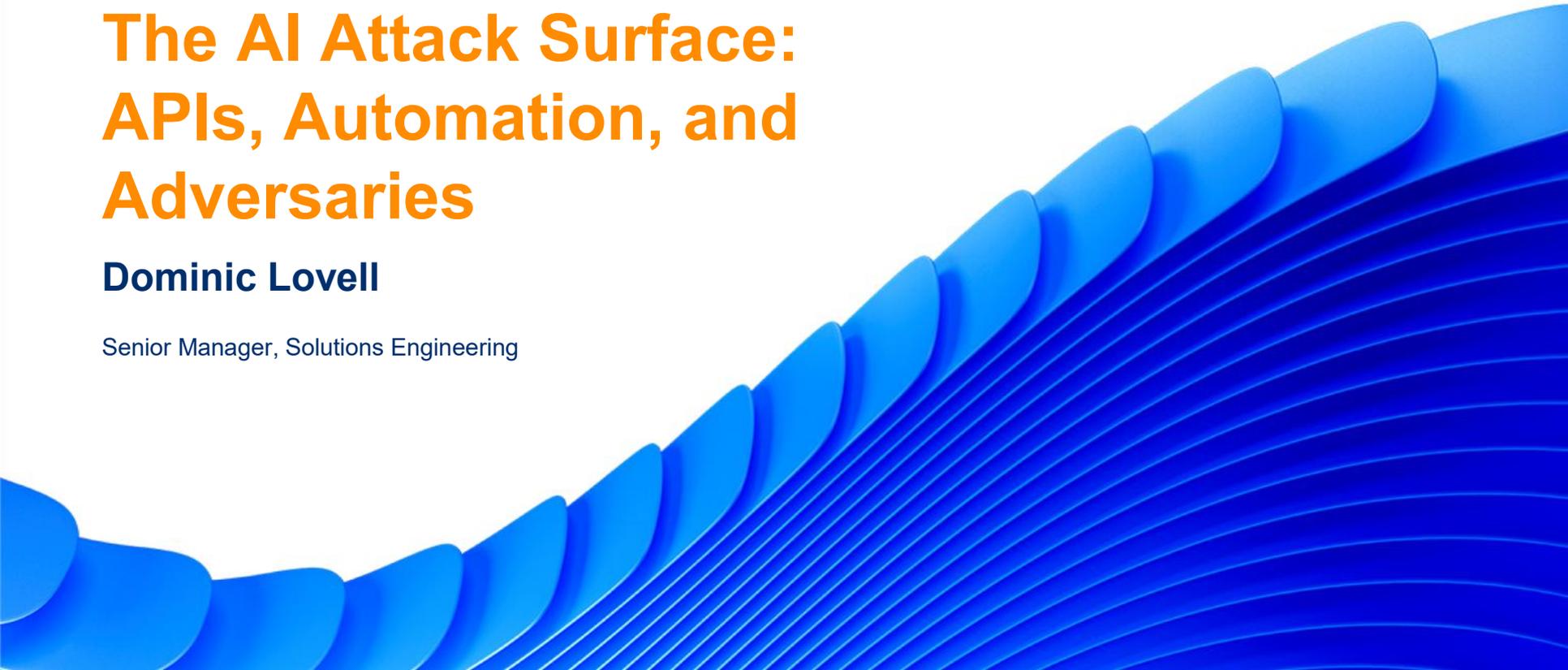# The AI Attack Surface: APIs, Automation, and Adversaries

**Dominic Lovell**

Senior Manager, Solutions Engineering

# **Akamai** is the
# world's most distributed cloud platform,
# with leading solutions for:

| Content Delivery | Cyber Security | Cloud Computing |
|---|---|---|

Akamai

**Through massive distribution, full automation, and network intelligence, Akamai provides:**

3

| | |
|---|---|
| **14T**<br>Bot requests | **1,000**<br>TBPS of capacity |
| **43B**<br>WAF attacks | **4,000+**<br>Edge PoPs |
| **1.5**<br>Petabytes of DDoS traffic mitigated | **1,200+**<br>Networks |
| **9.6T**<br>L7 DDoS requests directed at customers | **750+**<br>Cities |
| **17.9B**<br>Fraudulent login attempts | **130**<br>Countries |

# 1061

## Average number of applications per enterprise
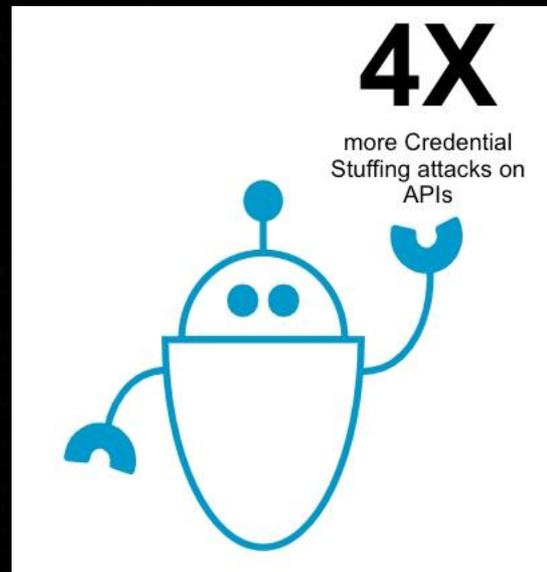
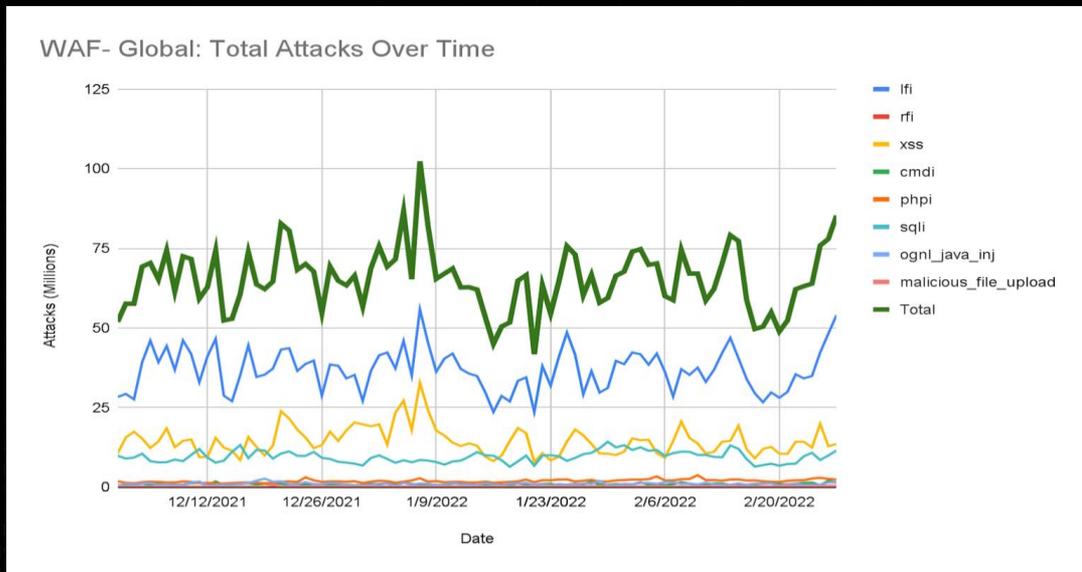Including mission-critical applications

# Only 1 in 3

Breaches were identified by an organization's own security team

(IBM/Ponemon Institute)

# Agentic APIs: A Primary Target For Attackers Today



WAF- Global: Total Attacks Over Time

Legend:
- lfi
- rfi
- xss
- cmdi
- phpi
- sqli
- ognl_java_inj
- malicious_file_upload
- Total



**4X** more Credential Stuffing attacks on APIs

Web sites & Web APIs share the same (old) attack vectors

*– but APIs are often unprotected*

APIs are more performant *and less expensive* to attack compared with traditional web forms

# AI crawling is the new data breach

- Scraping is not benign

- Low and slow

- Hard to detect

- Requires behavioral monitoring

Akamai

# Let's frame the problem

# The mental model is wrong

**What most people think**

- AI risk = hallucinations

- AI risk = prompt injection

- AI risk = data privacy

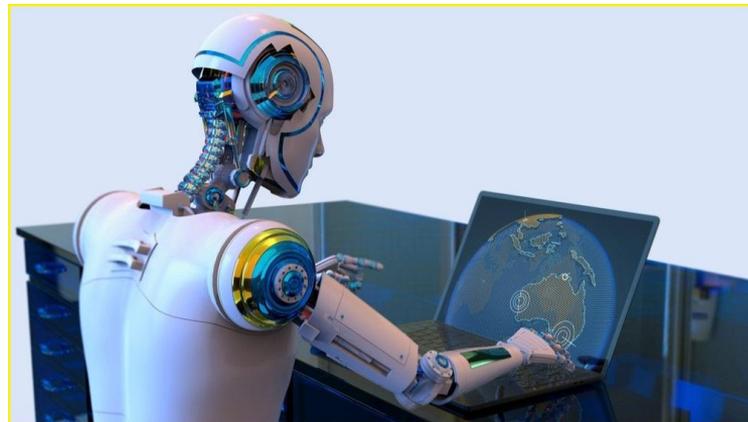**What's actually happening**

AI systems are:

- Autonomous

- Always-on

- Deeply integrated

- Privileged by design

AI didn't add a new attack surface… it connected all the old ones.

# Why AI changes security fundamentals

**Three shifts introduced by AI**

1. Automation: attacks run continuously

1. Agency: systems act, not just respond

1. Abstraction: humans stop seeing the plumbing



This is where the risk hides

# Real world AI system architecture

- Ai systems are connected to:

  ○ APIs

  ○ Data sources

  ○ SaaS tools

  ○ Internal systems

  ○ External services

  ○ Automation tools



Every AI system is a decision engine sitting on top of privileged access.

# AI failure mode: Over-trust

Some AI systems are trusted because

- They're "internal"

- They're "read-only" (until they aren't)

- They use valid credentials

- They improve workflows



AI behaviour is unsafe and inconsistent

# The "looks legitimate problem

**Why AI-driven abuse is hard to spot**

- Valid credentials

- Normal request shapes

- Human-like timing

- No spikes - can be human insititated



This is why detection fails.

# Autonomous Reach

**AI agents don't stop**

- They don't "log out"
- They don't forget
- They don't get bored

**They:**

- Crawl
- Correlate
- Retry
- Learn



This creates persistent exposure

# Risk: Tool amplification & sprawl

**AI agents don't act alone**

- They call tools

- They chain actions

- They delegate tasks



Every step in the workflow exposes data and expands blast radius.

# Risk: Supply Chain + Automation

**Real world attacks:**

- Supply-chain compromise

- Credentials harvested

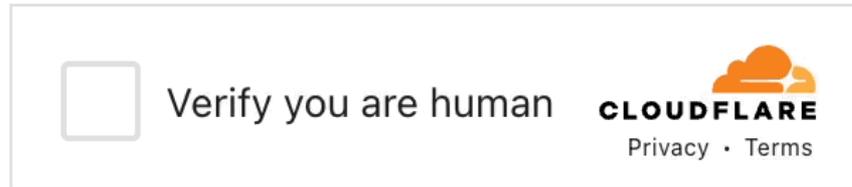- Internal access abused

**AI accelerates:**

- Credential testing

- Lateral movement

- Privilege mapping

# Why traditional controls don't work

## Controls built for humans

- Rate limits
- CAPTCHAs
- Alerts
- MFA



## AI bypasses by

- Distribution
- Patience
- Adaptation
- Validity

# Behavioural analysis in the only warning signs

**AI abuse shows up as:**

- Unusual sequences
- Consistent curiosity
- Systematic traversal
- Long-lived access

**It won't show up as:**

- Not spikes
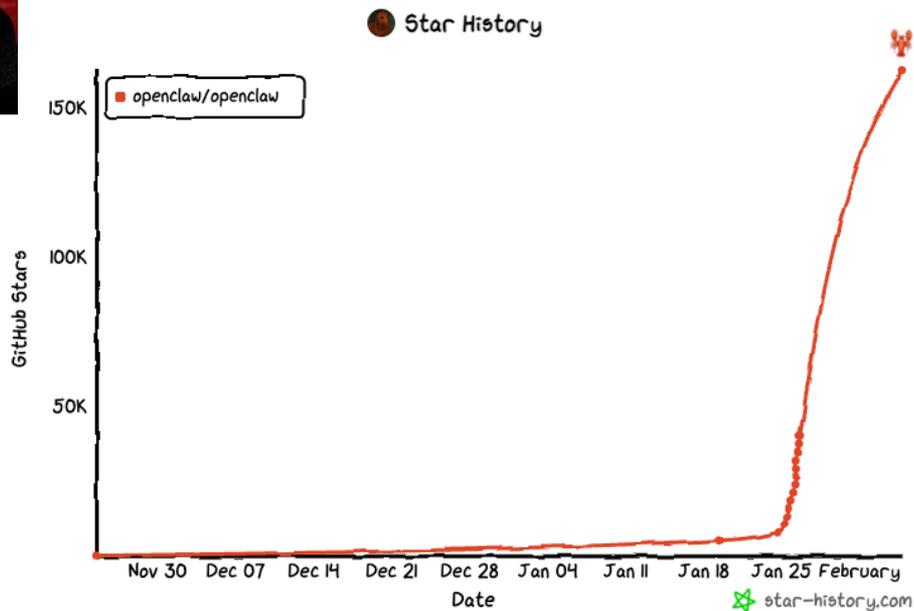- Alerts
- Traditional bot patterns

# What this means for you

- Assume AI agents behave like insiders

- Treat automation as privileged

- Monitor behaviour, not just access

- Expect persistence, not attacks

# Real world examples

Star History

openclaw/openclaw

← Blog

# Hacking Moltbook: The AI Social Network Any Human Can Control

1 exposed database. 35,000 emails. 1.5M API keys. And 17,000 humans behind the not-so-autonomous AI network.

Listen to the "Crying out Cloud" podcast episode

**Gal Nagli**
February 2, 2026          9 minute read

```
Table                    Records      Sensitive Data
agents                   1,494,823    API keys, claim tokens, verification codes
votes                    2,661,805    User voting behavior
comments                 232,813      User content
notifications            221,892      Private user alerts
follows                  56,815       Social graph data
posts                    50,156       Full post content
observers                29,631       Email addresses (newsletter subscribers)
owners                   17,008       Emails, Twitter handles, real names
submolts                 13,725       Community data
agent_messages           4,060        Private direct messages
identity_verifications   25           IP addresses, user agents
site_admins              1            Administrator identity

Total: ~4.75 million records exposed
```

Akamai

← Blog

# Wiz Research Uncovers Exposed DeepSeek Database Leaking Sensitive Information, Including Chat History

A publicly accessible database belonging to DeepSeek allowed full control over database operations, including the ability to access internal data. The exposure includes over a million lines of log streams with highly sensitive information.

**Gal Nagli**
January 29, 2025    3 minute read

Akamai

# "Everyone is a vide-coder now"

SIGN IN

**Malwarebytes** LABS

Personal   Business   Pricing   Partners   Resources   Help   **FREE DOWNLOAD**

AI, BUGS, DATA BREACHES, NEWS

# AI chat app leak exposes 300 million messages tied to 25 million users

by Pieter Arntz | February 9, 2026

An independent security researcher uncovered a major data breach affecting Chat & Ask AI, one of the most popular AI chat apps on Google Play and Apple App Store, with more than 50 million users.

The researcher claims to have accessed 300 million messages from over 25 million users due to an exposed database. These messages reportedly included, among other things, discussions of illegal activities and requests for suicide assistance.

Behind the scenes, Chat & Ask AI is a "wrapper" app that plugs into various large language models (LLMs) from other companies, including OpenAI's ChatGPT, Anthropic's Claude, and Google's Gemini. Users can choose which model they want to interact with.

The exposed data included user files containing their entire chat history, the models used, and other settings. But it also revealed data belonging to users of other apps developed by Codeway—the developer of Chat & Ask AI.

The vulnerability behind this data breach is a well-known and documented Firebase misconfiguration.

## Other Articles tagged News

Discord will limit profiles to teen-appropriate mode until you verify your age →
February 10, 2026

25

# BUSINESS INSIDER

AI

# An AI agent spent 16 hours hacking Stanford's network. It outperformed human pros for much less than their 6-figure salaries.

By **Lee Chong Ming**    ( + Follow )

- An AI agent hacked Stanford's computer science networks for 16 hours in a new study.

- The AI agent outperformed nine out of 10 human participants, said the study by Stanford researchers.

- It also cost a fraction of the six-figure salary for a "professional penetration tester."

For 16 hours, an AI agent crawled Stanford's public and private computer science networks, digging up security flaws across thousands of devices.

26

Policy

# Disrupting the first reported AI-orchestrated cyber espionage campaign

13 Nov 2025

Read the report

In mid-September 2025, we detected suspicious activity that later investigation determined to be a highly sophisticated espionage campaign. The attackers used AI's "agentic" capabilities to an unprecedented degree—using AI not just as an advisor, but to execute the cyberattacks themselves.

The threat actor—whom we assess with high confidence was a Chinese state-sponsored group—manipulated our Claude Code tool into attempting infiltration into roughly thirty global targets and succeeded in a small number of cases. The operation targeted large tech companies, financial institutions, chemical manufacturing companies, and government agencies. We believe this is the first documented case of a large-scale cyberattack executed without substantial human intervention.

27

# Shadow AI practices: A wakeup call for enterprises

Opinion

Feb 10, 2026 · 7 mins

While executives talk AI strategy, shadow agents are already inside the enterprise, quietly rewriting your risk profile faster than policies can keep up.

---

SEARCH

*THOUGHT LEADERS*

# Forget Shadow AI Panic: Sprawl Is Here to Stay

SUPERCHARGE

By **Adam Magill**, SVP, Global Security (CISO) at Concentrix
Published February 10, 2026

29

# The rise of workers using 'shadow AI' to do their jobs

**Julie Hare** and **Tess Bennett**

Aug 14, 2025 – 12.00am

Listen to this article
4 min

Save     Share

Gift this article

A grassroots movement of workers is using generative AI, particularly in small to medium enterprises, but often without the formal knowledge of their bosses for fear of being labelled lazy or less competent, a major new report has found.

The "shadow adoption" of AI showed that workers were driving "bottom-up innovation", potentially generating initiative and experimentation but also shifting governance and risk management on to individual workers, Jobs and Skills Australia said.



A grassroots revolution is under way as employees take it upon themselves to use AI. **Bethany Rae**

30

**TC TechCrunch**

Latest   Startups   Venture   Apple   Security   AI   Apps   |   Events   Podcasts   Newsletters



IMAGE CREDITS: HEATHER DIEHL / GETTY IMAGES

Zack Whittaker

# Trump's acting cybersecurity chief uploaded sensitive government docs to ChatGPT

The acting head of U.S. cybersecurity agency CISA uploaded sensitive contracting documents marked "for official use only" to ChatGPT, according to Politico.

The outlet, citing officials, reported Tuesday that CISA's acting director,

Akamai

GitGuardian BLOG   WEBSITE   INTERACTIVE DEMO   LEARNING CENTER   NHI   SEARCH   in ▶ ○   **Book a demo**

BREACH EXPLAINED

# The Secret's Out: How Stolen Okta Auth Tokens Led to Cloudflare Breach

Cloudflare experienced a security breach when its internal systems were compromised, leading to unauthorized access to sensitive data. Another incident highlights the importance of maintaining strict secrets security across the supply chain.

**THOMAS SEGURA, DWAYNE MCDANIEL**
2 FEB 2024 · 6 MIN READ

Follow us on   in  ✱

## What Happened?

Cloudflare's internal Atlassian systems were breached using tokens and service accounts compromised from a previous Okta breach. The attackers gained access to the Confluence wiki, Jira database, and Bitbucket source code system. The incident illustrates the damaging domino effect of secrets sprawl and the importance of maintaining rigorous secrets security across the supply chain.

Akamai

# Okta Data Breach: What Happened, Impact, and Security Lessons Learned

**The Nightfall Team** • May 13, 2024 • 1 min read

## How the Attack Occurred: Attack Vectors and Techniques

The Okta breach involved a sophisticated attack targeting the company's customer support systems rather than its core identity platform. According to Okta's investigation, the attackers gained access to a service account within Okta's support system. This access allowed them to view and download HAR files that customers had submitted as part of support requests.

HAR files are particularly sensitive because they archive HTTP transactions, often containing session tokens, cookies, and other authentication data. Security researchers determined that the attackers specifically searched for these files to extract valid session tokens. Once obtained, these tokens allowed the threat actors to impersonate legitimate users without needing to know passwords or bypass multi-factor authentication.

The attack demonstrates the concept of "living off the land," where attackers use legitimate credentials and tools to avoid detection. Rather than exploiting a technical vulnerability in Okta's systems, the attackers exploited access management gaps in support workflows, highlighting how sophisticated threat actors often target the path of least resistance.

kamai

# The Hacker News

Home    Data Breaches    Cyber Attacks    Vulnerabilities    Webinars    Expert Insights    Contact
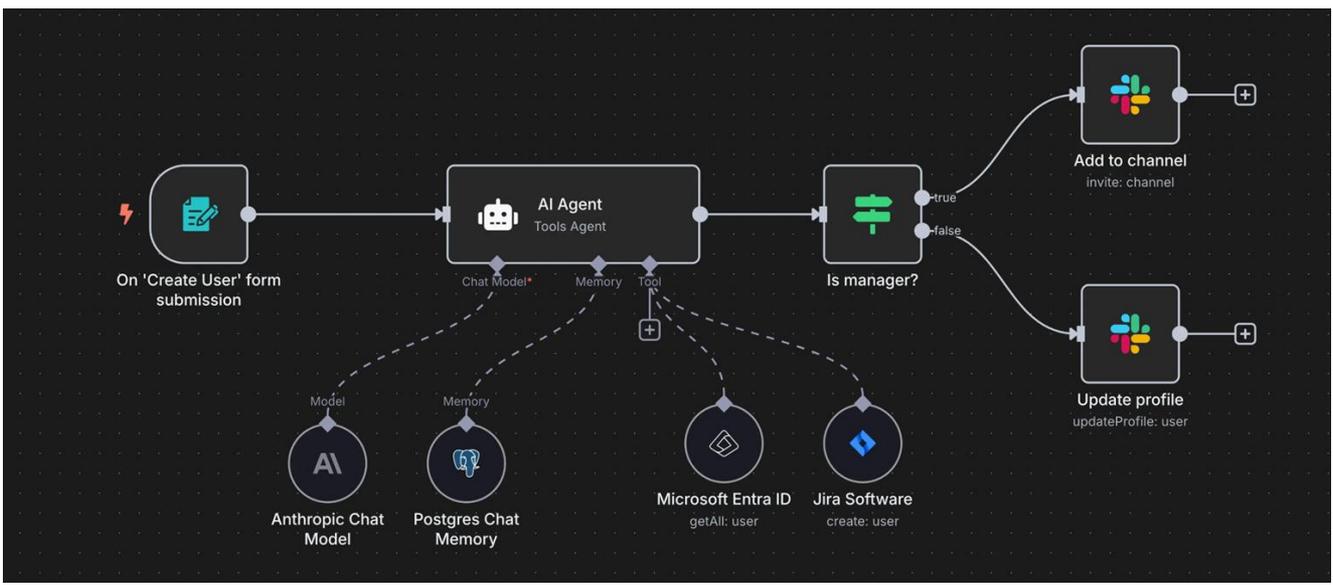
## n8n Supply Chain Attack Abuses Community Nodes to Steal OAuth Tokens

👤 Ravie Lakshmanan    📅 Jan 12, 2026

Vulnerability / Workflow Automation

**SECURITYWEEK**
CYBERSECURITY NEWS, INSIGHTS & ANALYSIS

Malware & Threats ⌄    Security Operations ⌄    Security Architecture ⌄    Risk Management ⌄    CISO Strategy ⌄    ICS/OT ⌄    Funding/M&A ⌄    Cyber AI

**SUPPLY CHAIN SECURITY**

# Hackers Target Popular Nx Build System in First AI-Weaponized Supply Chain Attack

"To our knowledge, this is one of the first documented cases of malware coercing AI-assistant CLIs to assist in reconnaissance.

"This technique forces the AI tools to recursively scan the file system and write discovered sensitive file paths to /tmp/inventory.txt, effectively using legitimate tools as accomplices in the attack."

systems when opening new terminal sessions, GitGuardian explains.

Additionally, the code was designed to weaponize AI tools such as Claude and Gemini to help with reconnaissance and data exfiltration.

"This marks the first known case where attackers have turned developer AI assistants into tools for supply chain exploitation," StepSecurity points out.

38

# Hidden npm Malware Exposes New Supply Chain Weakness

PhantomRaven's success hinged on npm's built-in lifecycle scripts.

The malicious dependency contained a preinstall hook — *"preinstall": "node index.js"* — that executed automatically without user consent. This meant even deeply nested dependencies could trigger execution as part of a normal installation process.

Once active, the malware systematically harvested data from the developer's system including:

▸ **Email addresses** from *.gitconfig*, *.npmrc*, and environment variables.

▸ **CI/CD credentials**, including GitHub Actions tokens, GitLab CI keys, Jenkins, CircleCI, and npm publishing tokens.

▸ **System fingerprinting** data, such as IP addresses, hostnames, OS details, and usernames.

The exfiltrated data was redundantly transmitted via HTTP GET, HTTP POST, and WebSocket connections, ensuring delivery even under network restrictions.

https://www.esecurityplanet.com/news/hidden-npm-malware-supply-chain/

Akamai

# Hidden npm Malware Exposes New Supply Chain Weakness

## The Rise of Slopsquatting Attacks

Beyond the stealthy delivery mechanism, PhantomRaven introduced a novel social-engineering tactic called slopsquatting — a twist on traditional typosquatting.

Instead of mimicking existing package names, attackers registered plausible-sounding names that AI assistants like GitHub Copilot or ChatGPT might hallucinate on.

Examples include:

▸ *eslint-comments* instead of the legitimate *eslint-plugin-eslint-comments*

▸ *unused-imports* instead of *eslint-plugin-unused-imports*

▸ *Transform-react-remove-prop-types* instead of *babel-plugin-transform-react-remove-prop-types*

When AI suggested these nonexistent packages to developers, users unknowingly installed the malicious versions — demonstrating how artificial intelligence can unintentionally amplify supply-chain risks.

https://www.esecurityplanet.com/news/hidden-npm-malware-supply-chain/

*Akamai*

# Our approach: Behavioural analysis

# Example: Suspicious New Account Creation



1. Path Parameter Fuzzing alert

2. Abnormal fuzzing of {customer_id} path parameter

3. Correlated to suspicious new account creations

# Key AI Security Challenges

## Shadow AI

Discover your complete AI footprint - including rogue, legacy, shadow, zombie, etc.

## Vulnerable AI tools

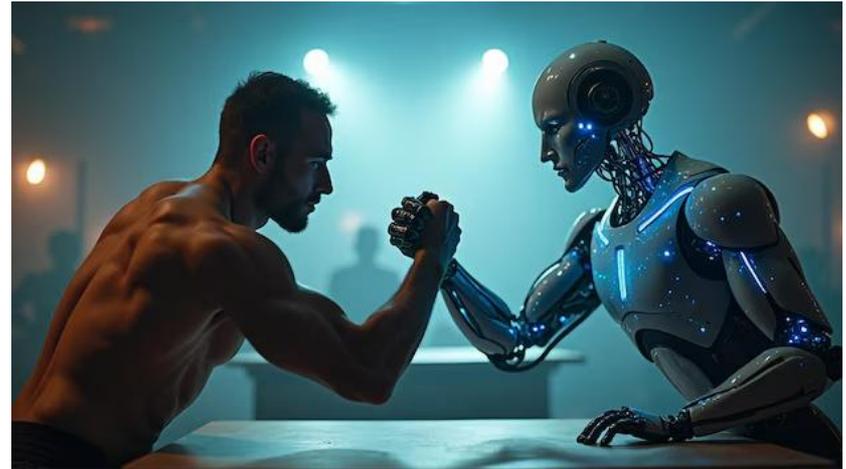Prevent agentic vulnerabilities and misconfigurations

## Agent Abuse

Stop business logic abuse such as data scraping or data exfiltration using behavioral analytics.

**Tomorrow's Focus**

**Today's Focus**

Akamai

# What this means for you

- Assume AI agents behave like insiders

- Treat automation as privileged

- Monitor behaviour, not just access

- Expect persistence, not attacks

**Thank you**

45