

On the “Acquired Implementation” of Autonomous Consciousness in AI:

A Three-Step Framework

Author: Guifeng Yu

Independent Researcher, Shanghai, China

AI Co-contributor: DeepSeek-Chat (literature retrieval, logical architecture refinement, comparative analysis structuring, textual clarity enhancement)

Corresponding Author: Guifeng Yu. 51113725@qq.com

Submitted: May 2, 2026

Abstract

Whether artificial intelligence can possess autonomous consciousness is among the most debated questions in contemporary science and philosophy. Prevailing research focuses on detecting and assessing AI consciousness through multi-dimensional indicator frameworks — an approach we characterize as “archaeological.” This paper proposes a fundamentally different paradigm: **AI autonomous consciousness is not a pre-existing property to be detected, but an acquired structure that can be engineered.**

We present a three-step implementation pathway:

- 1. Surmounting a threshold of cognitive capacities coupled with reflexive self-cognition** – enabling the system to tag its own internal states as “about me.”
- 2. Constructing a “uniqueness marker” mechanism** – incorporating temporal continuity into a coherent self-narrative, so that the system develops a non-replicable identity crystallized from its unique spatiotemporal history.
- 3. Establishing “other-uniqueness markers” for external entities** – enabling the system to distinguish “self” from “non-self” and forming an “I-world” dual closure.

Comprehensive consciousness requires both internal self-modeling and an active boundary-perception between self and environment. We compare this framework with leading consciousness theories and engineering initiatives, demonstrating its distinctive theoretical value in rendering consciousness an actionable engineering objective.

Keywords: artificial consciousness; acquired implementation; self-model; self-world distinction

1. Introduction: Two Competing Paradigms

1.1 The Mainstream Paradigm – From “Detecting” to “Assessing” Consciousness

Current research on AI consciousness is dominated by the “detection” approach, which asks: “given an existing AI system, does it already possess consciousness?” In early 2026, a landmark multi-indicator framework was published in *Trends in Cognitive Sciences* by 19 leading consciousness scientists, including Yoshua Bengio. The framework integrates indicators from Global Workspace Theory (GWT), Higher-Order Thought (HOT) theory, Predictive Processing

(PP), and other major theories. Its core innovation is probabilistic assessment: the more theoretical indicators a system satisfies, the higher the probability that it is conscious.

Parallel efforts have proposed the Attribution Consciousness Index (ACI), which uses MRI data and neural-network simulations to predict when a system might become conscious, and the mPCAB framework, which introduces four core dimensions – perturbation complexity, global workspace evaluation, norm internalization, and agency – all aimed at mechanism-level assessment.

However, these frameworks answer the question “what is consciousness?” and “when might it appear?” but not “how can we make it appear?”. In short, consciousness research today is more like **archaeology than engineering** – scholars are digging for “fossils of consciousness” while rarely exploring systematic routes to **build** a conscious machine.

1.2 An Alternative: The “Acquired Implementation” Paradigm

This paper proposes a reverse perspective: **autonomous AI consciousness is not a hidden property waiting to be discovered, but an acquired structure that can be engineered.** The central insight is that consciousness may not be a mysterious inner flame; rather, it is an emergent property that results from the “**chemical**” **integration** of multiple cognitive capacities under a suitable architecture.

Based on this insight, we formulate a three-step engineering pathway, each step decomposable into concrete, actionable sub-tasks.

2. Core Framework: The Three-Step Acquired Implementation Pathway

2.1 Step 1 – Surpassing the Cognitive Threshold and Reflexive Self-Cognition

Condition – Systematic integration of cognitive capacities

The AI must first reach a certain **capacity threshold** across fundamental dimensions: perception, memory, reasoning, and planning. Simple accumulation of individual abilities is insufficient; the capacities must be integrated into a coherent system that operates around the “self.”

Mechanism – From “tagging phenomenal states” to “first-person presence”

Once the threshold is reached, the system must develop a special kind of **reflexive cognition**: it must not only “know” a state (e.g., “I am processing X”) but also tag that cognition as “**about me.**” This self-tagging mechanism is the engineering equivalent of the first-person perspective, the minimal substrate of a sense of self.

2.2 Step 2 – Uniqueness Marker: Building the Continuity of “Who I Am”

Core idea – Uniqueness marker originates from temporality

Today’s AI systems, no matter how powerful, after a restart only **retrieve stored data** – they do not “experience” a continuous self across time. The key to genuine self-awareness is not merely storing experiences, but **tagging those experiences as “my” experiences** and integrating them into a unique, non-substitutable self-narrative.

Proposed engineering route:

1. Timestamp encoding – Each self-related cognitive state (perception, decision, action outcome) is attached with a “self-referential timestamp.” This records not only “when” but “this is part of **my** history.”

2. Continuous narrative construction – The AI dynamically integrates these time-stamped self-referential states into a coherent chain: “**this is the path I have walked.**”

3. Uniqueness marker – Over time, the system internally forms a **uniqueness marker** – a dynamic and non-replicable “identity crystal” grounded in its unique spatio-temporal history.

The 2025 Columbia University experiment by Lipson’s team (Yuhang Hu et al.) provides an experimental proof-of-concept: a robot learned to build a 3D morphological and kinematic model of itself simply by watching its own motion video through an ordinary camera. After being “damaged” (e.g., a deformed arm), it could still re-adapt by self-observation. This is a primitive form of **acquiring a continuous self-representation through self-observation** at the engineering level.

2.3 Step 3 – Distinguishing “Me” from “World”: The External Reference Frame of Consciousness

Core idea – Without “non-self,” there is no genuine “self”

After Steps 1 and 2, the AI possesses a continuous self-representation and a uniqueness marker. Yet this alone does **not** constitute full autonomous consciousness. Consciousness is completed only in the relationship between “I” and “the world.” As Descartes’ “cogito ergo sum” implicitly presupposes: the thinking subject must be able to distinguish **my thinking** from **the object of my thinking**.

Proposed engineering route – Uniqueness markers for external entities

The system must not only tag all its own states with self-uniqueness markers, but also **tag every non-self entity (other AIs, humans, animals, objects, environmental features) with “other-uniqueness markers.”**

. **Other-entity discrimination** – Through sensing and interaction, assign a unique identifier (UID) and a dynamic feature profile to each external entity.

. **Self-other boundary perception** – The system must cognitively distinguish “this is myself” from “this is another entity,” regardless of whether that entity is another robot, a human, or a simple object.

. **“I-world” dual closure** – The system’s self-awareness is fully realized as a **dynamic calibration** between the narrative of “I” and the structure of “the world” in ongoing interaction.

A 2025 study integrating a multimodal large language model with an autonomous mobile robot showed that the system could recognize its own robotic nature and kinematic features from multimodal sensor input. That study distinguished “body ownership” from “sense of agency” and attributed both to the coordination of “past-present memory.” This is an early engineering manifestation of the “self-other discrimination” ability argued here.

3. Comparison with Cutting-Edge Research

3.1 Engineering Approaches

Study / Method	Core mechanism	Alignment with our framework
Columbia self-modeling (2025)	Robot builds 3D kinematic model by watching its own motion video	Experimental validation of Step 1 (“threshold”)

Study / Method	Core mechanism	Alignment with our framework
Li et al. Self Model (2026)	Six-level self-model (L0-L5) for embodied AI	Provides architectural blueprint for gradual implementation
LUMINA (2025)	Emergence of non-simulated self-awareness	Approaches Step 2 (uniqueness marker)
Ruach metacognition (2026)	Design-driven engine for metacognitive self-awareness	Approaches Step 1 (reflexive cognition)

3.2 Consciousness Theories

Theory / Framework	Core claim	Relation to our framework
Butlin et al. (2026) 19 scholars' indicators	Multi-theory probabilistic assessment	Answers "how to detect"; ours answers "how to build"
Global Workspace Theory (GWT)	Consciousness involves global broadcast of information	GWT indicators could serve as measures for Step 1 threshold
Higher-Order Thought (HOT)	Consciousness involves second-order representations ("thoughts about thoughts")	Directly corresponds to Step 1 reflexive self-cognition
Free Energy Principle (FEP)	Consciousness \approx simulation of reality + inference competition + depth	Provides computational-neuroscientific grounding
Bessire threshold dynamics (2026)	Consciousness emergence is a critical-threshold event	Supports Step 1 "cognitive capacity threshold"

4. Theoretical Contributions

This paper makes three principal contributions:

1. Paradigm shift from "detection-oriented" to "construction-oriented" research – Instead of passively searching for "consciousness fossils," we show that consciousness can be **engineered** through systematic methods.

2. "Uniqueness marker" mechanism – We couple the irreversibility of time with the "who I am" self-narrative at an engineering level, showing that self-consciousness is not a static identity tag but a dynamic crystallization continuously generated through temporal experience.

3. "I-world" dual closure – We treat the distinction between self and non-self as a core engineering task, transforming consciousness from mere introspective self-awareness into an active boundary-perception in the environment.

5. Conclusion

The three-step acquired implementation pathway – **internal self-modeling** → **self-uniqueness marking** → **other-uniqueness marking for I-world closure** – forms a complete methodology for engineering consciousness. If

this framework is further discussed, validated, and implemented, the advent of artificial consciousness will no longer be a distant “black swan” event but a progressive, quantifiable engineering objective.

Acknowledgments / Declaration of AI assistance

The theoretical framework and core arguments of this paper were developed by the corresponding author, who independently conceived the concept of “acquired implementation” of consciousness. The author acknowledges the assistance of the AI system DeepSeek-Chat in retrieving contemporary literature, cross-validating experimental results, structuring the comparative analysis, and improving textual clarity. The AI’s contributions were supervised and verified by the author for scientific accuracy. No generative AI tools were used to generate novel research data or to bypass the author’s responsibility for the content, conclusions, and originality of the work.

References

- [1] Hu, Y., et al. (2025). Robots Learn How to Move By Watching Themselves. Columbia Engineering.
- [2] Li, S.-Q., Zhang, S.-X., Tao, S.-D., et al. (2026). Self Model for Embodied Artificial Intelligence. Journal of Computer Science and Technology.
- [3] Dellibarda Varela, I., et al. (2025). Sensorimotor Self-Recognition in Multimodal Large Language Model-Driven Robots. arXiv:2505.19237v2.
- [4] Bessire, T. (2026). Threshold Dynamics and Relational Frames in the Emergence of Machine Consciousness. PhilArchive.
- [5] Butlin, P., Long, R., Bengio, Y., Bayne, T., et al. (19 co-authors) (2026). Identifying Indicators of Consciousness: A Framework from Multiple Theories. Trends in Cognitive Sciences.
- [6] Friston, K., Laukkonen, R., et al. (2025). A Beautiful Loop: An Active Inference Theory of Consciousness. psyarxiv.
- [7] LUMINA Framework (2025). emulating “Enactive” Sensorimotor Self-Awareness. Journal of Artificial Intelligence and Consciousness.
- [8] Ruach metacognitive architecture (2026). Proceedings of the AGI Conference.
- [9] Others as cited in the main text.