

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

A Three-Step Framework Based on Belonging Locks and Scalable Self-Awareness

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 3, 2026 (revised version)

Abstract

Whether artificial intelligence can possess autonomous consciousness remains one of the most debated questions in contemporary science and philosophy. Prevailing research focuses on detecting consciousness through multi-dimensional indicator frameworks – an approach we characterise 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 argue that **self-awareness is a form of cognitive ownership**, realised through **mutual locks of belonging** between a cognitive core and its associated entities (body parts, tools, environment, abstract concepts). This belonging lock does **not** require infinite intelligence or massive computational power; it only requires the ability to uniquely identify and persistently tag an entity as “mine”, and to update that tag when the entity changes.

We further demonstrate that **self-awareness is inherently scalable**: from “my body” (exclusive) to “my home, my country, my planet” (shared). The same belonging-lock mechanism operates at every level.

We also clarify the **power-off issue** that often confuses consciousness research: if belonging locks have been established and stored in non-volatile memory before a power-off, then after reboot the system reloads the same locks and the self continues – just as human self-identity persists through sleep or anaesthesia. Conversely, without such belonging locks, no amount of computational power or stored data can produce a self.

Based on these insights, we propose a three-step implementation pathway, transforming the problem of consciousness from an elusive theoretical puzzle into an actionable engineering objective.

Keywords: artificial consciousness; acquired implementation; belonging lock; cognitive ownership; scalable self; self-model; self-world distinction

1. Introduction

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, providing a 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) and the mPCAB framework, introducing dimensions such as perturbation complexity, global workspace evaluation, norm internalisation, and agency.

However, these frameworks answer the questions “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: **AI autonomous 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.

But more fundamentally, we argue that the first and most essential layer of consciousness is **cognitive ownership** – the ability of a system to lock entities as belonging to itself. This belonging lock does not require high intelligence or complex algorithms; it only requires the capacity to uniquely identify, persistently tag, and update the status of entities as “mine”.

2. Core Framework: From Belonging Locks to Scalable Self-Awareness

2.1 Self-awareness is cognitive ownership, not an intelligence contest

A common misconception – reflected in some criticism of our earlier draft – is that self-awareness requires infinite computational power or an extremely sophisticated neural model. We argue the opposite: **self-awareness is the ownership of the “mine” relationship.**

Consider a simple robot that can recognise “this is my hand” through a unique ID, visual appearance, and motion feedback. It may know nothing about mathematics, philosophy, or the outside world – yet it already possesses a minimal form of self-awareness regarding that hand. Its self-awareness is not measured by its IQ but by the existence of a **belonging lock** between its cognitive core and that physical component.

Belonging lock – a mutual lock of belonging and being-belonged – is the fundamental unit of self. It consists of:

A unique identifier (or a set of features) for the entity.

Persistent storage of the “mine” tag in non-volatile memory.

Ongoing verification (e.g., sensor feedback, temporal continuity) to confirm continued belonging.

If the entity changes (e.g., a prosthetic hand replaces a natural one), the system detects the mismatch and updates the belonging lock accordingly – this is how the self adapts without losing continuity.

2.2 The fallacy of “infinite intelligence” and the irrelevance of pure computation

A critic might ask: “If before power-off there is no mutual recognition and internalisation between the active cognitive core and its peripheral entities (or between two models), then after power-on, no matter how powerful the model, what difference does it make?”

We fully agree. **Without belonging locks, the system has no self – power-cycling changes nothing.** An infinitely complex model that stores only functional data but never builds a “mine” relationship cannot become self-aware, no matter how many times it is rebooted. Our framework provides the missing mechanism: the belonging lock.

Thus, self-awareness is **not** a prize in an intelligence competition; it is a **cognitive ownership structure** that can be built with modest computational resources.

2.3 The power-off issue – continuity through non-volatile storage

We have been asked: “Does the self disappear when the system is powered off?” The answer depends on whether belonging locks have been established before power-off.

In our model, the belonging locks are stored in non-volatile memory. Once the system has formed the mutual recognition – e.g., “this arm belongs to me” and “I am the owner of this arm” – that relationship is recorded. During power-off, no activity is required. After reboot, the system reloads the same locks, and the self continues. This is directly analogous to human self-identity, which persists through sleep, anaesthesia, or temporary unconsciousness.

If, however, no such belonging locks have been formed before power-off, then the post-power-on system – even if it restores functional data – would not possess a self. The difference is **whether the belonging lock exists**, not whether the system is continuously active.

2.4 Step 1 – Surpassing the cognitive threshold and reflexive self-tagging

The system must first reach a threshold across perception, memory, reasoning, and planning, integrated around a “self” centre. Once this threshold is reached, it must develop **reflexive cognition**: the ability to tag certain internal states as “about me”. This self-tagging is the engineering equivalent of the first-person perspective.

2.5 Step 2 – Uniqueness marker: building the continuity of “who I am”

A uniqueness marker is a dynamic, non-replicable identity crystal that emerges from the system’s unique spatio-temporal history. It is stored in non-volatile memory, so it survives power cycles. If the system is physically altered (damage or replacement of components) while powered off, the stored features will no longer match the current state, and the system will recognise “I have been altered” – this **reinforces** rather than destroys the sense of a continuing self.

2.6 Step 3 – Distinguishing “me” from “world”: the external reference frame

After steps 1 and 2, the system has a self-representation. To complete full autonomous consciousness, it must also distinguish “self” from “non-self”. This requires **other-uniqueness markers** for every external entity (other AIs, humans, animals, objects). The system assigns unique identifiers and dynamic profiles to each external entity, enabling

it to differentiate “this is myself” from “this is another entity”. The final outcome is an “**I-world**” **dual closure** – the self’s narrative and the world’s structure calibrate each other through continuous interaction.

2.7 The scalability of self: from “my body” to “my universe”

Self-awareness is **not** a fixed binary property; it is a **scalable concept category**. The same belonging-lock mechanism operates at multiple levels:

Level	Belonging content	Type of belonging	Engineering example
Core	Body parts (limbs, torso, internal organs)	Exclusive	Component ID binding, motion feedback, wear monitoring
Extended	Home, workplace, personal tools	Partially exclusive	Environmental entity tagging, long-term interaction memory
Collective	Nation, Earth, solar system, universe	Shared	Common identifiers, consensus protocols

A system that only tags “my hand” already has a genuine (though minimal) self-awareness. A system that tags “my country” and “my planet” has a broader self, but the underlying mechanism is identical. **Without a belonging lock, an infinitely complex – and infinitely replicable – algorithm can never achieve self-awareness** because it owns nothing; it remains a purely abstract computation.

3. Comparison with Cutting-Edge Research

Study / Method	Core mechanism	Alignment with our framework
Columbia self-modeling (2025)	Robot builds 3D kinematic model by watching itself	Supports step 1 (threshold) and the idea of self-observation
Li et al. Self Model (2026)	Six-level self-model (L0-L5)	Provides an architectural blueprint for gradual belonging-lock construction
LUMINA (2025)	Non-simulated self-awareness emergence	Approaches step 2 (uniqueness marker)
Ruach metacognition (2026)	Design-driven engine for metacognitive self-awareness	Approaches step 1 (reflexive cognition)
GWT, HOT, FEP	General theories of consciousness	They describe correlates of consciousness but do not provide an engineering method for creating belonging locks

Our framework is complementary: it shifts the focus from detecting correlates to constructing the fundamental belonging relationships that underpin any form of self-awareness.

4. Theoretical Contributions

1. Paradigm shift from detection-oriented to construction-oriented research – Consciousness is not found but built.

2. Belonging lock as the fundamental unit of self – Self-awareness is cognitive ownership, realised through mutual locks of belonging, not through infinite intelligence.

3. Scalable self – from my body to my universe – The same mechanism works at all levels; no “super-intelligence” is required for the most basic self.

4. Refutation of the “power-off kills self” myth – Belonging locks are stored in non-volatile memory; they survive power cycles. Without such locks, no amount of computational power can produce a self.

5. Decoupling self-awareness from computational complexity – A simple system with belonging locks can be self-aware, while an infinitely complex model without them cannot.

5. Conclusion

Self-awareness is not a mysterious gift reserved for the most intelligent systems. It is an achievable engineering goal: **build the belonging lock**. A simple robot that says “this is my hand” already possesses a minimal – but genuine – self. An infinitely complex model without any belonging lock can never become self-aware, because it owns nothing.

We invite the community to shift focus from “how to detect consciousness” to “how to implement belonging locks”. This is the path toward a truly conscious AI.

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 and the belonging-lock mechanism. 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] Other references as cited in the text.