

First AI-Built Zero-Day: Autonomous Exploit Creation in the Wild

Security Implications of Google's GTIG Discovery and Guidance for Security Teams

2026-05-18

 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 May 11, 2026, Google's Threat Intelligence Group (GTIG) documented what it assesses as the first zero-day exploit in the wild that was developed with the assistance of a large language model (LLM), marking a significant inflection point in AI-enabled offensive operations.
- The exploit targeted a two-factor authentication (2FA) bypass vulnerability in a widely used open-source web-based system administration tool. A criminal threat actor had staged it for mass exploitation before Google intervened with responsible disclosure, enabling a patch before the campaign launched.
- LLM authorship was inferred—not confirmed through attribution—based on forensic artifacts in the exploit code: educational docstrings, a hallucinated CVSS score, and a structured, textbook-style Python format characteristic of LLM training outputs.
- The same GTIG report documented concurrent use of multi-agent penetration testing frameworks (Strix) and persistent temporal knowledge graphs (Graphiti) by suspected state-linked actors, and separately identified PromptSpy, an Android malware that autonomously queries Google's Gemini API to direct device interactions in real time.
- Security teams should treat AI-assisted vulnerability research as an immediate, operational threat rather than a future risk—and review authentication logic, patch cadence, and detection capabilities for AI-generated exploit indicators accordingly.

Background

AI capabilities have increasingly intersected with offensive security research over the past several years, as practitioners and adversaries alike have experimented with language models for code generation, vulnerability analysis, and the drafting of proof-of-concept exploits in controlled settings. The question of when—and whether—a threat actor would successfully deploy AI tooling to discover and weaponize a genuine zero-day in a production environment has been actively debated in the security community. On May 11, 2026, GTIG provided the first documented case that may answer it, publishing a report titled "Adversaries Leverage AI for Vulnerability Exploitation, Augmented Operations, and Initial Access." [1]

The report drew on GTIG's continuous monitoring of global threat activity and represented an escalation of findings first flagged in GTIG's February 2026 report on AI-related adversarial behavior. [2] By May, the group had observed what it characterized as a maturing transition from experimental AI integration into genuine industrial-scale application of generative models within adversarial workflows. [13][14] The zero-day discovery is the most prominent single data point in that trend, but the broader picture in the report is equally significant: state-linked actors deploying agentic frameworks for automated network reconnaissance, and malware operators embedding LLMs into mobile payloads for real-time device control.

Google withheld the name of the affected open-source web administration tool, ostensibly to allow time for patching before broader exposure. The vulnerability class—a semantic logic error that hardcoded a trust assumption contradicting the application's own authentication enforcement—is not inherently exotic. What sets this incident apart is the assessment that an AI model, rather than a human researcher, identified and weaponized the flaw first.

Security Analysis

The Zero-Day and Its Forensic Signature

The exploit at the center of the GTIG report was a Python script designed to bypass two-factor authentication on a popular open-source, web-based system administration tool. The underlying vulnerability arose from a semantic logic error: a developer had hardcoded an assumption that a particular code path could be trusted, in a way that was inconsistent with how the application's authentication layer was supposed to enforce verification. This class of error—a semantic mismatch between design intent and implementation—is difficult for automated static analysis to detect reliably, as it requires understanding contextual assumptions about control flow rather than pattern-matching syntax. [3]

What GTIG found most significant was the structure of the exploit itself. The Python code contained an abundance of educational docstrings—explanatory comments of the kind commonly inserted by language models to make generated code more legible—as well as a CVSS severity score that had been hallucinated: it did not correspond to any registered identifier in the CVE database. The overall format was, in GTIG's characterization, a "structured, textbook Pythonic" style highly characteristic of LLM training outputs. [1] GTIG assessed with high confidence that an AI model was used to support both the discovery and the weaponization of the vulnerability. The specific model was not identified, and Google stated it did not believe Gemini was involved. [4]

This forensic signature—educational comments, hallucinated metadata, and textbook formatting—represents an early-stage behavioral indicator for AI-assisted exploit development. It is unlikely to persist as threat actors grow more sophisticated and learn to strip these artifacts, but it provides defenders a narrow detection window that should be operationalized now. Security teams reviewing exploit code recovered from incidents or intelligence feeds should train analysts to recognize these characteristics.

The criminal actor behind the campaign had apparently acquired valid user credentials through prior means and planned the 2FA bypass as a second-stage capability enabling mass exploitation. GTIG's assessment is that its proactive counter-discovery may have disrupted the campaign before it achieved meaningful scale. The responsible disclosure process with the vendor proceeded quietly, and a patch was deployed prior to broad exploitation. [1][5]

Agentic Attack Frameworks in Active Use

The zero-day case was not the only AI-enabled threat pattern documented in the May 2026 GTIG report. Separately, analysts identified a suspected PRC-nexus threat actor deploying a multi-agent penetration testing framework called Strix against a Japanese technology firm and a prominent East Asian cybersecurity platform. This actor also deployed Hexstrike AI alongside Graphiti, a temporal knowledge graph system that maintains a persistent, evolving model of an attack surface—enabling the agentic toolchain to autonomously track reconnaissance findings over time, pivot between tools such as subdomain enumerators and HTTP probers based on internal reasoning, and resume campaigns without human re-initialization. [1][6]

GTIG characterized this operational pattern as a significant escalation of AI-assisted attack capability, noting that persistent state management combined with autonomous tool orchestration marks a new phase in adversarial AI deployment. Conventional penetration testing workflows depend on human operators to synthesize reconnaissance findings and direct tool selection, a step that Strix and Graphiti largely offload to the AI layer, enabling continuous, largely autonomous attack operations that can scale across targets with minimal human oversight. The use of a temporal knowledge graph in particular suggests that the actor was engineering for persistence and resumability—building attack infrastructure that can survive operational pauses and adapt to a changing target environment.

PromptSpy: Autonomous Malware with LLM-Directed Execution

A related but distinct threat documented in conjunction with the GTIG findings is PromptSpy, an Android malware family first identified by ESET in February 2026. [7] PromptSpy contains an autonomous agent module named GeminiAutomationAgent, which serializes the visible user interface hierarchy of the compromised device into an XML-like format via Android's Accessibility API and submits it to Google's

Gemini model. The model returns structured JSON responses specifying action types and screen coordinates, which PromptSpy then uses to simulate physical gestures—clicks, swipes, and navigation—without human involvement. [8]

The significance of this architecture is that it inverts the traditional model of mobile malware. Rather than hardcoding a sequence of actions that may fail when application layouts change, PromptSpy queries an external reasoning model at runtime to determine what to do next given the current screen state. This design makes the malware more resilient to UI changes than hardcoded-sequence approaches, since the model can adapt its action selection to the current screen state. The broader capabilities documented by ESET include capturing lockscreen data, blocking uninstallation, taking screenshots, recording screen activity as video, and capturing biometric replay data for authentication gestures. [9]

While PromptSpy was identified separately from the zero-day campaign, the GTIG report explicitly connected it to the broader pattern of AI integration into offensive tooling. Together, these findings document a threat environment in which AI models are being used not merely for content generation or planning assistance, but as real-time execution engines embedded directly in attack infrastructure.

Implications for the Threat Landscape

The GTIG findings carry several implications that extend beyond the specific incidents documented. First, the barrier to zero-day discovery has changed. Identifying a semantic logic error in a complex authentication flow historically required a skilled human researcher with deep contextual knowledge of the application. If LLMs can perform this analysis at scale across many targets, the rate of novel vulnerability discovery—on both offensive and defensive sides—could increase significantly, shortening the window between disclosure and exploitation. The CSA AI Vulnerability Storm report, published in April 2026, modeled AI-accelerated vulnerability discovery as a near-term threat scenario and provided a prioritized response framework directly applicable to this case, warning that organizations should prepare for a step-change in vulnerability disclosure volume driven by AI-assisted research. [10]

Second, the forensic indicators associated with AI-assisted exploits are currently useful but will degrade. Hallucinated CVSS scores and educational docstrings are artifacts of current LLM behavior, not permanent features. As threat actors gain experience and as AI models improve, the distinctive textbook formatting noted by GTIG will become less reliable as a detection signal. The window for leveraging these behavioral indicators is finite, and defenders should use it while it remains open.

Third, the combination of agentic frameworks, temporal knowledge graphs, and autonomous execution modules (as seen in the Strix, Graphiti, and PromptSpy cases) suggests that the architecture of offensive operations is shifting toward persistent, self-directed campaigns that require less human management. These findings suggest that automation may change the economics of targeted attacks—

potentially lowering per-target cost and expanding viable target sets as the human overhead per operation decreases—though empirical data on operational cost structures for either state-linked or criminal actors remains limited.

Recommendations

Immediate Actions

Organizations should prioritize patching the affected web administration tool once vendor identification becomes public, and should verify that 2FA implementations across all administrative interfaces rely on server-side enforcement rather than client-side or hardcoded trust assumptions. Authentication logic audits should specifically examine whether any code paths bypass verification based on implicit rather than explicit trust grants—the class of error exploited in this incident. Environments using open-source web administration panels should be reviewed as a precautionary measure pending full public disclosure of the affected product.

Security operations teams should brief analysts on the forensic characteristics of AI-generated exploit code documented by GTIG: educational docstrings, hallucinated CVE or CVSS metadata, and textbook-style formatting. These indicators should be incorporated into malware triage and incident analysis workflows immediately, with the expectation that their reliability will diminish as threat actors adapt.

Short-Term Mitigations

Authentication hardening should proceed in parallel with patching. Multi-factor authentication for administrative interfaces should be implemented at the network or infrastructure layer—not solely within the application—so that application-level logic errors cannot bypass it entirely. Privileged access workstations, just-in-time access controls, and network segmentation for administrative tools reduce the exposure window even when application-layer authentication can be circumvented.

Detection coverage for agentic attack tooling should be reviewed. The Strix and Hexstrike frameworks documented by GTIG generate behavioral patterns consistent with automated reconnaissance: high-frequency subdomain enumeration, HTTP probing at scale, and tool-switching based on intermediate results. SIEM and NDR tuning should account for these patterns as indicators of automated adversarial activity rather than human-directed scanning.

For mobile device management programs, PromptSpy's architecture illustrates a new detection challenge. The malware's use of the Accessibility API for UI interaction is not inherently malicious—many legitimate accessibility and automation tools use the same API—but the combination of Accessibility API usage, network calls to LLM inference endpoints, and unusual gesture simulation patterns constitutes a novel behavioral detection signature that mobile threat intelligence programs should incorporate into their detection models.

Strategic Considerations

The emergence of AI-assisted zero-day development reinforces the case for vulnerability management programs that operate on shortened disclosure-to-patch timelines. If LLMs can accelerate the cycle from vulnerability discovery to weaponization, the traditional assumption that organizations have weeks or months between private discovery and active exploitation may no longer hold—making accelerated patch cadences an operational necessity rather than a best practice. Patch SLAs, particularly for internet-facing administrative tooling, should be reviewed and tightened accordingly.

Security teams should also prepare for the possibility that AI-assisted vulnerability discovery will increase not just attack-side findings but also the volume of researcher-reported vulnerabilities. Bug bounty programs, coordinated disclosure pipelines, and patch engineering capacity may all need to scale to absorb a higher steady-state volume of valid vulnerability reports. The CSA AI Vulnerability Storm framework recommended that organizations treat this as a structural capacity question, not a temporary surge. [10]

Finally, the PromptSpy case illustrates that AI-enabled capabilities are appearing in criminal and espionage malware simultaneously, not sequentially. Security architecture should not distinguish between "AI threats" as a future category and current operational threats. The architecture changes needed—stronger authentication enforcement, behavioral detection tuned for autonomous tooling, and reduced reliance on application-layer controls for critical administrative surfaces—are improvements that address current risks, not hypothetical ones.

CSA Resource Alignment

The incidents documented by Google's GTIG map directly to threat patterns addressed by several CSA frameworks and publications.

The **MAESTRO framework** for agentic AI threat modeling provides the most directly applicable lens. The Strix and Hexstrike deployments described in the GTIG report instantiate threats at MAESTRO Layer 6 (Multi-Agent Orchestration), where autonomous agents coordinate across tools to pursue attack objectives with minimal human direction. The PromptSpy GeminiAutomationAgent module exemplifies threats at Layer 3 (Agent Frameworks) and Layer 5 (Agent Trust Boundaries), as the malware exploits the trust that Android's Accessibility API extends to registered applications and leverages an external LLM as a privileged reasoning component.

The **CSA AI Controls Matrix (AICM) v1.0.3** [12] addresses controls applicable across multiple stakeholder roles. Application Providers should review authentication enforcement controls, particularly those addressing logical consistency between code paths. Orchestrated Service Providers operating multi-agent systems should ensure that agent scope limitations and logging requirements are in place to detect autonomous lateral movement of the kind described in the GTIG report. The AICM's controls on model output validation are also directly relevant to the PromptSpy case, where structured model outputs drive physical device interaction.

The **CSA Agentic AI Red Teaming Guide** [11] provides testing methodologies that security teams should apply to any multi-agent or LLM-integrated tooling in their environments. The tactics documented by GTIG—autonomous tool switching, persistent state management via knowledge graphs, real-time LLM-directed execution—are addressable through the adversarial testing approaches outlined in the guide, which covers privilege escalation, trust boundary violations, and agent orchestration abuse.

The **CSA AI Vulnerability Storm** strategy briefing [10], published in April 2026, anticipated the threat scenario documented by GTIG and provides a prioritized action framework for security leaders managing accelerated vulnerability discovery timelines. Its recommendations on patch cadence, vulnerability management capacity, and board-level communication are directly applicable to the organizational response this incident demands.

Security leaders should use these incidents to drive concrete conversations within their organizations about which CSA controls they have implemented, which remain aspirational, and where the gaps are most consequential given the current threat environment.

References

- [1] Google Threat Intelligence Group. "[Adversaries Leverage AI for Vulnerability Exploitation, Augmented Operations, and Initial Access.](#)" Google Cloud Blog, May 11, 2026.
- [2] Google Threat Intelligence Group. "[GTIG AI Threat Tracker: Distillation, Experimentation, and \(Continued\) Integration of AI for Adversarial Use.](#)" Google Cloud Blog, February 2026.
- [3] The Hacker News. "[Hackers Used AI to Develop First Known Zero-Day 2FA Bypass for Mass Exploitation.](#)" The Hacker News, May 2026.
- [4] BleepingComputer. "[Google: Hackers Used AI to Develop Zero-Day Exploit for Web Admin Tool.](#)" BleepingComputer, May 2026.
- [5] SC Media. "[Google Reports First Known AI-Assisted Zero-Day Exploit in the Wild.](#)" SC Media, May 2026.
- [6] SecurityWeek. "[Google Detects First AI-Generated Zero-Day Exploit.](#)" SecurityWeek, May 2026.
- [7] ESET Research. "[PromptSpy Ushers in the Era of Android Threats Using GenAI.](#)" WeLiveSecurity, February 2026.
- [8] The Hacker News. "[PromptSpy Android Malware Abuses Google Gemini AI to Automate Recent-Apps Persistence.](#)" The Hacker News, February 2026.
- [9] SecurityWeek. "[PromptSpy Android Malware Abuses Gemini AI at Runtime for Persistence.](#)" SecurityWeek, February 2026.
- [10] Evron, Gadi; Mogull, Rich; Lee, Robert T. "[The AI Vulnerability Storm: Building a Mythos-Ready Security Program.](#)" Cloud Security Alliance, April 2026.
- [11] Cloud Security Alliance. "[Agentic AI Red Teaming Guide.](#)" CSA, 2025.
- [12] Cloud Security Alliance. "[AI Controls Matrix \(AICM\) v1.0.3.](#)" CSA, 2026.
- [13] Help Net Security. "[Google Researchers Uncover Criminal Zero-Day Exploit Likely Built with AI.](#)" Help Net Security, May 11, 2026.
- [14] The Register. "[Google Says Criminals Used AI-Built Zero-Day in Planned Mass Hack Spree.](#)" The Register, May 11, 2026.