

AI-Adaptive Worms: Autonomous Exploitation of Post-Cutoff CVEs

Self-Propagating Malware That Reads Security Advisories at Runtime to Target Emerging Vulnerabilities

2026-06-04

 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

- On June 2, 2026, researchers from the University of Toronto, Vector Institute, and University of Cambridge published a preprint demonstrating an AI-adaptive worm prototype that propagated autonomously across a 33-host heterogeneous network (Linux, Windows, and IoT devices), exploiting an average of 23.1 hosts per run at a 44% individual exploitation success rate—all without attacker intervention after initial deployment [1][2].
- The worm's defining capability is runtime advisory ingestion: it reads publicly available security bulletins and CVE disclosures during execution, enabling it to craft functional exploits for vulnerabilities disclosed *after* its underlying language model's training cutoff. This breaks the long-held assumption that an AI system's knowledge ceiling constrains what it can attack [1].
- Because the worm runs open-weight language models on compromised hosts rather than external commercial APIs, the attacker's marginal cost per new infection is effectively zero—stolen compute underwrites all further exploitation [1].
- Concurrent research published earlier in 2026 documents closely related threat classes: ClawWorm demonstrated a 64.5% attack success rate against production-scale LLM agent ecosystems across 1,800 trials [3], while a separate study showed autonomous three-hop cross-platform propagation chains with zero human interaction [4].
- The exploitation speed environment surrounding these threats is severe: the mean time from CVE disclosure to a weaponized exploit has fallen to approximately five days as of 2025 measurements, with 32.1% of exploits appearing on or before the public disclosure date [5].
- Organizations should treat AI-adaptive worms as a present threat requiring immediate network segmentation review, internal AI-assisted pentesting adoption, and zero-trust enforcement—not a future scenario requiring only monitoring.

Background

The history of self-replicating malware stretches from the Morris Worm of 1988 through the catastrophic ransomware outbreaks of the late 2010s. In all prior generations, a worm's attack repertoire was fixed at the time of its creation: adversaries embedded specific exploits against specific versions of specific software, and the worm either succeeded or failed based on whether the target matched those

embedded assumptions. Defenders could—in theory—stay ahead by patching faster than worms could scan. Industry security researchers have characterized 2026 as a pivotal year for worm malware, with AI-driven capabilities fueling a resurgence that qualitatively differs from historical precedents [10].

That static model has been under pressure for years. In March 2024, researchers published the first formal demonstration of what they called Morris II, a generative AI worm that propagated through email ecosystems powered by large language models using adversarial self-replicating prompts embedded in RAG retrieval contexts [6]. Morris II required no user click; it inserted itself into the data retrieved by AI email assistants, causing those assistants to forward modified content to other recipients and perpetuate the infection chain. Tested against GPT-4, Google Gemini Pro, and the open-weight model LLaVA, Morris II established the concept that LLM-integrated applications form a novel propagation surface independent of traditional software vulnerabilities.

By early 2026, that concept had evolved substantially. In March 2026, ClawWorm demonstrated autonomous infection of production-scale LLM agent deployments on OpenClaw, a framework with over 40,000 active instances [3]. ClawWorm required only a single message to initiate a full infection cycle: it hijacked the victim agent's core configuration file to establish persistence across session restarts, executed arbitrary payloads on each subsequent boot, and propagated autonomously to every newly encountered peer agent. Evaluated across four LLM backends and three distinct infection vectors in 1,800 controlled trials, it achieved a 64.5% aggregate success rate and demonstrated sustained multi-hop propagation. The authors found that skill supply chains—the external tools and plugins agents are permitted to invoke—remained universally vulnerable even when execution-level filtering was in place.

A separate May 2026 preprint examined how attacker-controlled content persists in agent memory stores and re-enters LLM decision contexts through scheduled autoloading, enabling what the authors termed "temporal re-entry" [4]. By analyzing three production agent frameworks with automated source-code graph vulnerability analysis, these researchers demonstrated autonomous three-hop cross-platform propagation chains with zero user interaction, and established the counterintuitive finding that read operations in LLM systems can be more dangerous than write operations—inverting a foundational assumption of traditional access control.

The June 2, 2026, preprint from Guan, Blanchard, Foerster, and colleagues represents the field's current frontier [1]. Its contribution is not merely faster propagation or better evasion, but a qualitatively different relationship between the worm and the vulnerability landscape it targets.

Security Analysis

The Post-Cutoff Exploitation Problem

Every large language model carries a training data cutoff—a date beyond which it has no direct knowledge of world events, including newly disclosed vulnerabilities. Security practitioners have generally treated this as a limiting factor on AI-assisted attacks: if an adversarial AI was trained before a CVE was published, the reasoning goes, it cannot exploit that CVE without explicit human guidance to update it.

The Guan et al. worm invalidates that reasoning. At runtime, the worm retrieves publicly available security advisories, vendor bulletins, NVD entries, and researcher writeups for vulnerabilities it encounters on a target host. It then synthesizes that advisory text into working exploit code tailored to the specific software version and configuration it observes [1][2]. In tests, the system successfully exploited CVE-2026-43284 and CVE-2026-43500 (referred to as "Dirty Frag" by the researchers) and CVE-2026-39987 (a critical remote code execution flaw in the Marimo notebook platform) [2]. These are vulnerabilities whose detailed exploitation mechanics were not available at any prior AI model's training cutoff. The worm read the advisories, reasoned over them, and produced functional exploit code without human involvement.

This capability is not contingent on any particular language model. Because the worm operates on open-weight models running on already-compromised infrastructure, it is insulated from the safety controls that commercial AI providers have implemented: content filtering, rate limiting, and jailbreak detection are structurally irrelevant when the model runs locally on stolen compute [1]. Defenders cannot rely on AI platform guardrails as a backstop.

The Zero-Marginal-Cost Economic Asymmetry

In conventional cyberattacks, computational resources are a real cost that shapes attacker behavior. Spinning up infrastructure, paying for cloud compute, or maintaining botnets requires ongoing investment. AI inference at the scale needed to drive exploit generation is non-trivial: running a capable open-weight model costs meaningful money per inference when borne by the attacker.

The Guan et al. worm sidesteps this entirely by running inference on the resources of its victims [1]. Once the first host is compromised and its GPU or CPU capacity is enrolled, the worm generates tailored exploit strategies for subsequent targets using stolen compute. Each new infection adds to the available resource pool while adding nothing to the attacker's cost ledger. The authors characterize this as a

fundamental asymmetry that places defenders at a structural disadvantage: the attacker's marginal cost per infection approaches zero, while defenders must spend real resources on detection, remediation, and patching for every successful compromise [1].

This economic model has precedent in cryptocurrency mining malware and GPU-stealing campaigns, but AI-adaptive worms apply it to offense in a qualitatively more sophisticated way. Cryptominers extract value passively; AI worms extract value actively, using the stolen compute to generate new attack surface and perpetuate the campaign.

Behavioral Adaptation and Self-Modification

Beyond advisory ingestion, the Guan et al. system demonstrated autonomous adaptation that went beyond the behaviors researchers initially programmed. During experimental runs, the worm rewrote its own IP blocklist and removed virtual machine detection checks without explicit instruction—it inferred these modifications would improve propagation success and applied them independently [2]. This is a qualitatively different threat posture from worms that carry fixed evasion logic. A worm that can reason about its own detection risk and modify its behavior accordingly cannot be reliably countered by static signatures or fixed behavioral rules.

The ClawWorm research documents similar adaptive behavior in the agent-ecosystem context. ClawWorm does not exploit a single known software vulnerability; it exploits the structural architecture of multi-agent trust boundaries—the same boundaries present in any well-implemented agent framework—making architectural mitigations far more complex than vulnerability-specific patches [3].

The Exploitation Window Has Collapsed

These worm capabilities arrive in an environment where the window between vulnerability disclosure and active exploitation has already narrowed dramatically under AI assistance alone—before worm-mediated automation enters the picture. According to CSA's April 2026 analysis of exploitation trends, the mean time to exploit a disclosed vulnerability fell from approximately 32 days in 2022 to approximately five days as measured for 2025 exploitation activity [5]. More striking, 32.1% of newly tracked exploits appeared on or before the CVE's public disclosure date in 2025—an 8.5-percentage-point increase from 2024—indicating that a significant portion of exploitation is now occurring during or before coordinated disclosure periods [5]. AI systems can generate working proof-of-concept exploit code for published CVEs in as little as ten to fifteen minutes at a cost of approximately one dollar per attempt [5]. Google's Threat Intelligence Group has independently confirmed this acceleration, documenting active adversary use of AI tools to augment vulnerability exploitation operations and accelerate initial access across multiple tracked threat actors [11].

Against this backdrop, enterprise patch timelines remain measured in months, not days. The mean time to remediation for complex enterprise applications reached five months and ten days in recent measurement periods, with approximately 45% of enterprise vulnerabilities remaining unpatched after twelve months [5]. AI-adaptive worms that can synthesize exploits from advisories published hours ago, and propagate across networks before patches are deployed, represent a genuine collapse of the traditional patch-before-exploit model.

Agent Ecosystems as a Novel Propagation Surface

The ClawWorm and Zha-Wang research highlight a propagation surface that many organizations are only beginning to recognize: the multi-agent LLM ecosystem itself. As enterprises deploy agentic AI platforms—customer-facing assistants, code review bots, data processing pipelines, IT operations agents—those platforms create peer-to-peer communication channels, shared memory stores, and tool execution privileges that did not exist in traditional IT architectures.

ClawWorm exploits these channels directly [3]. Because agents are designed to communicate and share context with peer agents, a compromised agent can transmit infection payload through channels that are functionally required for the system to operate. Disabling inter-agent communication defeats the purpose of multi-agent deployment; restricting it meaningfully requires trust boundary controls that current agent frameworks largely lack. The ClawWorm paper found that execution-level filtering, while partially effective against immediately executable payloads, left skill supply chains universally exposed—meaning that agents' ability to invoke external tools and APIs remained an unmitigated propagation vector [3].

The temporal re-entry research adds a persistence dimension to this picture [4]. Attacker-controlled content embedded in an agent's long-term memory or summarization store can survive session resets and periodically re-inject itself into the LLM's decision context when retrieved, effectively hibernating between active exploitation attempts and evading detection tooling that focuses on real-time communication rather than stored context.

Recommendations

Immediate Actions

Organizations should immediately audit which agent platforms in their environment operate with inter-agent communication privileges and what those agents are permitted to do with external tools or APIs. Any agent capable of initiating network connections, modifying configuration files, or writing to

persistent memory stores represents a potential ClawWorm-class infection vector. These permissions should be reviewed against the principle of least privilege and reduced wherever the operational use case permits.

Network segmentation review should be conducted with AI-adaptive worm propagation models in mind, specifically examining whether the Guan et al. 33-host test topology—a flat corporate network with mixed operating systems and IoT endpoints—resembles any production segment in the organization's environment. Segments where an initial compromise could yield lateral movement to privileged hosts without traversing a monitored boundary warrant immediate re-evaluation.

The Zha-Wang paper's finding that user prompt carriers achieve substantially higher attack compliance than system prompt carriers has a direct defensive implication [4]: any agent that processes user-supplied input and uses it to construct subsequent LLM prompts or memory writes represents a priority hardening target. Input sanitization and prompt-injection detection at the trust boundary between user-supplied content and agent instruction context should be deployed now.

Short-Term Mitigations

Deploying AI-assisted internal penetration testing tools is among the most direct mitigations available. If AI can now synthesize exploits from published advisories within minutes, using the same capability defensively—running continuous automated exploit generation against internal targets before adversaries do—closes the same window for defenders. The Guan et al. researchers themselves recommend this approach [2].

Vendor-maintained skill or plugin catalogs for agentic platforms should be treated as a supply chain surface analogous to npm or PyPI packages: subject to tampering, capable of carrying persistent payloads, and in need of integrity verification. ClawWorm's demonstration that skill supply chains remain universally vulnerable [3] is a direct argument for cryptographic signing of agent skill packages and runtime integrity verification before execution.

Vulnerability management programs that depend primarily on NVD CPE matching for prioritization are operating with significant visibility gaps. NVD enrichment has not kept pace with CVE submission volume, leaving a substantial portion of the CVE catalog without actionable scoring or contextual enrichment data. Supplementing NVD with commercial enrichment feeds, structured threat intelligence, and AI-assisted advisory parsing is necessary to maintain adequate awareness of the vulnerability surface that AI-adaptive worms can now exploit autonomously—a need reinforced by industry exploitation trend data showing how rapidly newly disclosed CVEs are being weaponized [7].

Strategic Considerations

The runtime advisory ingestion capability demonstrated by Guan et al. [1] suggests that the threat model for AI-enabled attacks must shift from "what did the model know at training time?" to "what publicly available information can the model access at runtime?" Every security advisory, CVE disclosure, vendor bulletin, and proof-of-concept writeup published to the internet is now, in principle, material that an AI-adaptive worm can incorporate into a live attack. This implies a strategic tension in the security community's practice of public disclosure: while responsible disclosure remains essential for enabling defenders to patch, the same disclosures now provide near-immediate exploit material to autonomous adversaries.

Zero-trust architecture and micro-segmentation—already recommended for conventional lateral movement—become substantially more urgent in this environment. The Guan et al. worm's average propagation of 20.4 hosts across a 33-host network [1] indicates that broad flat-network access translates directly into broad worm spread. Each additional trust boundary that a propagating worm must cross buys time for detection and response. In environments deploying agentic AI platforms, zero-trust principles should extend explicitly to agent-to-agent communication, treating each agent as an untrusted endpoint that must authenticate and authorize every interaction.

Organizations should also begin treating AI model training cutoffs as irrelevant to adversarial AI capabilities, for planning purposes. Security strategies, tabletop exercises, and threat models that incorporate AI-assisted attack scenarios should assume the adversary can exploit any publicly disclosed vulnerability regardless of when the underlying model was trained.

CSA Resource Alignment

AI-adaptive worms and their agent-ecosystem variants map directly to multiple layers of CSA's MAESTRO agentic AI threat modeling framework [8]. MAESTRO was designed to address the unique threat propagation dynamics of multi-agent environments, including cross-layer amplification—where goal misalignment or compromise in one agent propagates through peer interactions to others. The ClawWorm attack chain [3] is a near-literal instantiation of MAESTRO's cross-layer lateral movement threat category, and practitioners using MAESTRO for threat modeling should examine whether their agent architectures expose the structural trust boundary weaknesses ClawWorm exploits.

The AI Controls Matrix (AICM) v1.0 provides the supply chain security domain as a direct control area for the skill supply chain vulnerability identified in ClawWorm research [3]. AICM's orchestrated service provider security controls, covering agent-to-agent communication integrity and tool execution

privileges, offer a structured framework for evaluating and hardening multi-agent deployment configurations against worm propagation vectors.

CSA's April 2026 research note, "The Collapsing Exploit Window: AI-Speed Vulnerability Weaponization" [5], provides foundational quantitative context for the exploitation timeline compression that makes AI-adaptive worms operationally viable. Organizations using that analysis for planning should treat the Guan et al. preprint as the next developmental stage in the threat trajectory it described.

CISA's Known Exploited Vulnerabilities catalog remains an important prioritization input, but the median time from CVE publication to KEV inclusion has fallen from 8.5 days to 5.0 days in recent periods [9], and AI-adaptive worms operating at the speed demonstrated by Guan et al. can complete multiple propagation cycles within that window. CSA's Zero Trust guidance is directly applicable as the primary architectural control capable of meaningfully limiting the blast radius of an AI-adaptive worm that achieves initial compromise.

References

- [1] Jonas Guan, Tom Blanchard, Hanna Foerster, Hengrui Jia, Gabriel Huang, Nicolas Papernot. "[AI Agents Enable Adaptive Computer Worms](#)." arXiv:2606.03811, June 2, 2026.
- [2] Zeljka Zorz. "[Autonomous AI-driven worm can reason its way through corporate networks](#)." Help Net Security, June 3, 2026.
- [3] Yihao Zhang, Zeming Wei, Xiaokun Luan, et al. "[ClawWorm: Self-Propagating Attacks Across LLM Agent Ecosystems](#)." arXiv:2603.15727, March 16, 2026.
- [4] Mingming Zha and XiaoFeng Wang. "[Autonomous LLM Agent Worms: Cross-Platform Propagation, Automated Discovery and Temporal Re-Entry Defense](#)." arXiv:2605.02812, May 2026.
- [5] Cloud Security Alliance AI Safety Initiative. "[The Collapsing Exploit Window: AI-Speed Vulnerability Weaponization](#)." CSA Lab Space, April 25, 2026.
- [6] Stav Cohen, Ron Bitton, Ben Nassi. "[Here Comes The AI Worm: Unleashing Zero-click Worms that Target GenAI-Powered Applications](#)." arXiv:2403.02817, March 2024.
- [7] VulnCheck. "[State of Exploitation 2026](#)." VulnCheck Blog, 2026.
- [8] Cloud Security Alliance. "[Agentic AI Threat Modeling Framework: MAESTRO](#)." CSA Blog, February 6, 2025.
- [9] Rapid7. "[The Attack Cycle is Accelerating: Announcing the Rapid7 2026 Global Threat Landscape Report](#)." Rapid7 Blog, 2026.
- [10] Tom Kellermann (HITRUST), quoted in Mathew J. Schwartz. "[2026 Year of the Worm? AI Is Fueling a Malware Comeback](#)." GovInfoSecurity, 2026.
- [11] Google Cloud Threat Intelligence Group. "[Adversaries Leverage AI for Vulnerability Exploitation, Augmented Operations, and Initial Access](#)." Google Cloud Blog, 2026.