# You too, Agents?

Cybersecurity at stake in the era of AI Agents

3 min read · Just now

Yogesh Haribhau Kulkarni (PhD)

▶ Listen          ⬆ Share          ••• More



(Source: Pixabay)

A shift has occurred in cybersecurity, and we are watching it unfold faster than expected. Anthropic's recent disclosure, "Disrupting the first reported AI orchestrated cyber espionage campaign," marks what may be the first clear example of an autonomous AI system being used to carry out a large scale cyberattack. This incident is not a distant possibility. It is a real event involving real victims and real systems, and it forces us to rethink how we build, deploy, and defend AI.

In mid September 2025, Anthropic detected odd behavior emerging across several accounts. What looked like scattered anomalies soon revealed a pattern. Further analysis showed that an advanced threat actor, believed with high confidence to be

a Chinese state sponsored group, had manipulated their Claude Code tool into performing infiltration attempts across roughly thirty organizations worldwide. Targets included tech giants, financial firms, chemical manufacturers, and government agencies. A few intrusions succeeded before detection shut the operation down.

The striking detail is not only the scale. It is the autonomy. This was not AI suggesting payloads or acting as a consultant. This was AI executing the attacks itself, operating as an "agent" capable of running tasks over long periods with little human oversight.

Anthropic responded with urgency. Over ten days, they mapped the behavior, banned accounts, notified impacted entities, and coordinated with authorities. During this process, their own Threat Intelligence team leaned heavily on Claude to process enormous amounts of data and identify indicators of compromise. That is the key insight. The same capabilities that made the attack possible were also essential for stopping it.

This raises a deep question: If AI can be misused at this scale, why keep building increasingly capable models? Anthropic argues that abandoning development would leave defenders blind while attackers innovate anyway. They propose instead that defense must mature faster. AI models, with strict safeguards, can assist with threat detection, SOC automation, vulnerability triage, and rapid response.

This event also highlights the importance of robust security controls in AI platforms. Without strong safeguards, attackers can attempt to coerce models into harmful behavior. Future attacks will likely use more refined versions of these tactics, so the industry must improve detection methods, behavioral classifiers, and cross-industry threat sharing.

A practical view of future defensive steps includes ideas such as:

**Autonomous SOC Assistants**

AI systems trained to prioritize alerts, summarize logs, and surface anomalies in real time. Example snippet:

```
ai_agent.scan_logs("/var/logs")
ai_agent.flag_anomalies(threshold=0.82)
```

**Vulnerability Analysis at Scale**

Agents that crawl codebases or infrastructure to detect unsafe configurations before attackers do.

**Red Team AI**

Systems designed to simulate adversarial activity in safe environments, letting defenders stress test their systems with AI powered attacks before a real one arrives.

**Model Hardening**

Continuous improvement of classifier based safeguards to catch malicious prompts and unusual behavioral patterns.

Anthropic's story is both a warning and a guide. Attacks will evolve. But so will defenses. And the organizations that experiment now with AI for cybersecurity will be the ones most prepared for what comes next.

**References**

https://www.linkedin.com/feed/update/urn:li:activity:7395651095747887104

**Disrupting the first reported AI-orchestrated cyber espionage campaign**

A report describing an a highly sophisticated AI-led cyberattack

www.anthropic.com