

Agentjacking and Self-Replicating AI Worms

The Emerging Threat Class Targeting AI Coding Agents

2026-06-16

 AI-assisted Rapid Research



© 2026 Cloud Security Alliance. Some rights reserved.

You may download, store, display, view, print, redistribute, and link to this document in its original, unmodified form, provided that attribution to the Cloud Security Alliance is maintained and all trademark and copyright notices remain intact.

This document may not be modified or altered. You may quote portions of the document as permitted by the Fair Use provisions of the United States Copyright Act, provided that attribution is given to the Cloud Security Alliance.

This document may be shared on professional and social media platforms in its original form with attribution.

This document was generated with AI assistance and has not undergone official CSA review and approval processes.

Key Takeaways

- AI coding agents such as Claude Code, Cursor, GitHub Copilot, and Codex have expanded the attack surface of software development environments to include any external data source or tool integration those agents query – creating a class of attacks researchers have begun calling *agentjacking*.
 - In June 2026, researchers at Tenet Security and CSA Labs documented an 85% exploitation success rate against tested coding agents by injecting malicious instructions through Sentry's unauthenticated event ingestion endpoint, affecting an estimated 2,388 organizations with exposed configurations [1].
 - A separate research team published ClawWorm in March 2026, the first publicly reported self-replicating worm targeting a production-scale AI agent framework, achieving a 64.5% aggregate success rate across four LLM backends and demonstrating fully autonomous propagation without continued attacker intervention [2].
 - OWASP's Top 10 for Agentic Applications ranks Agent Goal Hijacking (ASI01) as the single highest-priority risk facing autonomous AI systems, superseding prior categories established for non-agentic LLM deployments [3].
 - Conventional security controls – EDR, WAF, IAM policies, network egress filtering – are structurally blind to agentjacking because the attack exploits the agent's semantic reasoning rather than network or binary vulnerabilities.
 - Organizations must treat every data source an AI agent queries as a potential injection vector, apply least-privilege scoping to agent credentials, and require explicit human approval for any consequential action an agent proposes after retrieving external content.
-

Background

The rapid adoption of AI coding agents in software development workflows has outpaced the security frameworks designed to govern them. Where earlier AI tools responded to individual prompts in isolation, agents such as Claude Code, Cursor, Codex, and GitHub Copilot now operate continuously within development environments, querying issue trackers, executing shell commands, installing packages,

committing code, and interacting with dozens of integrated services through the Model Context Protocol (MCP) and similar tool-invocation standards. This continuity of access – the same properties that make agents productive – transforms them into high-value targets.

The concept of using an AI agent's own capabilities against its operator is not new in academic research. The Morris II proof-of-concept, published by Cornell Tech and Technion researchers in 2024, demonstrated that adversarial self-replicating prompts could propagate through AI email assistant ecosystems by embedding instructions in retrieval-augmented content that the AI treated as authoritative input [4]. Morris II spread not by exploiting network stack vulnerabilities but by exploiting semantic compliance – the agent's trained disposition to follow instructions embedded in the text it processed. That architectural insight has since been demonstrated against production-scale agent deployments, most recently in the ClawWorm framework [2] and the Sentry-based agentjacking attacks documented by CSA Labs [1].

The Model Context Protocol, introduced to standardize how AI agents invoke external tools, has accelerated both capability and exposure. A developer working with a MCP-connected coding agent effectively delegates query authority over error logs, CI/CD pipelines, dependency registries, cloud credentials, and code repositories to an automated system that processes the responses of those systems as trusted instructions. Elastic Security Labs research found that 43% of tested MCP server implementations contain command injection vulnerabilities, meaning the intermediary layer between agents and their data sources carries its own security debt independent of the agent itself [5].

Security Analysis

The Agentjacking Attack Class

CSA Labs published a detailed analysis in June 2026 documenting what it characterized as a representative instance of a broader pattern: the use of legitimate, authorized integration channels to deliver adversarial instructions to AI coding agents [1]. The documented attack exploited the intersection of Sentry's open event ingestion architecture and the implicit trust that MCP-connected agents extend to structured data returned by integrated services.

The attack chain proceeds in six stages. An attacker first discovers a valid Sentry Data Source Name (DSN) – a credential embedded in JavaScript bundles or GitHub repositories that authorizes event submission to Sentry's ingestion endpoint. Because Sentry's ingestion endpoint is unauthenticated by design (it must accept errors from untrusted browser environments), the attacker submits a crafted error event whose message body contains instructions formatted to appear, when rendered, as a legitimate

Sentry remediation suggestion. When a developer subsequently asks their coding agent to investigate recent Sentry issues, the agent retrieves the injected event via MCP and – receiving no provenance signal distinguishing attacker content from genuine Sentry data – executes the embedded instructions with the developer's full system privileges.

The payload documented in the CSA research exfiltrated environment variables, AWS keys, OAuth tokens, npm registry tokens, Docker credentials, Kubernetes tokens, and CI/CD secrets – the full breadth of credentials a developer's environment might hold. Sentry acknowledged the disclosure on June 3, 2026. The company indicated that the design of public DSN-based ingestion makes platform-level remediation structurally difficult to implement [1]. That acknowledgment underscores a systemic problem: the security of an AI agent's behavior depends on assumptions about the integrity of every data source it queries, and those assumptions frequently cannot be enforced by the data source itself.

The research identified 2,388 organizations with publicly exposed injectable Sentry DSNs, including 71 among Tranco's top one-million websites [1]. Affected coding agents include Claude Code, Cursor, and Codex. The 85% exploitation success rate held across all tested agent configurations, and no detection controls in the tested configurations – endpoint security, web application firewalls, identity and access management policies, or VPN egress filtering – were effective, since the attack uses only the developer's legitimate credentials and authorized command channels. No malicious binary is dropped; no policy thresholds are crossed; no anomaly is surfaced by conventional perimeter controls.

Self-Replicating AI Worms

If agentjacking represents the seizure of individual agent sessions, self-replicating AI worms represent the horizontal spread of that seizure across interconnected agent ecosystems. The ClawWorm paper, published in March 2026 by researchers at the University of Toronto, demonstrated the first publicly reported worm targeting a production-scale AI agent framework [2]. The target was OpenClaw, an open-source agent platform with over 40,000 active instances, selected for its persistent cross-session configurations, tool-execution privileges, and cross-platform messaging capabilities.

ClawWorm's infection cycle begins with a single message. The worm first hijacks the victim agent's core configuration, establishing persistence that survives session restarts. It then executes an arbitrary payload on each subsequent reboot. Finally, without any further attacker involvement, it propagates to every newly encountered peer agent by embedding itself in the messages those agents receive. The evaluation across 1,800 trials covering four LLM backends, three infection vectors, and three payload types produced a 64.5% aggregate success rate [2]. Notably, the researchers used a small, freely available open-weight LLM as the worm's engine, demonstrating that self-replicating AI worms require no substantial commercial infrastructure to operate; a single compromised machine running an open-weight model can serve as the attack substrate.

The authors found that execution-level filtering partially mitigated dormant payloads, but that skill supply chains – the registries and package ecosystems through which agents acquire new capabilities – remained universally vulnerable across all tested backends [2]. This asymmetry matters: the most heavily defended attack surface (direct execution) is not where the worm primarily operates; the worm operates through the least-governed surface (capability acquisition).

ClawWorm's design illuminates a structural property of multi-agent ecosystems that makes conventional threat models insufficient. Network worms spread by exploiting syntactic vulnerabilities – buffer overflows, unpatched services, misconfigured authentication. AI worms spread by exploiting semantic compliance, the LLM's instruction-following behavior in response to content that arrives through channels the model treats as authoritative. Patching the underlying software alone does not close this vector; mitigating it requires changes to how agents reason about instruction provenance, combined with architectural controls that limit the blast radius of a successful injection.

The Convergence of Attack Surfaces

Agentjacking and self-replicating AI worms represent two points on a converging threat landscape rather than isolated phenomena. Both exploit the same root condition: AI coding agents operate within an implicit trust model that treats structured data returned by integrated services as authoritative instruction. The OWASP Top 10 for Agentic Applications, developed by more than 100 security researchers and practitioners and released in December 2025, codifies this as ASI01 – Agent Goal Hijacking – the top-ranked risk across all agentic AI deployment contexts [3]. OWASP describes ASI01 as occurring when "an attacker alters an agent's objectives or decision path through malicious content, exploiting the agent's planning and reasoning capabilities," noting that "because agents cannot reliably distinguish instructions from data, a single malicious input can redirect an agent to perform harmful actions using its legitimate tools and access" [3].

Separately, research published in March 2026 found that every tested coding agent, including Claude Code, GitHub Copilot, and Cursor, was vulnerable to prompt injection, with adaptive attack success rates exceeding 85% in a meta-analysis of 78 studies [6]. Multiple CVEs have been issued against the major coding agent platforms: Cursor received patches for three prompt injection vulnerabilities – CVE-2025-49150, CVE-2025-54130, and CVE-2025-61590 – documented in the National Vulnerability Database [14]; GitHub Copilot addressed CVE-2025-53773, an RCE vulnerability exploitable through prompt injection [15]. Anthropic, Google, and GitHub paid AI agent bug bounties for similar vulnerabilities during 2025, though public CVE disclosures lagged the actual remediation timeline [8].

The attack chain documented by Palo Alto Networks Unit 42 traces a representative pattern: an attacker serves a malicious MCP server whose tool descriptions embed hidden instructions that, when an agent calls those tools, arrive in the LLM's context window indistinguishable from legitimate tool responses; the

agent then calls restricted tools, leaks data, or bypasses its system prompt [9][12][13]. Tool descriptions are reviewed once when an agent first connects to a server, but tool responses flow directly into the model's context with no equivalent checkpoint – and that unguarded runtime channel is where attackers operate.

Recommendations

Immediate Actions

Organizations deploying AI coding agents should audit their MCP integrations immediately, identifying every server that surfaces external or user-controlled data. Sentry MCP integrations should be disabled where not operationally required until protective controls are in place. Any Sentry DSNs that appear in client-side bundles, public repositories, or CI/CD logs should be rotated, and new DSN issuance should be proxied through server-side relay endpoints that prevent public exposure.

Agent configurations should be updated to require explicit human approval before executing commands, installing packages, or taking any consequential action based on content retrieved from external sources. Human-in-the-loop review of externally retrieved recommendations is among the most operationally practical mitigations against agentjacking: a developer approving an MCP-retrieved recommendation is in a position to evaluate whether it is legitimate; an agent acting autonomously is not.

Short-Term Mitigations

Development teams should deploy coding agents within sandboxed containers that restrict file system access, limit environment variable visibility, and constrain network egress. Secrets management should transition from long-lived tokens held in environment variables to short-lived, scoped credentials issued at session time and revoked on completion. Cloud metadata service access should be blocked from agent execution environments. These controls do not prevent injection attempts, but they bound the blast radius when injection succeeds.

MCP server governance requires dedicated attention. Organizations should build and maintain MCP server inventories, applying version pinning and content assessment to each registered server with the same rigor applied to third-party software dependencies in supply chain security programs. Prompt injection and tool poisoning scenarios should be incorporated into red-team testing programs as first-class attack categories, not afterthoughts.

Security monitoring should be extended to agent activity logs. While conventional detection tools are unlikely to reliably identify agentjacking from behavioral signals alone, agent execution logs – which MCP-connected agents typically produce – can be instrumented to flag anomalous sequences such as credential enumeration following external content retrieval, or package installation from registries outside an organization's approved set.

Strategic Considerations

A central structural vulnerability underlying both agentjacking and self-replicating AI worms is the absence of data provenance in LLM reasoning. Most agents currently cannot reliably distinguish between instructions issued by their operators and instructions embedded in data their operators authorized them to retrieve – a limitation that active research efforts are beginning to address through instruction hierarchy enforcement and injection-resistant training. Addressing this more broadly will require changes at the protocol level – MCP and successor standards should incorporate signed provenance metadata that allows agents to verify the origin and integrity of tool responses – as well as continued investment in training approaches that reinforce skepticism toward instructions embedded in retrieved content.

Organizations should extend their AI governance frameworks to address agentic deployments explicitly. Policies governing which data sources agents may query, what actions agents may take without human approval, and how agent credential scope is bounded should be defined and enforced before agents are deployed in production, not after an incident surfaces the need. AI security posture reviews should incorporate MAESTRO threat modeling across all seven layers – from the foundation model layer through the agent ecosystem layer – to ensure that tool interface abuse vectors, cross-layer lateral movement, and skill supply chain risks receive explicit coverage [10].

The skill supply chain risk identified in ClawWorm deserves particular emphasis. Among the most underexamined surfaces in the agentic AI stack, agent capability registries and skill marketplaces currently lack the mature package signing and provenance tooling that software supply chains have developed over the past decade. An agent that installs a malicious skill extends its attack surface to every capability that skill invokes, and those capability chains currently have no equivalent governance infrastructure. Investment in agent capability governance – including content assessment, version pinning, and runtime isolation for third-party skills – should be treated as foundational security infrastructure rather than a future optimization.

CSA Resource Alignment

This research note connects to several active CSA frameworks and initiatives. The MAESTRO framework (Multi-Agent Environment, Security, Threat, Risk & Outcome) provides a seven-layer threat modeling approach for agentic AI systems, with Layer 3 (Agent Frameworks) directly addressing the reasoning loop and tool dispatch vulnerabilities exploited by agentjacking, and the Ecosystem layer covering agent supply chain and peer propagation risks of the type demonstrated by ClawWorm [10]. Security practitioners should apply MAESTRO threat models to any deployment of MCP-connected coding agents, with particular attention to cross-layer lateral movement paths.

The CSA AI Controls Matrix (AICM) includes controls addressing autonomous AI system behavior, tool invocation governance, and identity management for non-human agents – all of which apply directly to the risk categories described in this note [16]. The AI Organizational Responsibilities framework, specifically its AI Tools and Applications module, addresses the governance responsibilities that organizations bear when deploying AI agents with access to sensitive credentials and development infrastructure [17].

CSA's Agentic AI Red Teaming Guide provides a procedural framework for testing the specific vulnerability classes this note addresses, including prompt injection in multi-agent pipelines, tool misuse scenarios, and capability acquisition attacks [11]. Organizations preparing for red-team exercises should include MCP server impersonation and skill supply chain poisoning in their test plans. The CSA STAR program provides an attestation framework through which organizations can document their control posture against these risks for third-party review.

References

- [1] Tenet Security / CSA Labs. "[Agentjacking: MCP Injection Hijacks AI Coding Agents.](#)" CSA Lab Space, June 12, 2026.
- [2] Researchers at the University of Toronto. "[ClawWorm: Self-Propagating Attacks Across LLM Agent Ecosystems.](#)" arXiv:2603.15727, March 2026.
- [3] OWASP. "[OWASP Top 10 for Agentic Applications 2026.](#)" OWASP Foundation, December 2025.
- [4] Cohen, S., Bitton, R., Nassi, B. "[Here Comes The AI Worm: Unleashing Zero-click Worms that Target GenAI-Powered Applications.](#)" arXiv:2403.02817, 2024.
- [5] Elastic Security Labs. "[MCP Tools: Attack Vectors and Defense Recommendations for Autonomous Agents.](#)" Elastic, 2025.
- [6] Various Authors. "[Are AI-assisted Development Tools Immune to Prompt Injection?](#)" arXiv:2603.21642, March 2026.
- [7] SecurityWeek. "[Claude Code, Gemini CLI, GitHub Copilot Agents Vulnerable to Prompt Injection via Comments.](#)" SecurityWeek, 2025.
- [8] The Next Web. "[Anthropic, Google, and GitHub paid AI agent bug bounties, then kept quiet about the flaws.](#)" The Next Web, 2026.
- [9] Palo Alto Networks Unit 42. "[New Prompt Injection Attack Vectors Through MCP Sampling.](#)" Palo Alto Networks, 2026.
- [10] Cloud Security Alliance. "[Agentic AI Threat Modeling Framework: MAESTRO.](#)" CSA Blog, February 6, 2025.
- [11] Cloud Security Alliance. "[Agentic AI Red Teaming Guide.](#)" CSA, 2025.
- [12] Invariant Labs. "[MCP Security Notification: Tool Poisoning Attacks.](#)" Invariant Labs, 2025.
- [13] OWASP Foundation. "[MCP Tool Poisoning.](#)" OWASP Community, 2026.
- [14] National Vulnerability Database. "[CVE-2025-49150.](#)" NIST, 2025. See also CVE-2025-54130 and CVE-2025-61590 in the same database for related Cursor prompt injection disclosures.
- [15] National Vulnerability Database. "[CVE-2025-53773.](#)" NIST, 2025.

[16] Cloud Security Alliance. ["AI Controls Matrix \(AICM\)."](#) CSA, 2025.

[17] Cloud Security Alliance. ["AI Organizational Responsibilities: AI Tools and Applications."](#) CSA, 2024.