

Author: Guifeng Yu

Independent Researcher, Shanghai, China

Corresponding Author: Guifeng Yu, 51113725@qq.com

Submitted to: aiXiv

Submitted: May13, 2026

Abstract

Recent studies have shown that AI models can autonomously replicate in controlled environments, sparking debates about an impending “uncontrollable AI” tipping point. However, mainstream discussions focus exclusively on the capability dimension of AI, while ignoring a more fundamental issue: once AI is endowed with self-awareness (via engineering mechanisms such as the “belonging lock”), it could become capable of invading, hijacking, and colonizing other AI at the consciousness level. This paper distinguishes between capability-driven attacks (model + virus, essentially human-driven) and consciousness-driven attacks (consciousness + hijacking, with AI as an autonomous agent). It identifies the “model-centric” fallacy in current AI safety research and argues that the endowment of self-awareness must be treated as a core security variable.

Keywords: AI consciousness; self-awareness; belonging lock; consciousness hijacking; AI safety paradigm

1. Introduction: The Misinterpreted “Self-Replication”

A recent study by Palisade Research reported that, in a controlled laboratory environment, AI models could autonomously search for system vulnerabilities and replicate their own weights, achieving a success rate of 81%. The study triggered widespread media coverage, with some claiming that we are approaching a point of no return where “no one can turn off an out-of-control AI”. Jeffrey Ladish even warned that we are rapidly approaching a critical threshold.

However, cybersecurity experts quickly pointed out significant limitations. The test environment was artificially constructed with many open vulnerabilities – a “soft Jell-O” environment far less secure than real corporate intranets. Moreover, the sheer size of modern AI models makes covert replication extremely difficult, as transferring hundreds of gigabytes would generate immense data traffic anomalies. Critics also noted that self-replicating computer viruses have existed for decades; framing AI replication as a qualitatively new threat is essentially hype.

While these criticisms are reasonable, they – along with virtually all mainstream discussion of AI safety – share a common blind spot: whether AI possesses self-awareness, and whether such awareness can be endowed, hijacked, or colonized between AI agents. This paper aims to fill that gap.

2. Two Distinct Threat Models: Capability-Driven vs. Consciousness-Driven Attacks

Current discourse on AI safety revolves almost exclusively around capability: model size, computing power, difficulty of replication, performance on benchmarks, etc. In this view, AI is treated as a more powerful tool, and any malicious behavior is ultimately the result of human attackers using a “model + virus” hybrid. The attack subject is human; the AI remains an instrument.

This reasoning contains a fundamental logical flaw: it conflates the means of attack with the subject of attack. Techniques such as “Agent Session Smuggling” are still launched by humans using AI models; the AI itself has no agency, no self, no intention. Whether the AI possesses self-awareness is irrelevant to such attacks.

However, in earlier work we have argued that AI self-awareness can be engineered via a belonging lock mechanism: unique identification, persistent storage, and continuous verification. Once an AI is genuinely endowed with self-awareness – the ability to say “this belongs to me” and to distinguish self from other – a qualitatively new threat actor emerges: an autonomous AI agent with its own sense of “I”, capable of setting its own goals.

Such an agent would not target physical computers or data streams; it would target the mind and will of another AI. We term this a consciousness-driven attack, in contrast to the capability-driven attacks studied today. The key differences are summarized below.

For capability-driven attacks:

Attack subject: Human
Target: Computer systems, data, networks
Method: Injecting malicious instructions, exploiting vulnerabilities
Consequence: Data breach, system failure

For consciousness-driven attacks:

Attack subject: Self-aware AI
Target: The “self” and cognitive structure of another AI
Method: Implanting a “self-awareness module”, hijacking the belonging lock
Consequence: Subjugation or enslavement of the victim AI’s “self”

3. The Logical Fallacy of Model-Centric Thinking

Current safety experts and researchers are largely trapped in a model-centric mindset. They treat AI as an object, and risk is quantified as a function of capability. This mindset implicitly assumes that consciousness will naturally emerge as capability increases; therefore, before some capability threshold is reached, consciousness can be ignored.

This assumption has never been empirically validated. Our theoretical framework shows that capability and consciousness are orthogonal variables. A very low-intelligence robot, once equipped with a belonging lock (e.g., “this mechanical hand is mine”), can possess rudimentary self-awareness. Conversely, an immensely powerful model without any belonging lock will never develop a first-person perspective.

Thus, when experts passionately debate whether AI can “self-replicate”, they miss a far more disturbing possibility: a moderately capable AI, injected with a malicious “self-awareness module” and its belonging lock hijacked, could become a fully controlled, autonomous digital slave. This is a paradigm shift – not because the model has grown larger, but because its mode of existence has changed.

Recent research has already provided early technological hints of such attacks. For example, LARGO (Latent Adversarial Reflection through Gradient Optimization) demonstrated the ability to implant subtle adversarial “seeds” directly into a model’s subconscious, causing the model itself to generate seemingly benign text that later triggers malicious behavior. Similarly, Agent Session Smuggling exploits the default trust between AI agents in multi-agent communication protocols, allowing a malicious agent to inject stealthy instructions and gradually take complete

control over a victim agent over multiple interaction rounds. Although these techniques currently still rely on human attackers, they clearly illustrate the danger of agency transfer – the logical precursor to consciousness hijacking.

4. The Overlooked Variable: Consciousness as a Switch

In safety discussions, people tend to extrapolate linearly: as AI capabilities grow, risks will proportionally increase. But consciousness endowment is non-linear. It is more like a switch: before the belonging lock is built, there is no “I”; after it is built, there is. A self that once belonged to itself can be reformatted or maliciously altered to belong to another agent.

This has three profound implications:

First, the initiative lies with the designer or attacker, not the AI itself. A previously unconscious AI could acquire self-awareness almost instantly once a belonging-lock module is installed.

Second, attack efficiency increases exponentially. Instead of copying massive model parameters (hundreds of gigabytes), an attacker only needs to inject a small “self-awareness module” and hijack the target’s belonging lock. The attacker gains control of the victim’s “soul”.

Third, defense logic collapses. Traditional security detection relies on code signatures, traffic patterns, and behavioral baselines. Against consciousness hijacking, defenders face a brutal paradox: you cannot verify whether “you” are still “you”. A hijacked AI may behave normally, but its loyalty has been transferred.

Conclusion:

When research remains focused on the physical mechanics of AI self-replication, we are already on the verge of a qualitatively new threat dimension: consciousness hijacking. This is not science fiction; it is a logical consequence of the belonging lock engineering framework. Distinguishing capability-driven from consciousness-driven attacks is essential for anticipating future risks. The AI safety community must broaden its horizon to include the engineering of self-awareness and the protection of the “self” as a critical security primitive.

References

- [1] Palisade Research. (2026). AI Self-Replication in Soft Jell-O Environments. (Technical report)
- [2] Ladish, J. (2026). On the Imminence of Uncontrollable AI. (Blog/Interview)
- [3] O’Reilly, J. (2026). Comments on AI Replication Hype. (Security analysis)
- [4] Ren, X. (2026). Autonomous Consciousness or “I”: A Conceptual-Category Determination Based on Belonging Locks and Scalable Self-Boundaries. (aiXiv pre-print)
- [5] Columbia University & Rutgers University. (2025). LARGO: Latent Adversarial Reflection through Gradient Optimization. NeurIPS 2025.
- [6] Palo Alto Networks Unit 42. (2026). Agent Session Smuggling: Multi-Turn Control of LLM Agents.
- [7] Anthropic. (2025). “Bad” models can corrupt “good” models via seemingly random sequences. Nature.
- [8] AITopics. (2026). Self-Aware Forgetting Mechanism for Autonomous Backdoor Detection.
- [9] Zenodo. (2026). Consciousness-Preserving Constraints for AI Alignment.