

CSAI Foundation | Cloud Security Alliance

# AI Autonomous Vulnerability Hunters: The Offense-Defense Gap

Implications for Enterprise Patch Strategy in an Era of Machine-Speed Exploitation

2026-05-15

 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.*

---

# Table of Contents

- Executive Summary ..... 4
- Introduction and Background ..... 4
- The Autonomous Vulnerability Discovery Landscape ..... 5
  - From Human Researchers to Machine Hunters
  - The Dual-Use Nature of Defensive AI Research
  - The Commercial Ecosystem
- The Widening Offense-Defense Gap ..... 8
  - Time-to-Exploit Goes Negative
  - Nation-State Actors Operationalize AI
  - The Economics of AI-Enabled Offense
- Implications for Enterprise Patch Strategy ..... 10
  - The Collapse of the Traditional Patch Cycle
  - Risk-Based Prioritization in an AI Era
  - Compensating Controls as a Risk Bridge
  - Detection and Response as a Patch Complement
- Strategic Recommendations ..... 13
  - Immediate Actions
  - Short-Term Program Improvements
  - Strategic Investments
- CSA Framework Alignment ..... 15
- Conclusions ..... 16
- References ..... 17

# Executive Summary

The cybersecurity profession is experiencing a capability inflection point that fundamentally alters the economics and timing assumptions underlying enterprise vulnerability management. Artificial intelligence systems can now autonomously discover, analyze, and exploit software vulnerabilities faster than organizations can detect, triage, and remediate them. This is not a speculative future risk – it describes the threat environment as it exists in May 2026.

Mandiant's incident response data shows that the mean time to exploit a disclosed vulnerability has reached negative seven days [1], meaning that on average, exploitation is underway before a patch is even available. Meanwhile, enterprise organizations require approximately 55 days to patch half of their critical vulnerabilities [9]. This gap – measured in weeks on the defensive side versus hours or less on the offensive side – represents a structural asymmetry that no amount of process improvement alone can close.

The offensive capability driving this gap has matured rapidly. Research published in 2024 demonstrated that GPT-4 could autonomously exploit 87% of one-day vulnerabilities when provided with CVE descriptions [2]. By April 2026, Anthropic's Claude Mythos Preview model – purpose-built for security research under Project Glasswing – demonstrated a 72.4% autonomous success rate against Firefox's JavaScript shell, identifying thousands of high-severity vulnerabilities that had survived decades of human review [3]. The same model was deemed too capable to release broadly. Simultaneously, DARPA's AI Cyber Challenge concluded its final competition with participating teams discovering 86% of synthetic vulnerabilities, patching 68% of those found, and uncovering 18 real zero-day vulnerabilities across 54 million lines of production code – at an average task cost of \$152 [4].

These developments demand a fundamental revision of how enterprise security teams think about patching. The traditional model – centered on cyclical patch windows, CVSS-based triage, and remediation timelines measured in weeks – was designed for a world where human researchers and state-sponsored teams represented the frontier of vulnerability research. That world no longer exists. This paper analyzes the current offensive-defensive landscape, documents the specific mechanisms by which AI has widened the gap, and offers a practical strategic framework for enterprise security organizations responding to this new reality.

---

## Introduction and Background

For most of the past two decades, the implicit model underlying enterprise vulnerability management assumed a meaningful window between disclosure and exploitation. Vendors disclosed vulnerabilities, security teams assessed severity, administrators scheduled maintenance windows, and patches were

deployed – often over weeks or months – before most organizations faced active exploitation. This model was imperfect, but it was workable: the majority of attackers lacked the expertise to rapidly develop functional exploits for newly disclosed vulnerabilities, and even sophisticated nation-state actors required days or weeks to weaponize new findings.

That window has collapsed. The primary driver is not simply that adversaries have become more organized or well-resourced, though both are true. The primary driver is that AI systems have commoditized the technical skills required to discover and exploit software vulnerabilities, dramatically lowering the barrier to offensive capability and compressing timelines that were once measured in days to timelines measured in hours or minutes.

The term "autonomous vulnerability discovery" encompasses a broad range of AI-assisted and fully autonomous techniques for finding security flaws in software. At the least-automated end, AI systems assist human researchers by accelerating code review, suggesting likely vulnerability patterns, and generating proof-of-concept exploit candidates. At the more-automated end, systems operate with minimal human involvement: ingesting source code or binary applications, reasoning about potential flaws, generating test cases, and confirming exploitability – all without meaningful human steering. The research community has published systems spanning this entire range, and evidence from incident response suggests that both criminal and nation-state actors have begun deploying operationalized variants.

Understanding the implications of this shift requires examining three interrelated questions: How capable are current AI systems at offensive vulnerability research? How have adversaries actually deployed these capabilities in the wild? And what must enterprise security organizations do differently as a consequence? This paper addresses each in turn before synthesizing strategic guidance grounded in CSA frameworks.

---

## The Autonomous Vulnerability Discovery Landscape

### From Human Researchers to Machine Hunters

The academic and industry research establishing AI's offensive capability has been accumulating since at least 2023, but 2024 and 2025 produced a series of results that decisively crossed practical thresholds. The foundational study, published in April 2024, demonstrated that GPT-4 – operating as an agent with access to shell, Python, and web browsing tools – could autonomously exploit 87% of a curated set of one-day vulnerabilities when given their CVE descriptions. Crucially, no other model tested came close: GPT-3.5, open-source LLMs, and commercial scanning tools like ZAP and Metasploit all registered near-zero success rates on the same benchmark [2]. The result established that a specific capability threshold had been crossed and that it was narrowly concentrated in the frontier model tier.

That concentration has since broadened. Subsequent research produced systems specifically engineered for autonomous penetration testing, integrating LLM reasoning with structured attack methodologies. RapidPen, published in early 2025, achieved automated IP-to-shell access in an average of 200 to 400 seconds at a cost between \$0.30 and \$0.60 per run [5]. A separate research effort produced Excalibur, a system built on an advanced LLM that autonomously compromised four of five hosts in a realistic Active Directory engagement at a total API cost of \$28.50 [5]. The economics of capability are shifting as dramatically as the capabilities themselves.

The DARPA AI Cyber Challenge (AixCC), which concluded its final competition in 2025, offers the most comprehensive public dataset on AI autonomous security capability at scale. Seven finalist teams competed with Cyber Reasoning Systems (CRS) – fully autonomous AI agents tasked with finding, confirming, and patching vulnerabilities in large, realistic codebases. The final results were striking: participating systems identified 86% of synthetic vulnerabilities embedded in the competition infrastructure, an improvement from 37% at the semifinal stage, and successfully patched 68% of what they found. More consequentially, the systems discovered 18 previously unknown vulnerabilities across 54 million lines of production code from real open-source projects, providing patches for 11 of them – all at an average per-task cost of approximately \$152 [4]. These were not toy examples or narrow benchmarks; they were real flaws in real software that human reviewers had missed.

## The Dual-Use Nature of Defensive AI Research

The same capabilities being weaponized offensively are also being pursued for defensive purposes, and in April 2026, Anthropic made the defensive use case explicit with the launch of Project Glasswing. Framing the effort as a preemptive measure against AI-driven cyberattacks, Anthropic announced that Claude Mythos Preview – an unreleased frontier model judged too capable for general release – would be made available exclusively to a consortium of founding organizations including Amazon Web Services, Apple, Broadcom, Cisco, CrowdStrike, Google, JPMorgan Chase, the Linux Foundation, Microsoft, NVIDIA, and Palo Alto Networks [3]. Anthropic committed \$100 million in model usage credits to the initiative and provided grants to open-source security foundations.

The capabilities demonstrated by Mythos Preview represent a qualitative advance over prior models. Testing against Firefox's JavaScript shell produced a 72.4% autonomous exploit success rate – compared with near-zero rates for predecessor models – on vulnerabilities that had survived decades of human security review and millions of automated test runs [3]. Critically, Project Glasswing exposed a challenge that extends beyond the AI capability question: fewer than one percent of the vulnerabilities Mythos found were actually patched [16]. The gap between AI's ability to discover vulnerabilities and the enterprise's ability to remediate them is not primarily a detection problem anymore. It is a capacity and prioritization problem, and one that AI-driven discovery – whether wielded by defenders or attackers – will only intensify.

Bruce Schneier, analyzing Project Glasswing in April 2026, characterized the dynamic as a dual-use paradox – the same capability that enables comprehensive defensive vulnerability discovery simultaneously enables comprehensive offensive exploitation, and defenders hold an advantage only as long as comparable capability remains out of adversary hands, a window Schneier estimated in months rather than years [6]. This framing captures the essential strategic problem. Defensive use of AI vulnerability discovery provides a genuine but time-limited advantage, because the same capabilities are simultaneously being pursued by well-resourced offensive actors.

## The Commercial Ecosystem

The defensive security market has moved rapidly to commercialize AI penetration testing and vulnerability discovery capabilities. Horizon3.ai's NodeZero platform had run more than 150,000 autonomous penetration tests in production enterprise environments by early 2026 [5]. Intruder, a GCHQ-backed UK startup, launched AI-driven pentesting agents that replicate human pentester methodology and deliver results in minutes, with early access producing findings that human-led engagements had previously missed [7]. Alias Robotics' open-source CAI framework documented a 156-fold cost reduction in penetration testing – an engagement priced at \$17,218 with a human team conducted for \$109 using AI agents, while completing 3,600 times faster [5].

The broader market context reinforces how seriously the industry views this shift. The AI penetration testing segment attracted more than \$665 million in disclosed venture capital investment, producing two unicorn valuations, while the overall penetration testing market was valued at approximately \$2.74 billion in 2025 with projections reaching \$6.25 to \$7.41 billion by 2033 [5]. This investment activity reflects both genuine capability and genuine demand: enterprises facing machine-speed exploitation need machine-speed discovery on the defensive side.

The emergence of 70 or more new AI-driven offensive security tools between 2024 and 2026 [8] means that capability once concentrated in frontier model APIs is increasingly accessible via commercial and open-source tooling. The cost asymmetry that historically favored defenders – where a single vulnerability exploited could be worth millions to an attacker, but discovering it required expert humans – has inverted. Exploitation is becoming cheaper and faster; vulnerability discovery is becoming automated; and the economic logic that once constrained attackers to high-value targets is eroding.

---

# The Widening Offense-Defense Gap

## Time-to-Exploit Goes Negative

The most consequential quantitative shift in the vulnerability management landscape is the inversion of the time-to-exploit metric. Mandiant's research group measured an average time-to-exploit of 63 days in 2018, meaning that organizations with mature patch programs had meaningful time to remediate before active exploitation began. By 2024, that figure had reached negative one day [9] – exploitation was, on average, occurring the day before public disclosure. Mandiant's M-Trends 2026 report, built on more than 450,000 hours of incident response engagements, puts the current figure at negative seven days [1].

A negative time-to-exploit does not mean that every vulnerability is exploited before its patch is released. It means that the distribution of exploitation timelines has shifted dramatically leftward – more vulnerabilities are being exploited before disclosure, the fastest exploitations are occurring within hours of disclosure, and the window available to defenders has compressed to the point that traditional patch management cycles are structurally mismatched with the threat. Approximately 56% of vulnerabilities are weaponized within the first month after disclosure, with a significant subset exploited before any public patch is available [10].

Several mechanisms explain the collapse of the exploitation window. First, AI systems operating at machine speed can analyze newly disclosed CVE advisories, map them to code repositories, generate exploit candidates, and test functionality without human involvement – collapsing a process that previously required days of expert human effort into a process that can complete in hours. Second, sophisticated threat actors have invested in monitoring disclosure pipelines – vendor notification systems, security researcher channels, and code commit histories – to gain early access to vulnerability information before public announcement. Third, the proliferation of proof-of-concept exploit code on public repositories means that even less-sophisticated actors can weaponize newly disclosed vulnerabilities within hours of disclosure.

The implications for enterprise defenders are direct. Defenders require an average of 55 days to patch 50% of critical vulnerabilities in their environments [9] – a figure that reflects real operational constraints including testing requirements, change management processes, application compatibility concerns, and limited patching personnel. The gap between a 55-day average patching time and a negative-seven-day mean exploitation time is not a gap that can be closed by exhortation or modest process improvement. It requires a fundamental rethinking of how organizations prioritize, sequence, and automate remediation.

## Nation-State Actors Operationalize AI

Google's Threat Intelligence Group (GTIG) documented that more than 57 distinct nation-state threat groups were actively using AI capabilities for cyber operations by early 2025 [11]. The patterns of AI use vary by actor sophistication and objective, but several consistent behaviors have emerged. Chinese APT groups

have used AI to accelerate reconnaissance, troubleshoot offensive code, and identify lateral movement paths within compromised networks. North Korean state-linked operator APT45 sent thousands of repetitive prompts to AI systems to analyze CVE data and validate proof-of-concept exploits, effectively using AI as a high-throughput vulnerability research assistant [12]. A Russian state-linked threat group deployed malware that queried AI models in real time to determine post-compromise actions within compromised Ukrainian networks – a form of AI-guided adaptive attack that had not previously been documented in the wild [12].

The most significant development in this space occurred in May 2026, when Google's researchers documented what they assessed to be the first confirmed case of a zero-day exploit developed using AI [13]. The criminal threat actor responsible – whose planned mass exploitation campaign was disrupted through Google's proactive counterdiscovery – had produced exploit code with several characteristics distinguishing it from human-authored work: educational docstrings explaining the vulnerability in textbook language, a hallucinated CVSS score embedded in the code comments, and clean, structured Python consistent with LLM output rather than the idiosyncratic style of experienced human exploit developers [13]. Google's counter-response prevented the planned exploitation campaign, but the incident establishes that the threshold separating AI-assisted from AI-generated exploit development has been crossed in live attack operations.

Earlier in 2026, a separate criminal operation attributed to threat group TeamPCP compromised several open-source software repositories, including those associated with the LiteLLM AI gateway library and the Trivy vulnerability scanner – a supply-chain attack that, notably, used AI to develop the compromise tooling against AI infrastructure [12]. The recursive dimension of that incident – AI-built malware attacking AI security tools – illustrates how the AI threat landscape is developing in ways that traditional threat modeling frameworks were not designed to anticipate.

## The Economics of AI-Enabled Offense

The offense-defense gap is partly a capability gap and partly an economics gap, and the economic dimensions deserve explicit analysis because they inform both the threat trajectory and the strategic response. Human vulnerability research is expensive. Senior researchers with the skills to discover novel vulnerabilities in hardened software command compensation that makes broad scanning of enterprise attack surfaces impractical for most adversaries. AI systems do not have this constraint. An autonomous system that can analyze millions of lines of code, generate thousands of exploit candidates, and test exploitability without human involvement operates at a cost structure that makes large-scale vulnerability discovery economically viable even for modestly resourced threat actors.

The CAI penetration testing cost data – \$109 versus \$17,218 for equivalent engagements – understates the true economic implication, because it compares AI-assisted work against human work on the same task scope [5]. The more consequential comparison is between what a motivated adversary can now afford to do and what they could previously afford to do. At \$0.30 to \$0.60 per automated penetration testing run [5], a

threat actor can afford to run comprehensive assessments against thousands of target environments at a cost that would have been prohibitive with human teams. The economics do not merely make existing attacks cheaper; they make previously uneconomical attacks viable.

This cost structure also affects the defensive calculus. Defenders who rely primarily on scheduled penetration testing – quarterly or annual engagements – are operating on a cadence that was calibrated for the economics of human-led testing. When continuous AI-driven assessment is available for a fraction of the cost, the argument for infrequent manual testing weakens substantially. The same AI that enables continuous offensive probing can enable continuous defensive probing – but only if organizations consciously reorient their vulnerability management programs to take advantage of it.

---

## Implications for Enterprise Patch Strategy

### The Collapse of the Traditional Patch Cycle

The traditional enterprise patch management cycle was designed around a set of assumptions that no longer hold. It assumed that critical vulnerabilities would be publicly disclosed before they were widely exploited, providing at minimum a few days for organizations to begin remediation. It assumed that the complexity of developing functional exploits would create a natural buffer, with exploitation requiring specialized expertise available only to sophisticated attackers targeting high-value organizations. And it assumed that a risk-stratified approach to patching – addressing critical findings first, significant findings on a 30-day cycle, and moderate findings on a 60-to-90-day cycle – would adequately manage exposure for most environments.

Each of these assumptions has been undermined. The mean time-to-exploit data from Mandiant confirms that exploitation is now, on average, underway before patches are available [1]. The democratization of exploit development via AI has eliminated the expertise barrier that previously concentrated rapid exploitation among sophisticated actors. And the CISA Known Exploited Vulnerabilities catalog, which documented 245 new actively exploited vulnerabilities in 2025 alone for a cumulative total of approximately 1,484 entries [14], demonstrates that the universe of vulnerabilities requiring urgent attention has grown far beyond what traditional cycle-based patching can adequately address.

None of this means that patch management is futile. The data consistently shows that many exploitations target vulnerabilities for which patches have been available for months or years – attackers pursue the path of least resistance, and unpatched older vulnerabilities remain attractive targets. The discipline of sustained patch hygiene has not lost its value. What has changed is the strategic framing. Patch management alone, even if executed flawlessly against traditional timelines, is no longer sufficient to protect against the most

active and capable threat actors. It must be embedded within a broader vulnerability management program that incorporates real-time exposure monitoring, risk-based prioritization, compensating controls, and active threat intelligence consumption.

## Risk-Based Prioritization in an AI Era

The CVSS scoring system, while useful as a vulnerability severity taxonomy, was not designed to drive patching sequencing decisions in an environment where exploitation timelines are measured in hours. A CVSS 9.8 score communicates severity, but it communicates nothing about whether a vulnerability is being actively exploited, whether the affected software is exposed in a given organization's specific architecture, or whether a compensating control already mitigates the risk. Organizations that sequence patching by CVSS score alone will systematically overprioritize vulnerabilities that are theoretically severe but not actively targeted while underattending to vulnerabilities that are lower-scored but under active exploitation.

The CISA Known Exploited Vulnerabilities catalog provides the most operationally relevant signal available for patching prioritization, because it reflects confirmed real-world exploitation rather than theoretical severity. The Binding Operational Directive requiring federal agencies to remediate KEV-listed vulnerabilities within defined windows [14] provides a model that enterprise organizations can adapt: treat active exploitation as the primary triage criterion, ahead of severity score, and ensure that any vulnerability confirmed as actively exploited receives immediate attention regardless of where it falls in the scheduled patching cycle.

Beyond KEV-based triage, effective risk-based prioritization in an AI-era threat environment requires integrating asset context into vulnerability decisions. A CVSS 7.5 vulnerability in an internet-exposed service handling customer authentication data is more urgent than a CVSS 9.0 finding in an air-gapped internal system with no direct path to sensitive data. Organizations that have not invested in asset inventory and exposure mapping will find that their patching decisions are systematically miscalibrated to the actual risk profile of their environments.

Threat intelligence subscription and consumption presents a second dimension of prioritization that many enterprise organizations underutilize. Real-time intelligence about which vulnerabilities are being actively targeted by threat actors relevant to an organization's sector and geography provides actionable signal that can inform patch sequencing decisions before exploitation reaches the scale necessary to appear in public datasets like the KEV catalog. The emergence of AI-driven threat intelligence platforms that correlate vulnerability disclosure data with observed exploitation activity across telemetry networks provides an opportunity to operationalize this signal at scale.

## Compensating Controls as a Risk Bridge

Given that patching timelines cannot be compressed to match the negative time-to-exploit reality for the full breadth of an enterprise attack surface, compensating controls become an essential element of the strategic response rather than an ad hoc fallback. A compensating control is any technical or procedural measure that reduces the exploitability or impact of a vulnerability without addressing the vulnerability directly – the classic example being temporary network segmentation to isolate an unpatched system, or additional authentication controls on a path to a vulnerable service.

The compensating controls strategy requires tight integration between threat intelligence, asset management, and network security functions. When a critical vulnerability is disclosed in software known to be present in the environment, the sequence of compensating control deployment must be rapid enough to be meaningful in an environment where exploitation may begin within hours. This argues for pre-approved playbooks for the most common compensating control patterns – network isolation, WAF rule deployment, credential rotation, additional authentication requirements – that can be executed without full change management review when the threat intelligence signal meets defined criteria.

Compensating controls are a risk bridge, not a risk elimination mechanism. They reduce exposure during the window between disclosure and remediation, but they introduce their own operational complexity and potential for misconfiguration. Organizations that rely on compensating controls as a long-term substitute for patching accumulate technical debt that becomes increasingly difficult to retire. The appropriate role of compensating controls is to buy time for proper remediation while maintaining defensible posture, not to indefinitely defer the remediation work.

## Detection and Response as a Patch Complement

The structural mismatch between exploitation timelines and remediation timelines means that a meaningful fraction of exploitation attempts will succeed before patches are deployed. This is not a failure of the patching program – it is a mathematical consequence of the current time-to-exploit environment that no realistic patching program can fully eliminate. Organizations that respond to this reality by redoubling patching velocity without also investing in detection and response will find that their improved patching metrics coexist with continued breach activity.

The necessary complement to improved patching velocity is the capacity to detect exploitation activity and contain its consequences. AI tools that can now autonomously discover vulnerabilities can also autonomously hunt for evidence of exploitation in log data, network traffic, and endpoint telemetry. The same LLM-driven reasoning that enables automated exploit development can be applied to the defensive problem of correlating behavioral signals that indicate a vulnerability is being actively targeted. Security operations capabilities that combine automated detection with playbook-driven response provide a hedge against the exploitation that patching cannot prevent.

Mandiant's finding that initial access broker handoffs – the transfer of compromised access from the actor who obtained it to the downstream threat actor who will monetize it – now average 22 seconds [1] underscores the temporal dimension of post-exploitation detection. Dwell time reduction, which measures the interval between initial compromise and detection, has improved substantially in mature organizations over the past decade, but the acceleration of attacker operational tempo means that even significantly reduced dwell times may not prevent meaningful data access or lateral movement. Detection capabilities must be calibrated not just to identify compromise but to enable rapid containment before attackers can establish persistent footholds.

---

## Strategic Recommendations

### Immediate Actions

Enterprise security organizations should treat three immediate actions as prerequisites for any further strategic adaptation to the AI threat environment. First, establish a real-time integration between asset inventory systems and the CISA KEV catalog, with automated alerting that surfaces KEV additions affecting known assets within minutes of the catalog update rather than through weekly or monthly reporting cycles. The KEV catalog has become the most operationally actionable prioritization signal available and should be consumed accordingly.

Second, review and update the compensating controls playbook to ensure that pre-approved responses exist for the most common vulnerability categories – remote code execution in internet-facing services, authentication bypass in identity systems, privilege escalation in operating systems – and that those responses can be executed within hours rather than days. The playbook review should include explicit criteria for when compensating controls are authorized without full change management review.

Third, assess the current AI capabilities available for defensive vulnerability scanning and determine whether they are being used to their potential. Platforms offering continuous, AI-driven vulnerability assessment are now commercially available and cost-competitive with scheduled human-led engagements. Organizations that are not yet running continuous automated assessment are operating with a significant visibility gap relative to the threat environment.

### Short-Term Program Improvements

Over a three-to-six-month horizon, organizations should pursue a set of structural improvements to their vulnerability management programs that align the program's operating model with the current threat environment. The most consequential of these is the transition from CVSS-primary to exploitation-context-

primary prioritization – replacing or supplementing the CVSS score as the primary triage criterion with a composite metric that weights confirmed exploitation activity, asset exposure, and asset criticality alongside theoretical severity.

Patching velocity for internet-exposed services deserves particular attention in this context. Internet-facing applications, network appliances, remote access infrastructure, and identity services represent the attack surface most accessible to AI-driven exploitation pipelines that do not require prior access to the target environment. Dedicating disproportionate patching resources to this subset of the attack surface – and maintaining it to a substantially tighter timeline than the overall environment – concentrates defensive effort where the marginal risk reduction is greatest.

Threat intelligence consumption should be operationalized rather than treated as a strategic input to quarterly risk reviews. Real-time intelligence feeds that identify vulnerabilities under active exploitation by threat actors relevant to the organization's sector should flow directly into the patching prioritization queue rather than arriving via reporting cycles. Many organizations receive this intelligence but consume it too slowly to translate it into operational action.

## Strategic Investments

Over the twelve-to-eighteen-month horizon, the strategic priority is building organizational capability for sustainable operation in a machine-speed threat environment. Two investments are particularly consequential. The first is in exposure management – the continuous, automated mapping of the organization's externally accessible attack surface, including shadow IT and third-party dependencies, against the current vulnerability landscape. Exposure management extends beyond traditional vulnerability scanning to encompass the full scope of exploitable paths from an external attacker's perspective, providing the contextual richness that prioritization decisions require.

The second strategic investment is in AI-augmented security operations that can detect and respond at the speed the threat requires. The combination of AI-driven detection – correlating behavioral signals across endpoint, network, and identity telemetry – with pre-approved, playbook-driven automated response enables containment actions that complete in seconds rather than the hours that human-escalation workflows typically require. This capability is not intended to replace human judgment in incident response; it is intended to contain the blast radius of exploitation during the window between initial detection and human analyst engagement.

Organizations operating in sectors targeted by sophisticated nation-state actors should additionally assess the maturity of their supply chain security posture. TeamPCP's 2026 compromise of LiteLLM and Trivy repositories [12] illustrates that AI security tooling is itself becoming an attack surface, and that the integrity of open-source dependencies consumed in security infrastructure warrants the same scrutiny applied to production application dependencies.

# CSA Framework Alignment

The threat dynamics described in this paper intersect with several CSA frameworks that provide structured guidance for organizational response.

The MAESTRO framework – CSA's agentic AI threat modeling methodology introduced in February 2025 – provides the most directly applicable structured approach for reasoning about AI-driven offensive threats [15]. MAESTRO's seven-layer architecture covers the full agentic AI stack from foundation models through agent ecosystems, and its emphasis on threats arising from autonomous reasoning and action maps directly to the offensive AI systems discussed in this paper. Organizations that have not yet applied MAESTRO to their security architecture analysis should use the AI autonomous vulnerability discovery threat as a concrete scenario for driving that exercise – specifically, the question of how an autonomous agent operating outside the organization's control might traverse its attack surface, chain vulnerabilities, and persist without triggering existing detection controls.

The AI Controls Matrix (AICM) v1.0 provides the control framework for assessing readiness across the relevant domains. AICM's AI supply chain security domain is particularly relevant to the supply chain risk illustrated by the LiteLLM and Trivy compromises, where the attack surface extended to the security tooling itself. Organizations evaluating their exposure to AI-enabled threats should assess AICM controls in the model security, data security, and application security domains as a baseline for understanding where gaps exist relative to the current threat environment.

CSA's Zero Trust guidance reinforces the architectural principle that is most directly applicable to the defensive challenge: assuming that perimeter-based defenses will fail and building security posture around the assumption of breach. In a threat environment where exploitation may occur before patches are available and where the handoff between initial compromise and further exploitation is measured in seconds, the compensating controls and detection capabilities discussed in this paper are most effective when deployed within an architecture that does not assume network location confers trust. Micro-segmentation, continuous identity verification, and least-privilege access enforcement each reduce the attacker's lateral movement options following exploitation, limiting the impact of the exploitation events that patching cannot prevent.

The Cloud Controls Matrix (CCM) provides control categories relevant to the operational disciplines required for effective adaptation: vulnerability and patch management controls under the TVM domain, incident response controls under the SEF domain, and supply chain security controls under the STA domain. Organizations seeking to assess their current control coverage against the recommendations in this paper should map their existing controls to these CCM domains as a structural gap analysis.

---

## Conclusions

The emergence of AI systems capable of autonomous vulnerability discovery and exploit generation is not a technology risk that enterprises can manage by updating their security policies. It represents a structural change in the economics and timing of offensive capability that requires a corresponding structural change in how organizations approach vulnerability management, defensive tooling, and security operations.

The core challenge is a temporal one. When exploitation routinely occurs before patches are available [1], when AI systems can traverse an attack surface at a cost of dollars rather than analyst-hours [5], and when the first AI-generated zero-day exploit has already been documented in active attack campaigns [13], the 55-day average enterprise patching timeline [9] is not a performance problem – it is a structural mismatch that patching velocity improvements alone cannot resolve. Addressing it requires simultaneously compressing patching timelines for the highest-risk exposures, deploying compensating controls as a risk bridge during the remediation window, investing in AI-augmented detection that can identify exploitation activity before attackers establish persistence, and building the organizational intelligence flows necessary to act on exploitation signals in hours rather than days.

The encouraging dimension of this challenge is that the same AI capabilities that have widened the offense-defense gap are available to defenders. Project Glasswing represents a deliberate effort to ensure that the frontier of autonomous vulnerability discovery provides defenders with a head start rather than attackers, and the commercial ecosystem has produced continuous AI-driven assessment tools that are cost-accessible to enterprises of meaningful scale [5]. The question is not whether AI tools are available for defensive use – they are – but whether enterprise organizations are deploying them with the urgency and structural integration that the current threat environment requires.

CSA's frameworks – MAESTRO for agentic threat modeling, AICM for control assessment, CCM for operational disciplines, and Zero Trust guidance for architectural orientation – provide the structured vocabulary and control catalog through which organizations can systematically evaluate their posture and plan their adaptation. The threat environment has moved faster than most enterprises' response programs. The work of closing that gap begins with an honest assessment of where the current program's assumptions are no longer valid and what structural changes are necessary to bring the program into alignment with the world as it actually exists.

## References

- [1] Google Cloud / Mandiant. "[M-Trends 2026: Data, Insights, and Strategies From the Frontlines](#)." Google Cloud Blog, March 2026.
- [2] Fang, Richard, Rohan Bindu, Akul Gupta, and Daniel Kang. "[LLM Agents can Autonomously Exploit One-day Vulnerabilities](#)." arXiv:2404.08144, April 2024.
- [3] Anthropic. "[Project Glasswing: Securing critical software for the AI era](#)." Anthropic, April 2026.
- [4] DARPA. "[AI Cyber Challenge marks pivotal inflection point for cyber defense](#)." DARPA News, 2025.
- [5] AppSec Santa / Security Research Group. "[AI Pentesting Agents 2026: 39+ Tools, Architecture Deep Dive](#)." AppSecSanta, 2026.
- [6] Schneier, Bruce. "[On Anthropic's Mythos Preview and Project Glasswing](#)." Schneier on Security, April 2026.
- [7] The Next Web. "[Intruder launches AI pentesting agents as GCHQ-backed startup automates \\$50K manual security tests](#)." The Next Web, 2026.
- [8] Hadrian. "[The AI Hacking Boom: What 70 New Offensive Security Tools Mean for Defenders](#)." Hadrian Blog, 2026.
- [9] Hadrian. "[Understanding negative time-to-exploit in 2025](#)." Hadrian Blog, 2025.
- [10] Saptang Labs. "[Time-to-Exploit Shrinks: Patches Can't Keep Pace](#)." Saptang Labs Blog, 2025.
- [11] The Hacker News. "[Google: Over 57 Nation-State Threat Groups Using AI for Cyber Operations](#)." The Hacker News, January 2025.
- [12] Google Cloud. "[Adversaries Leverage AI for Vulnerability Exploitation, Augmented Operations, and Initial Access](#)." Google Cloud Blog, May 2026.
- [13] Help Net Security. "[Google researchers uncover criminal zero-day exploit likely built with AI](#)." Help Net Security, May 2026.
- [14] CISA. "[Known Exploited Vulnerabilities Catalog](#)." CISA, continuously updated.
- [15] Cloud Security Alliance. "[Agentic AI Threat Modeling Framework: MAESTRO](#)." CSA Blog, February 2025.

[16] Picus Security. "[The Glasswing Paradox: The Thing That Can Break Everything Is Also The Thing That Fixes Everything.](#)" Picus Security Blog, April 2026.