



Flaws in OpenClaw AI Agent Risk Data Exfiltration

Description

China's National Computer Network Emergency Response Technical Team (CNCERT) has issued a warning regarding potential security vulnerabilities associated with OpenClaw, an open-source, self-hosted autonomous artificial intelligence (AI) agent, previously known as Clawdbot and Moltbot. This alert underscores the importance of cybersecurity in the rapidly evolving AI landscape.

CNCERT explained that OpenClaw's weak default security settings could be exploited by malicious actors to gain control over its operations. The platform's design allows it to execute tasks autonomously, which, in conjunction with its privileged system access, poses significant risks. One major concern is prompt injection, where harmful instructions hidden in web content can lead the agent to disclose sensitive data.

This type of attack, known as indirect prompt injection (IDPI) or cross-domain prompt injection (XPJA), involves manipulating benign AI functions, such as web page summarisation, to execute malicious commands. Such tactics can allow adversaries to manipulate AI systems for various nefarious purposes, including evading content filters and generating biased responses.

OpenAI has acknowledged the evolution of these prompt injection threats, highlighting that as AI agents become capable of browsing the web and acting on users' behalf, they also present new vulnerabilities for exploitation.

The risks associated with OpenClaw have already manifested in research findings from PromptArmor, revealing that link preview features in messaging applications could be misused to exfiltrate data via indirect prompt injection.

CNCERT raised additional alarms about other potential issues, including inadvertent loss of critical information and the possibility of attackers uploading malicious capabilities to the platform. Such breaches could severely impact vital sectors like finance and energy, leading to significant business disruptions and data leaks.

In response, authorities recommend that users strengthen network controls, isolate the service, and download only from trusted sources. Concurrently, Chinese authorities are limiting the use of OpenClaw among state-run enterprises and military personnel to mitigate these security threats.

Comprehension Questions



Multiple Choice

1. What is the previous name of OpenClaw?

- Option: AI Guard
- Option: Clawdbot
- Option: CyberShield
- Option: DataSentry

2. What type of attack involves harmful instructions hidden in web content?

- Option: SQL Injection
- Option: Cross-Site Scripting
- Option: Prompt Injection
- Option: DDoS Attack

3. Which team issued a warning about OpenClaw?

- Option: Global Security Team
- Option: CNCERT
- Option: OpenAI Advisory Board
- Option: AI Safety Group

4. What is one suggested measure to improve security for OpenClaw?

- Option: Use default settings
- Option: Strengthen network controls
- Option: Allow all users access
- Option: Download from any sources

5. What is the risk associated with OpenClaw related to sensitive data?

- Option: Increased data storage
- Option: Data encryption issues
- Option: Data exfiltration
- Option: Virus transmission

6. Which sectors could be significantly impacted by breaches in OpenClaw?

- Option: Retail and Hospitality
- Option: Finance and Energy
- Option: Healthcare and Education
- Option: Entertainment and Media



True-False

7. OpenClaw is a paid software application.
8. CNCERT has raised alarms about OpenClaw due to its weak default security settings.
9. Prompt injection can lead to the disclosure of sensitive data.
10. Chinese authorities encourage the unrestricted use of OpenClaw in state-run enterprises.
11. Indirect prompt injection is a type of cyber attack.
12. CNCERT has recommended users to isolate the service for better security.

Gap-Fill

13. OpenClaw was previously known as _____ and Moltbot.
14. One major concern with OpenClaw is _____ injection.
15. The attack type involves manipulating benign AI functions such as web page _____
to execute malicious commands.
16. Chinese authorities are limiting the use of OpenClaw among state-run _____ to
mitigate security threats.
17. CNCERT noted the risk of inadvertent loss of critical _____ related to OpenClaw.
18. Authorities recommend that users only download OpenClaw from _____ sources.

Answer

Multiple Choice: 1. Clawdbot 2. Prompt Injection 3. CNCERT 4. Strengthen network controls 5. Data exfiltration 6. Finance and Energy

True-False: 7. False 8. True 9. True 10. False 11. True 12. True

Gap-Fill: 13. Clawdbot 14. prompt 15. summarisation 16. enterprises 17. information 18. trusted

CATEGORY



1. Business - LEVEL6

Date Created

2026/03/15

Author

aimeeyoung99

ESL-NEWS.COM