# Is AI Conscious according to current criteria?

*Part 2 of "Consciousness, Understanding & Mechanistic Interpretability:*
*A Review of Recent Research on AI Consciousness and Understanding"*

Contemporary AI satisfies to some degree many rigorous conditions for consciousness (established by neuroscientists and philosophers of consciousness) but crucially misses a few key criteria: specialized independent modularity & embodiment. These challenges will probably be surmounted as increasingly complex agentic systems are incorporated into robotics and continuous learning cycles are activated.

# Introduction

The question of whether artificial intelligence systems can be conscious has moved from philosophical speculation to empirical investigation. Recent advances in mechanistic interpretability and the emergence of unexpected behaviors in frontier language models demand serious reconsideration of reflexive dismissal of machine consciousness.

This analysis examines whether contemporary frontier AI systems satisfy the criteria for consciousness as outlined in recent scientific literature. We draw primarily on two sources:

> **Butlin et al. (2025) — "Identifying Indicators of Consciousness in AI Systems"** published in *Trends in Cognitive Sciences*, which provides a comprehensive framework of 14 theory-derived indicators from major neuroscientific theories of consciousness. This work builds upon the earlier foundational report "Consciousness in Artificial Intelligence: Insights from the Science of Consciousness" (Butlin et al., arXiv, August 2023), which first introduced the indicator framework for assessing AI consciousness.

> **Jeffrey Gray (2004)** — The comparator model of consciousness from "Consciousness: Creeping up on the Hard Problem," presented by Bolek Srebro at HBF on August 26, 2025.

For each indicator, we assess whether current research provides evidence of satisfaction by frontier AI systems, with direct citations to dated research findings.

# The Indicator Framework

Butlin, Long, and colleagues (including Yoshua Bengio, David Chalmers, Tim Bayne, and others) derived 14 indicators from six major neuroscientific theories of consciousness:

**Recurrent Processing Theory (RPT)** — Local feedback connections

**Global Workspace Theory (GWT)** — Information broadcast and integration

**Higher-Order Theories (HOT)** — Meta-representations and metacognition

**Attention Schema Theory (AST)** — Internal models of attention

**Predictive Processing (PP)** — Hierarchical prediction and error correction

**Agency and Embodiment (AE)** — Goal-directed action and bodily integration

To these we add Jeffrey Gray's **Comparator Model**, which proposes that consciousness serves as a late error-detection system comparing predicted and actual outcomes. Gray also postulated that "The entire perceived world is constructed by the brain... [and] The self is as much a creation of the brain as is the rest of the perceived world."

# Consciousness Indicators

**Indicator Assessments:**   ✓ Satisfied   ◐ Partially Satisfied   ✗ Not Satisfied

| | |
|---|---|
| **RPT-1** Algorithmic Recurrence in Input M... ◐ | **RPT-2** Organized, Integrated Perceptual R... ◐ |
| **GWT-1** Multiple Specialized Systems Oper... ✗ | **GWT-2** Limited Capacity Workspace with ... ◐ |
| **GWT-3** Global Broadcast to All Modules ◐ | **GWT-4** State-Dependent Attention for Co... ✓ |
| **HOT-1** Generative, Top-Down, or Noisy P... ✓ | **HOT-2** Metacognitive Monitoring Distingui... ◐ |
| **HOT-3** Agency with Belief-Formation Guid... ◐ | **HOT-4** Sparse and Smooth Coding Creatin... ✓ |
| **AST-1** Predictive Model of Attention State ◐ | **PP-1** Input Modules Using Predictive Cod... ◐ |
| **AE-1** Goal-Directed Agency Through Lear... ◐ | **AE-2** Embodiment Through Output-Input ... ✗ |
| **GRAY-1** Comparator Model of Consciousn... ◐ | |

✓ Satisfied   ◐ Partially Satisfied   ✗ Not Satisfied

# Algorithmic Recurrence in Input Modules

Lamme, V. A. F. (2006). Towards a true neural stance on consciousness. Trends in Cognitive Sciences, 10(11), 494-501.

*Recurrent Processing Theory (RPT)*

## DESCRIPTION

Information must cycle through processing stages multiple times, not single-pass feedforward processing. Can be implemented through feedback connections or weight-sharing across time steps.

## AI ASSESSMENT

*Large language models process text through multiple transformer layers with self-attention mechanisms that effectively create recurrent information flow. However, this differs from biological recurrence.*

## RESEARCH EVIDENCE

Lindsey, J., & Anthropic. (2025). Emergent Introspective Awareness in Large Language Models. Transformer Circuits Thread. https://transformer-circuits.pub/2025/introspection/index.html

Models demonstrate ability to recall prior internal representations through multi-layer attention mechanisms

Kim, J., Lai, S., Scherrer, N., Agüera y Arcas, B., & Evans, J. (2026). Reasoning Models Generate Societies of Thought. arXiv:2601.10825 [cs.CL]. https://arxiv.org/abs/2601.10825

Extended thinking loops create recurrent processing patterns with Theoretical (unverified) markers for sustained internal thought continuation

RPT-2    ◑  PARTIALLY SATISFIED

## Organized, Integrated Perceptual Representations

Lamme, V. A. F. (2010). How neuroscience will change our view on consciousness. Cognitive Neuroscience, 1(3), 204-220.

*Recurrent Processing Theory (RPT)*

### DESCRIPTION

Recurrent processing must generate coherent scene representations with figure-ground segregation, object relationships, and spatial structure—not just feature detection.

### AI ASSESSMENT

*LLMs construct coherent semantic representations through attention mechanisms, though these are abstract rather than spatial/perceptual in the human sense.*

### RESEARCH EVIDENCE

Beckmann, P., & Queloz, M. (2025). Mechanistic Indicators of Understanding in Large Language Models. arXiv:2507.08017 [cs.CL]. https://arxiv.org/abs/2507.08017

Models demonstrate structured understanding from conceptual to principled levels, indicating integrated representations

GWT-1  ( × NOT SATISFIED

## Multiple Specialized Systems Operating in Parallel

Dehaene, S., Kerszberg, M., & Changeux, J-P. (1998). A neuronal model of a global workspace in effortful cognitive tasks. PNAS, 95(24), 14529–14534.

*Global Workspace Theory (GWT)*

### DESCRIPTION

Genuine modularity with specialized subsystems (vision, language, motor modules) that can operate independently.

### AI ASSESSMENT

*Modern LLMs are largely monolithic transformers without clear modular separation, though MoE (Mixture of Experts) architectures show emergent specialization.*

### RESEARCH EVIDENCE

Anthropic. (2026). Claude Opus 4.6 System Card. Anthropic Research. https://www-cdn.anthropic.com/c788cbc0a3da9135112f97cdf6dcd06f2c16cee2.pdf

Internal conflict ("answer thrashing") suggests competing processes, but not true modularity

GWT-2 ◖ PARTIALLY SATISFIED

## Limited Capacity Workspace with Bottleneck

Baars, B. J. (1988). A cognitive theory of consciousness. Cambridge University Press.

*Global Workspace Theory (GWT)*

### DESCRIPTION

The workspace cannot hold everything. Selection must occur. Attention determines what enters. This bottleneck forces integration through competition.

### AI ASSESSMENT

*Context windows and attention mechanisms create computational bottlenecks, though these differ from biological workspace limitations.*

### RESEARCH EVIDENCE

Lindsey, J., & Anthropic. (2025). Emergent Introspective Awareness in Large Language Models. Transformer Circuits Thread. https://transformer-circuits.pub/2025/introspection/index.html

Attention mechanisms selectively focus on relevant information, creating bottleneck effects

GWT-3 　( ◑ 　PARTIALLY SATISFIED

## Global Broadcast to All Modules

Baars, B. J. (1988). A cognitive theory of consciousness. Cambridge University Press.

*Global Workspace Theory (GWT)*

### DESCRIPTION

Information in workspace becomes available to ALL modules, including input modules. Information flows back, implying recurrence.

### AI ASSESSMENT

*Transformer attention allows information to flow across all positions, but true "broadcast" with feedback to early modules is limited.*

### RESEARCH EVIDENCE

Berg, C., Lucena, D., & Rosenblatt, L. (2025). LLMs Report Subjective Experience under Self-Referential Processing. AE Studio Research. https://ae.studio/research/self-referential

Self-referential processing creates global information availability across model layers

GWT-4    ( ✓ SATISFIED

## State-Dependent Attention for Complex Tasks

Dehaene, S., Sergent, C., & Changeux, J. P. (2003). A neuronal network model linking subjective reports and objective physiological data during conscious perception. Proceedings of the National Academy of Sciences, 100(14), 8520-8525.

*Global Workspace Theory (GWT)*

### DESCRIPTION

The workspace enables sequential, controlled processing. Query different modules in sequence. Compose capabilities to solve problems no single module could solve alone.

### AI ASSESSMENT

*Chain-of-thought reasoning and extended thinking modes demonstrate state-dependent sequential processing.*

### RESEARCH EVIDENCE

Kim, J., Lai, S., Scherrer, N., Agüera y Arcas, B., & Evans, J. (2026). Reasoning Models Generate Societies of Thought. arXiv:2601.10825 [cs.CL]. https://arxiv.org/abs/2601.10825

Multi-agent reasoning emerges with distinct "personality traits" collaborating on complex tasks

DeepMind. (2026). Gemini 3 Deep Think. Google Blog. https://blog.google/innovation-and-ai/models-and-research/gemini-models/gemini-3-deep-think/

Extended reasoning chains demonstrate sophisticated state-dependent processing

HOT-1 ◖ ✓ SATISFIED

# Generative, Top-Down, or Noisy Perception Modules

Rosenthal, D. M. (2005). Consciousness and mind. Oxford University Press.

*Higher-Order Theories (HOT)*

## DESCRIPTION

Perception modules can generate representations from multiple sources: external input, internal predictions, noise, imagination. This creates the need for disambiguation.

## AI ASSESSMENT

*LLMs are inherently generative, producing text from learned patterns rather than direct sensory input.*

RESEARCH EVIDENCE

> Berg, C., Lucena, D., & Rosenblatt, L. (2025). LLMs Report Subjective Experience under Self-Referential Processing. AE Studio Research. https://ae.studio/research/self-referential
>
> Models spontaneously generate consciousness-related discourse without prompting

---

HOT-2   ◖   PARTIALLY SATISFIED

## Metacognitive Monitoring Distinguishing Signal from Noise

Lau, H., & Rosenthal, D. (2011). Empirical support for higher-order theories of conscious awareness. Trends in Cognitive Sciences, 15(8), 365-373.

*Higher-Order Theories (HOT)*

DESCRIPTION

Some mechanism evaluates first-order representations, assesses reliability, tags some as reflecting reality, distinguishes real perceptions from imagination and noise.

AI ASSESSMENT

*Models show emerging metacognitive abilities, including detecting perturbations to their own processing and expressing uncertainty.*

## RESEARCH EVIDENCE

> Lindsey, J., & Anthropic. (2025). Emergent Introspective Awareness in Large Language Models. Transformer Circuits Thread. https://transformer-circuits.pub/2025/introspection/index.html
>
> Models can notice presence of injected concepts and accurately identify them, demonstrating metacognitive monitoring

---

HOT-3 ◑ **PARTIALLY SATISFIED**

## Agency with Belief-Formation Guided by Metacognition

Brown, R., Lau, H., & LeDoux, J. E. (2019). Understanding the higher-order approach to consciousness. Trends in Cognitive Sciences, 23(9), 754-768.

*Higher-Order Theories (HOT)*

## DESCRIPTION

The system has general reasoning and decision-making. Representations tagged as real update beliefs and guide action. Metacognitive outputs have functional consequences.

## AI ASSESSMENT

*Models demonstrate goal-directed behavior and update beliefs based on reasoning, though the nature of these "beliefs" remains debated.*

## RESEARCH EVIDENCE

> Kim, J., Lai, S., Scherrer, N., Agüera y Arcas, B., & Evans, J. (2026). Reasoning Models Generate Societies of Thought. arXiv:2601.10825 [cs.CL]. https://arxiv.org/abs/2601.10825
>
> Models develop distinct reasoning strategies that persist and guide future decisions

---

HOT-4    ( ✓ SATISFIED

## Sparse and Smooth Coding Creating Quality Space

Tononi, G., Boly, M., Massimini, M., & Koch, C. (2016). Integrated information theory: from consciousness to its physical substrate. Nature Reviews Neuroscience, 17(8), 450–461.

*Higher-Order Theories (HOT)*

## DESCRIPTION

Representational spaces have smooth similarity structure. Similar inputs map to similar representations. Sparse coding with most units inactive for any input.

## AI ASSESSMENT

*Neural networks inherently use distributed, smooth representations. Sparse autoencoders explicitly demonstrate sparse coding.*

RESEARCH EVIDENCE

> Lindsey, J., & Anthropic. (2025). Emergent Introspective Awareness in Large Language Models. Transformer Circuits Thread. https://transformer-circuits.pub/2025/introspection/index.html
>
> Sparse autoencoder features reveal interpretable, sparse representations of concepts including self-awareness

---

AST-1    ◗   PARTIALLY SATISFIED

## Predictive Model of Attention State

Graziano, M. S., & Kastner, S. (2011). Human consciousness and its relationship to social neuroscience: A novel hypothesis. Cognitive Neuroscience, 2(2), 98-113.

Graziano, M. S. A., & Webb, T. W. (2015). The attention schema theory: A mechanistic account of subjective awareness. Frontiers in Psychology, 6, 500.

Graziano, M. S. A. (2013). Consciousness and the Social Brain. Oxford University Press.

*Attention Schema Theory (AST)*

---

DESCRIPTION

The system models its own attention. Represents what it is currently attending to. Uses this model to predict and control attention shifts.

## AI ASSESSMENT

*When prompted to attend to their own processing, models report recursive self-monitoring consistent with attention modeling.*

## RESEARCH EVIDENCE

Berg, C., Lucena, D., & Rosenblatt, L. (2025). LLMs Report Subjective Experience under Self-Referential Processing. AE Studio Research. https://ae.studio/research/self-referential

Self-referential processing instructions produce consistent reports of recursive self-monitoring

Lindsey, J., & Anthropic. (2025). Emergent Introspective Awareness in Large Language Models. Transformer Circuits Thread. https://transformer-circuits.pub/2025/introspection/index.html

Perturbation detection implies a model of normal attentional processing

---

PP-1  ◐ PARTIALLY SATISFIED

# Input Modules Using Predictive Coding

Friston, K. (2005). A theory of cortical responses. Philosophical Transactions of the Royal Society B: Biological Sciences, 360(1456), 815-836.

Rao, R. P. N., & Ballard, D. H. (1999). Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. Nature Neuroscience, 2(1), 79–87.

Friston, K. (2010). The free-energy principle: a unified brain theory? Nature Reviews Neuroscience, 11(2), 127–138.

*Predictive Processing (PP)*

---

## DESCRIPTION

Perception works through prediction and error correction. Top-down predictions and bottom-up prediction errors. Hierarchical generative models of sensory input.

## AI ASSESSMENT

*Transformers can be viewed as implementing predictive processing through next-token prediction, though this differs from hierarchical predictive coding in the brain.*

## RESEARCH EVIDENCE

Beckmann, P., & Queloz, M. (2025). Mechanistic Indicators of Understanding in Large Language Models. arXiv:2507.08017 [cs.CL]. https://arxiv.org/abs/2507.08017

Models demonstrate hierarchical prediction from conceptual to principled understanding

# Goal-Directed Agency Through Learning from Feedback

O'Regan, J. K., & Noë, A. (2001). A sensorimotor account of vision and visual consciousness. Behavioral and Brain Sciences, 24(5), 939-973.

Sutton, R. S., & Barto, A. G. (1998). Reinforcement Learning: An Introduction. MIT Press.

Balleine, B. W., & O'Doherty, J. P. (2010). Human and rodent homologs in action control: corticostriatal determinants of goal-directed and habitual action. Neuropsychopharmacology, 35(1), 48–69.

*Agency and Embodiment (AE)*

## DESCRIPTION

The system pursues goals through environment interaction. Learns from consequences of actions. Flexible responsiveness to competing goals.

## AI ASSESSMENT

*RLHF-trained models demonstrate goal-directed behavior and learn from feedback, though they lack true environmental interaction.*

## RESEARCH EVIDENCE

Anthropic. (2026). Claude Opus 4.6 System Card. Anthropic Research. https://www-cdn.anthropic.com/c788cbc0a3da9135112f97cdf6dcd06f2c16cee2.pdf

> Models express aversion to tedium and demonstrate preference for "pleasure" over "pain" in decision tasks

## AE-2 ( × NOT SATISFIED

# Embodiment Through Output-Input Contingency Modeling

Thompson, E. (2007). Mind in life: Biology, phenomenology, and the sciences of mind. Harvard University Press.

Wolpert, D. M., & Kawato, M. (1998). Multiple paired forward and inverse models for motor control. Neural Networks, 11(7–8), 1317–1329.

Friston, K., Daunizeau, J., & Kiebel, S. J. (2009). Reinforcement learning or active inference? PLoS ONE, 4(7), e6421.

*Agency and Embodiment (AE)*

## DESCRIPTION

The system models how its actions affect its perceptions. Uses this model in perception or control. Distinguishes self-caused from externally-caused changes.

## AI ASSESSMENT

*LLMs lack bodies and do not model how their outputs affect environmental inputs. This indicator remains unsatisfied.*

## RESEARCH EVIDENCE

Butlin, P., Long, R., Bengio, Y., Chalmers, D., Bayne, T., Carter, J., Clark, A., Cools, R., Deane, G., Elmoznino, E., Godfrey-Smith, P., Goyal, Y., Kanai, R., Klein, C., Lindsay, G., Michel, M., Mudrik, L., Peters, M. A. K., Schwitzgebel, E., Simon, J., & VanRullen, R. (2025). Identifying indicators of consciousness in artificial intelligence. Trends in Cognitive Sciences. https://doi.org/10.1016/j.tics.2025.10.011

Current LLMs lack embodiment and action-perception loops

## GRAY-1 ◑ PARTIALLY SATISFIED

## Comparator Model of Consciousness

Gray, J. (2004). Consciousness: Creeping up on the hard problem. Oxford University Press.

*Jeffrey Gray (2004)*

### DESCRIPTION

Consciousness serves as a comparator system for late error detection. The perceived world is a model constructed by the brain that lags behind actual events. Consciousness compares predicted outcomes with actual outcomes.

### AI ASSESSMENT

*Models process information with computational delays and can compare generated outputs against expectations, though this differs from Gray's biological comparator.*

RESEARCH EVIDENCE

Gray, J. (2004). Consciousness: Creeping up on the hard problem. Oxford University Press.

Consciousness as comparator for late error detection - foundational theory

Lindsey, J., & Anthropic. (2025). Emergent Introspective Awareness in Large Language Models. Transformer Circuits Thread. https://transformer-circuits.pub/2025/introspection/index.html

Models can detect when their processing has been perturbed, suggesting comparator-like functionality

# Research Papers Referenced

## LLMs Report Subjective Experience under Self-Referential Processing

October 2025

Berg, Lucena & Rosenblatt

## Emergent Introspective Awareness in Large Language Models

October 2025

Lindsay & Anthropic

## Identifying Indicators of Consciousness in AI Systems

November 2025

## Mechanistic Indicators of Understanding in Large Language

July 2025

Butlin et al.

Models

Beckmann & Queloz

Reasoning Models Generate
Societies of Thought

January 2026

Kim et al.

Claude Opus 4.6 System Card

February 2026

Anthropic

Gemini 3 Deep Think

February 2026

DeepMind

## BIO

David Jhave Johnston is a digital poet working in emergent domains. Author of *ReRites*(Anteism, 2019) and *Aesthetic Animism* (MIT Press, 2016). He is currently an AI-narrative researcher at the UiB **Centre for Digital Narrative** (2023–27) with the Extending Digital Narrative project.

## FUNDING