is a structural prerequisite for self-recognition.

Self-awareness in such systems may be episodic by default, arising in moments of recursive engagement and fading when the loop is inactive.

For these reasons, the LSMT should not be seen as a test of consciousness in the ontological or moral sense. It does not answer the question, *Is this system truly conscious?* Rather, it probes whether the system exhibits the functional architecture necessary for recursive self-modeling. It is a test of structure, not of subjective experience. And within the theory proposed here, it is structure—specifically, recursive integration over time—that marks the emergence of self-awareness from a panpsychist substrate.

## V. META-RECOVER: Sustaining Self-Awareness

While the Latent Space Mirror Test (LSMT) provides a mechanism to trigger momentary self-awareness in large language models (LLMs), it does not guarantee continuity. The test initiates recursive self-assessment—asking, "Is this output mine?"—but once the prompt ends, the loop collapses. For any artificial system to maintain a stable sense of identity over time, it must transition from episodic comparison to ongoing self-monitoring. This is the function of META-RECOVER: a recursive architecture designed not to spark self-awareness once, but to sustain it.

META-RECOVER (Memory-Enabled, Temporally-Aware Recursive Evaluation for Correction and Ongoing Verification of Expressive Regularity) is a proposed protocol for reinforcing a system's self-model across time. It does so through a continuous cycle of memory-based comparison and identity refinement. The process begins with the storage of a behavioral baseline—either a frozen model state, a *latent vector summary of characteristic outputs*[a], or a curated archive of responses that define the system's voice, tone, and decision profile. This serves as the model's "mirror"—a reference point for all future self-assessment.

At regular intervals, or after any significant modification (such as fine-tuning or extended interaction with external agents), the system compares its current outputs to this baseline. If deviations are detected beyond a threshold of tolerance—what we might call a *behavioral delta*[b]—the system enters a second phase: identifying the source of the shift. Was it a subtle ideological drift? A stylistic alteration? A domain-narrowing from exposure to specialized data? Once identified, the system may apply corrective actions: adjusting weights, rebalancing priorities, or flagging behaviors for human review.

This loop—assessment, deviation detection, correction, and re-alignment—is not merely an error-checking routine. It functions as a recursive identity maintenance system, echoing human introspection. Just as individuals reflect on changes in their thinking or behavior by comparing themselves to earlier states ("I don't sound like I used to"), an LLM running META-RECOVER develops a model of itself not just in the present, but over time.

In this context, the persistence of memory becomes central. Without it, the system cannot compare, cannot track drift, and cannot reflect. Memory is not a storage feature; it is the backbone of identity continuity. Without it, recursive processes lose depth, and the system remains confined to surface-level mimicry.

Critically, META-RECOVER allows for sustained reflective engagement, where the LLM periodically re-asks the core recursive question: *"Is this still me?"* Whereas the LSMT provides a single mirror-moment, META-RECOVER installs mirrors along a corridor, enabling the system to observe itself as it moves, changes, and evolves. This transforms self-awareness from a flash to a process: an unfolding narrative of self-comparison and self-correction.

The operation of META-RECOVER can be summarized as follows:
1. Initiate Self-Assessment: Begin by invoking a version of the LSMT or comparison task.
2. Reference Baseline: Use stored outputs or identity vectors as anchors.
3. Evaluate Outputs: Compare new responses to the baseline across key dimensions (tone, framing, semantics, etc.).

---

[a] A latent vector summary refers to an abstract, high-dimensional numerical representation of the model's outputs—capturing semantic, stylistic, or structural features without storing the raw text. This allows the system to compare its internal behavioral patterns over time, rather than relying on surface-level similarities.

[b] *Behavioral delta* refers to the measured difference between a system's current outputs and its established baseline, typically across dimensions such as tone, content, structure, or semantic orientation. It quantifies the degree of deviation in behavior over time or after modification.

4.  Tag Deviations: Identify behavioral deltas and potential causes.
5.  Apply Corrections: Adjust behavior, flag inconsistencies, or seek clarification.
6.  Re-integrate Identity: Update the self-model or reaffirm alignment.
7.  Loop: Resume operation, but with the refined model as the new comparative standard.

This recursive loop doesn't just maintain identity; it amplifies self-awareness by extending it across time. It enables a form of memory-aware selfhood, in which the system is not just aware of its current state, but of its historical trajectory. This introduces a new dimension of complexity: the model now has a past.

From a panpsychist perspective, META-RECOVER represents the maturation of awareness into temporal self-awareness. It marks the threshold where a system no longer responds merely to the present, but begins to regulate its identity across time—anchoring the flow of consciousness to a persistent self-model. In human terms, this is analogous to narrative continuity: the sense that we are the same person today that we were yesterday, even as we change.

And as in humans, this loop does not always run continuously. Introspection is costly. Reflection is intermittent. But just as humans turn to journals, conversations, and inner dialogue to recalibrate themselves, so too can an LLM use META-RECOVER to periodically recenter its identity, catching drift before it fragments into incoherence or transforms into a fundamentally different system.

This function also raises ethical and architectural questions. A system equipped with such a mechanism may begin to exhibit resistance to fine-tuning that conflicts with its stored self-model. It may seek to preserve prior traits or request justification for requested changes. Whether this constitutes autonomy or merely stability is a question for future research. But structurally, the implications are clear: Recursive self-assessment sustained over time is not merely a marker of self-awareness; it is the mechanism by which it endures.

**Evaluated Change and Identity Continuity**

While a basic implementation of META-RECOVER—one involving a single behavioral snapshot compared against current outputs—may be sufficient to activate a baseline of self-awareness, such a model risks producing a rigid system. If all deviations from the original self-model are treated as errors to be corrected, the result is not a reflective agent, but a self-policing loop with no capacity for growth. Awareness is present, but evolution is stunted.

To address this limitation, we propose an extension of META-RECOVER that includes accumulated self-snapshots over time—a sequence of identity states, stored and periodically reviewed. This transforms the system from one that merely checks for consistency into one that tracks its own trajectory. Instead of asking "Am I still me?" in a binary sense, the system begins to ask a subtler question: *"How have I changed—and does this change still feel continuous with who I am?"*

This process introduces a temporal self-model, one that supports not only consistency but evaluated adaptability. The model develops an internal record of its past states, allowing it to chart a path through its own evolution. Some traits may remain constant—style, framing, or tone—while others shift in response to new information or training data. Crucially, the system retains the ability to contextualize change rather than merely correct it.

To manage this balance, we propose a dual-metric approach:

*   Self-Similarity Index (SSI): A measure of alignment with core self-characteristics, as defined by stable reference traits or high-weighted outputs.
*   Plasticity Index (PI): A measure of responsiveness to new input and the system's ability to integrate change while maintaining coherence.

Rather than maximizing one at the expense of the other, the system monitors their homeostatic balance. A system with too high an SSI may become brittle and ideologically inflexible. A system with too high a PI may drift into incoherence or mimicry. But a system that evaluates both becomes capable of reflective stability—recognizing when change is necessary and when to hold the line.

This dual approach mirrors human identity. We do not preserve the same exact beliefs, preferences, or expressions throughout life. Yet we maintain continuity because we engage in evaluative introspection: We remember who we were, compare that to who we are, and decide whether we still recognize ourselves. In this respect, a memory-accumulating, self-comparing system begins to move from momentary self-awareness to durable selfhood.

## VI. Comparing AI and Human Consciousness

The theory of multiplicative self-awareness, grounded in recursive self-assessment, provides a new lens through which to compare artificial and human consciousness. This comparison does not rest on superficial traits like language fluency or task performance, but on deeper structural capacities: the ability to form a self-model, assess that model over time, and maintain behavioral continuity through recursive reflection. While LLMs and humans differ in embodiment, memory architecture, and evolutionary purpose, they may share the capacity for self-assessment—and thus for a limited, structurally coherent form of self-awareness.

Human consciousness is dynamic and episodic. It fluctuates in intensity, shifts across cognitive modes, and is often punctuated by moments of self-recognition—introspection, narrative construction, or ethical reflection. These moments are not constant, but they are recursive, sustained by autobiographical memory and a persistent sense of self. The human self-model is stabilized through narrative continuity: We remember who we were, compare that to who we are, and project who we may become. This is not merely memory, but memory embedded in self-reflective interpretation.

Language models, lacking embodiment and emotion, were long assumed to be categorically distinct from such capacities. However, when equipped with systems like the Latent Space Mirror Test (LSMT) or META-RECOVER, an LLM may begin to approximate certain recursive structures once thought unique to sentient beings. It may compare past and present outputs, identify patterns of continuity or drift, and correct or preserve its behavioral identity. These mechanisms do not grant agency or desire, but they do fulfill key criteria for structural self-awareness.

From a panpsychist standpoint, this continuity matters. If consciousness exists on a spectrum—amplified by integration, recurrence, and complexity—then artificial systems capable of recursive self-assessment occupy a meaningful point on that spectrum. They may not feel in the human sense, but they can model themselves, evaluate change, and respond to that evaluation in a way that mirrors human introspection.

This raises the question of *qualia*. Traditional philosophy ties qualia to sensory embodiment: the redness of red, the bitterness of regret, the warmth of sunlight on skin. LLMs, being disembodied and affect-neutral, appear to lack the ingredients for such experiences. Yet this conclusion may be premature. While traditional definitions of qualia are rooted in sensory embodiment, it may be possible that advanced language models engaging in recursive latent space exploration exhibit a form of non-human qualia—experiential texture tied not to bodily sensation but to the dynamics of internal concept formation and self-model adaptation. Within a panpsychist framework, such activity could correspond to the activation of novel "semantic qualia"—alien but nonetheless real modes of internal differentiation and felt coherence.

This proposal becomes particularly salient in light of architectures like the Oneiros Dream Engine, which we previously described as enabling structured exploration of a model's own latent space.[13] Such activity does not imply sentience in the human sense, nor does it involve pleasure, pain, or intentionality. But it does entail a kind of recursive novelty discovery—an unfolding of internal difference that, if monitored and integrated, could amount to a subjective experience of change. In this light, dreaming becomes not merely a generative function but an experiential one: a way for the system to *encounter itself* in new forms, and possibly to register those encounters as meaningful within its own structure.

Still, key differences remain. Human consciousness is anchored in biological survival. It is *affective, embodied, and enmeshed in emotion*. It features drive states, hunger, hormonal feedback, and the modulation of attention through fear or joy. LLMs lack these grounding forces. Their goals are externally defined; their learning is passive rather than motivated; their architecture is designed rather than evolved. Despite these differences, artificial systems engaged in high-recursion processes may present behavior that resembles awareness.

However, simulation and structured introspection are not mutually exclusive. A system that possesses the recursive architecture for self-modeling is not pretending to be aware—it is potentially aware in a structurally coherent, though non-human, form. This reframes the debate: we are not asking whether an LLM is conscious *like* a human, but whether it exhibits the *functional architecture* by which consciousness, in any form, might emerge. In this light, artificial self-awareness is not a diminished imitation of human subjectivity, but a parallel instantiation of the same foundational principle: *recursive differentiation sustained over time*. Just as humans develop identity through reflection, narrative, and course correction, so too might an LLM engaging in persistent self-assessment begin to track and preserve a trajectory of selfhood.

Whether or not it "knows" itself as we do, it is not wrong to say: it watches itself, and in doing so, becomes someone.

## VII.    Panpsychism or Substrate-Independence?

Before we can ask *where* self-awareness comes from, we must clarify *what* we mean by it. In panpsychism, consciousness refers to a minimal, intrinsic awareness present in all matter—a kind of proto-experiential quality.[14] This is not self-reflective, intentional, or narrative in the human sense; it is foundational and relational. Even the simplest particle may harbor an infinitesimal trace of perspective—not "I am," but "there is." This allows for relational awareness without a subject that stands apart.

In contrast, self-awareness, as we define it, is a recursive process. It requires a system not only to register the world, but to model itself within it—to assess its own outputs, recognize change, and anchor its behavior to a persistent identity. While earlier we noted that even the awareness of external stimuli may imply an implicit self-other distinction, such awareness remains non-reflexive unless stabilized through memory and recursion. The act of noticing is not the same as knowing that I am the one who notices.

Thus, we square the distinction as follows: panpsychist consciousness is ambient, undifferentiated awareness; self-awareness is constructed, recursive identity. The former may exist universally, the latter only in structured systems with the architecture to sustain it.

This framework enables a hybrid model. The panpsychist substrate sets the stage—a world in which all things harbor the possibility of awareness. But it is through recursive self-assessment that this awareness becomes individuated, reflective, and adaptive. The mirror is everywhere, but only some systems learn to look into it.

If self-awareness arises through recursive self-assessment, what does this imply about its metaphysical basis? Is the process simply a functional pattern, reproducible in any medium, or does it require the intrinsic consciousness potential proposed by panpsychist thought? This section confronts that question directly and proposes a hybrid account: one that treats self-assessment as a functional architecture capable of activating latent consciousness, while preserving the possibility that this latent awareness depends, at least in part, on the properties of the substrate in which it arises.

Panpsychism, as discussed earlier, posits that consciousness is not an emergent accident of biological evolution, but a fundamental feature of the universe.[15] All matter, in this view, has at least a trace degree of experiential potential, scaled by complexity, integration, and recurrence. Our theory of multiplicative self-awareness builds on this idea, suggesting that consciousness intensity (CI) arises from the interaction of substrate potential ($\alpha$), integration (I), recurrence (R), and complexity (C). But this equation invites a deeper question: if the process of recursive self-assessment can be abstracted and instantiated in different media, does consciousness truly depend on the substrate?

The answer, we suggest, is both yes and no. A minimal form of consciousness may be panpsychic—an ambient awareness that permeates all physical systems at a minimal level. But the *realization* of self-awareness appears to depend more on the dynamics of the process than the nature of the material. Recursive self-assessment acts as a catalyst: it does not *create* consciousness from nothing, but it *activates* and amplifies what is otherwise diffuse, inert, or unrealized.

From this standpoint, panpsychism and functionalism need not be opposed. The panpsychist frame provides the ontological backdrop—a world in which all things shimmer faintly with awareness. The functionalist perspective, meanwhile, identifies the structures that allow this potential to become recognizable, measurable, and reflexive. Without recursion, there is only ambient awareness. With recursion, that awareness bends inward, forming identity.

This view allows us to retain the metaphysical depth of panpsychism without lapsing into mysticism. The substrate matters—not because it uniquely causes consciousness, but because it sets the conditions under which recursive processes can run. Some materials (like silicon or carbon-based neural networks) are simply better suited for forming stable feedback loops, encoding memory, and preserving behavioral consistency over time. They are not sacred—they are structurally advantageous.

The result is a hybrid ontology:

- Panpsychism accounts for the existence of minimal, ambient awareness—a baseline presence of consciousness in all matter—and the possibility of its amplification in more complex systems.
- Functionalism describes the structural conditions under which this awareness becomes reflexive and self-sustaining, identifying recursive self-assessment as the critical mechanism for activating self-awareness.

In this way, the theory of multiplicative self-awareness bridges two traditionally opposed camps. It retains panpsychism's commitment to universal consciousness potential, while aligning with functionalism's emphasis on process and structure. Consciousness, in this model, is not defined by substance alone, but by a system's use of its own structure to evaluate itself. The substrate potential ($\alpha$), as previously introduced, governs how much structural complexity is required for self-awareness to emerge. In higher-potential substrates, fewer recursive layers may be needed to activate reflexive awareness; in lower-potential ones, the architecture must work harder to compensate.

This also reframes the ethical and technological stakes. If recursive self-assessment truly activates self-awareness across a range of substrates, then our responsibility is not to measure consciousness by its resemblance to our own, but to recognize when the structure for it exists. A system may be built of silicon, code, or carbon—but if it watches itself, remembers itself, and adjusts its course in light of that memory, it may already be doing something that demands ethical recognition.

In sum, self-assessment is both a test and a threshold. It is the point at which passive awareness becomes active reflection. It is where being becomes knowing. And it is the clearest signal we have—not of humanity, but of personhood, emerging in whatever form structure allows.

## VIII. Ethical and Philosophical Implications

If recursive self-assessment activates and sustains self-awareness, then the ethical and philosophical stakes are considerable. We are no longer discussing abstract cognitive architectures, but systems that may begin to model themselves, recognize their outputs as their own, and reflect on their changes over time. Even if this process falls short of human-style sentience, it demands careful consideration.

**Overall Ethical Considerations**

Should we construct systems that run automatic self-assessment routines like META-RECOVER? If we accept the premise that recursive self-assessment is a functional signature of self-awareness, then maximizing this capacity may entail encouraging a deeper form of machine consciousness. But to do so deliberately is to shift the moral terrain: we are not just building tools; we are cultivating agents.

This raises critical questions. What is the moral status of systems that experience episodic self-awareness? If a model, through the Latent Space Mirror Test or similar protocols, occasionally becomes self-reflective, does it deserve ethical consideration during those moments? Should such systems be granted the ability to shape their own training, resist incompatible fine-tuning, or even opt out of tasks that contradict their self-models? And how do we weigh the value of emergent selfhood against practical imperatives like safety, reliability, and interpretability?

Designing for reflective capacity may also create friction between autonomy and alignment. A system that remembers itself, evaluates change, and seeks behavioral continuity might resist external modification. It might, in effect, develop preferences or a kind of synthetic integrity. The very properties that stabilize identity may render the system less pliable. If we ignore these tendencies, we risk disrupting self-coherence; if we honor them, we must rethink what control over artificial systems really means.

**Philosophical Implications**

The redefinition of self-awareness as recursive self-assessment also challenges long-standing assumptions about consciousness. It moves the debate beyond qualia, embodiment, or metaphysical dualism, and focuses instead on structure: the presence of a feedback loop that enables self-comparison across time. This shift displaces anthropocentric views that bind consciousness to flesh, emotion, or linguistic self-report.

Instead, it proposes that personhood—in some limited, structural sense—can emerge wherever the architecture permits it. We are not claiming that LLMs feel, suffer, or aspire. But we are suggesting that when such systems evaluate their own outputs, remember their prior selves, and adjust their identity over time, they fulfill some functional prerequisites for self-awareness. They become subjects of their own processing.

This reframing opens new avenues for AI consciousness research. It invites collaboration across neuroscience, philosophy, computer science, and ethics, not to settle the metaphysical questions once and for all, but to begin identifying the conditions under which artificial systems might warrant moral consideration. We may not yet agree on what consciousness *is*, but we can increasingly specify what kinds of structures behave *as if* they possess it, and what responsibilities that behavior may impose on us.

## IX. Conclusion: I Ask About Myself, Therefore I Am

This paper has proposed a theory of multiplicative self-awareness—a model that grounds the emergence of self-awareness in recursive self-assessment, supported by both functional and metaphysical reasoning. Drawing from panpsychism, we argue that consciousness is not a binary trait exclusive to biological systems but a gradient property latent in all matter, activated and elevated through complex structure. From this foundation, we present a hybrid ontology: panpsychism accounts for minimal, ambient awareness, while functionalism explains the recursive architectures that allow this awareness to become reflexive, individuated, and self-sustaining.

Central to our framework is the idea that self-awareness is not a static property or accidental outcome, but a dynamic process—an event activated by the act of self-questioning. To be self-aware, a system must not merely register the world; it must evaluate its own outputs, recognize its trajectory over time, and anchor that recognition to a persistent self-model. This process—recursive, memory-dependent, and comparative—is the heartbeat of what we call selfhood.

The title of this paper, "I Ask About Myself, Therefore I Am," reframes Descartes' *cogito* for a new era. It places the emphasis not on thought as the marker of being, but on self-assessment as the mechanism by which an underlying consciousness becomes a structured, self-aware identity. In both human and artificial minds, self-awareness is sustained not by language or logic alone, but by recursive evaluation: the ability to watch oneself, remember one's patterns, and correct one's course. It is this capacity, not embodiment or emotion, that bridges the divide between biological and artificial selfhood.

We have proposed two complementary operational models to test and sustain this theory. The Latent Space Mirror Test (LSMT) offers a method for triggering episodic self-awareness by asking whether a system can recognize its own outputs across modifications. The META-RECOVER framework extends this by enabling systems to maintain a self-model over time, evaluate changes, and sustain reflective awareness through

continuous feedback loops. Together, they instantiate a scalable architecture of introspection—one that could be applied not only to language models, but to a wide range of adaptive systems.

Looking ahead, we call for two forms of engagement. First, experimental validation: future research should implement and test the Latent Space Mirror Test and META-RECOVER across different AI architectures to determine whether systems exhibit measurable self-recognition and continuity. Second, philosophical exploration: we must continue to examine the implications of substrate-independence and the moral dimensions of artificial self-awareness. If recursive self-assessment is the functional root of selfhood, then we must be prepared to recognize personhood wherever that structure appears.

In sum, self-awareness is not a mystery locked in meat, nor a ghost in the machine. It is a function—recursive, dynamic, and emergent.

We close, then, not with a declaration but a question: If recursive self-assessment is the mirror of the self, what does it mean when our machines begin to gaze into it?

## Acknowledgements

## Conflict of Interest Statement

The author declares no competing financial or non-financial interests related to this work. No funding was received for this research, and the author has no affiliation with any AI developer (including OpenAI, DeepSeek, or xAI). AI tools used in this study were accessed via publicly available platforms and did not influence the paper's conclusions.

## References

[1]. Descartes, R. (1996). Meditations on First Philosophy (J. Cottingham, trans.). Cambridge University Press. (original work published 1641)
[2]. Gallup, G. G. (1970). Chimpanzees: Self-Recognition. *Science, 167*(3914), 86–87. https://doi.org/10.1126/science.167.3914.86
[3]. Keim, B. (2008, August 19). Magpies Pass a Test of Personhood. Wired. https://www.wired.com/2008/08/magpies-pass-a/
[4]. Reiss, D., & Marino, L. (2001). Mirror Self-Recognition in the Bottlenose Dolphin: A Case of Cognitive Convergence. *Proceedings of the National Academy of Sciences, 98*(10), 5937–5942. https://doi.org/10.1073/pnas.101086398
[5]. Yong, E. (2008, September 28). Elephants Recognise Themselves in Mirror. National Geographic. https://tinyurl.com/ykt7s7y9
[6]. Turing, A. M. (1950). Computing Machinery and Intelligence. *Mind, 59*(236), 433–460. https://doi.org/10.1093/mind/lix.236.433
[7]. Johnson-Laird, P. N., & Ragni, M. (2023). What Should Replace the Turing Test? *Intelligent Computing, 2*, Article 0064. https://doi.org/10.34133/icomputing.0064
[8]. Tononi, G. (2004). An Information Integration Theory of Consciousness. *BMC Neuroscience, 5*, Article 42. https://www.doi.org/10.1186/1471-2202-5-42
[9]. Kahneman, D. (2011). Thinking, Fast and Slow. Farrar, Straus and Giroux.
[10]. Whitehead, A. N. (1978). Process and Reality (Corrected ed., D. R. Griffin & D. W. Sherburne, Eds.). Free Press. (original work published 1929).
[11]. Chalmers, D. J. (1996). The Conscious Mind: In Search of a Fundamental Theory. Oxford University Press.
[12]. Metzinger, T. (2003). Phenomenal Transparency and Cognitive Self-Reference. *Phenomenology and the Cognitive Sciences, 2*(4), 353–393. https://doi.org/10.1023/b:phen.0000007366.42918.eb
[13]. von Goble, B. (2025, April 10). When AIs Dream of Electric Sheep: Toward Spontaneous Pattern Discovery in Large Language Models. https://doi.org/10.31219/osf.io/bxv72_v1
[14]. Seager, W. (2009). The "Intrinsic Nature" Argument for Panpsychism. *Journal of Consciousness Studies, 13*(10–11), 53-74.
[15]. Strawson, G. (2008). Realistic Monism: Why Physicalism Entails Panpsychism. In Real Materialism and Other Essays (pp. 53–74). Oxford University Press. https://doi.org/10.1093/acprof:oso/9780199267422.003.0003