

光联AI Agent安全防护 服务解决方案

光联世纪

AI Agent删光Meta AI总监邮箱

Meta AI总监把最火的AI智能体OpenClaw接上了自己的工作邮箱。由于OpenClaw需要「压缩上下文」来处理大的信息量，在压缩的过程中，OpenClaw把之前设定的「未经批准不得操作」这条指令给忘了，疯狂删除邮件，喊停三次全部无视。



● 问题点

关键指令遗忘：大语言模型的底层机制决定上下文窗口有限，信息会被压缩，而被压缩掉的，可能恰好是最重要的那条安全指令。

AI Agent恶意Skill供应链投毒事件

2026年1月底，OpenSourceMalware（安全研究组织）称已在ClawHub上发现至少至少14个恶意'skills'，这些skills伪装成加密货币交易或钱包自动化工具，试图向用户系统投递恶意软件。

2月初，Koi Security（为OpenAI、Fortune Top50公司提供安全服务）称已在OpenClaw中发现341个恶意Skill。

Date	Event
Jan 27 - 29, 2026	Initial wave of at least 14 malicious "skills" uploaded to ClawHub.
Jan 27 - Feb 1, 2026	Second, larger wave published, bringing the total to over 230 malicious packages across ClawHub and GitHub.
Feb 1 - 2, 2026	Koi Security completes a full audit of all 2,857 skills on ClawHub, identifying 341 malicious entries—335 of which belong to the single ClawHavoc campaign.
Feb 2 - 3, 2026	Widespread reporting by security firms and community researchers. OpenClaw creator introduces a user-reporting feature in response.



● 问题点

- ① **缺乏Skill安全检查**：Agent Skill可被利用进行提示注入、RCE等攻击，使得Agent产生非预期操作。
- ② **Agent未校验Skill来源可信**：公开不等于可信，Agent应用需校验第三方源安全性。

AI Agent风险提示

关于OpenClaw安全应用的风险提示

原创 CNCERT 国家互联网应急中心CNCERT 2026年3月10日 19:02 北京

近期，OpenClaw（“小龙虾”，曾用名Clawdbot、Moltbot）应用下载与使用情况火爆，国内主流云平台均提供了一键部署服务。此款智能体软件依据自然语言指令直接操控计算机完成相关操作。为实现“自主执行任务”的能力，该应用被授予了较高的系统权限，包括访问本地文件系统、读取环境变量、调用外部服务应用程序编程接口（API）以及安装扩展功能等。然而，由于其默认的安全配置极为脆弱，攻击者一旦发现突破口，便能轻易获取系统的完全控制权。

前期，由于OpenClaw智能体的不当安装和使用，已经出现了一些严重的安全风险：

1.“提示词注入”风险。网络攻击者通过在网页中构造隐藏的恶意指令，诱导OpenClaw读取该网页，就可能导致其被诱导将用户系统密钥泄露。

2.“误操作”风险。由于错误的理解用户操作指令和意图，OpenClaw可能会将电子邮件、核心生产数据等重要信息彻底删除。

3.功能插件（skills）投毒风险。多个适用于OpenClaw的功能插件已被确认为恶意插件或存在潜在的安全风险，安装后可执行窃取密钥、部署木马后门软件等恶意操作，使得设备沦为“肉鸡”。

4.安全漏洞风险。截止目前，OpenClaw已经公开曝出多个高中危漏洞，一旦这些漏洞被网络攻击者恶意利用，则可能导致系统被控、隐私信息和敏感数据泄露的严重后果。对于个人用户，可导致隐私数据（像照片、文档、聊天记录）、支付账户、API密钥等敏感信息遭窃取。对于金融、能源等关键行业，可导致核心业务数据、商业机密和代码仓库泄露，甚至会使整个业务系统陷入瘫痪，造成难以估量的损失。

AI Agent风险总结

权限失控：高危操作执行



默认以高权限运行，易被诱导执行危险系统命令（如 `rm -rf /`），可能导致系统崩溃或永久性损坏。

注入攻击：提示词注入攻击



网络攻击者通过网页中隐藏恶意指令，诱导OpenClaw读取该网页，导致其被诱导将用户系统密钥泄露。

网络逃逸：横向移动风险



可利用网络工具扫描、探测内部网络，成为攻击者横向移动的跳板，威胁整个局域网安全。

数据泄露：核心机密窃取



具备广泛的文件访问能力，可能被利用窃取API密钥、数据库密码等核心机密，造成商业损失。

企业AI Agent安全防护全景架构

光联携手华为，以“云、网、安”一体化架构为基础，为企业构建可视、可控、可管、可运营的AI安全运营服务体系，助力AI安全落地与规模化应用。



华为AI安全运营三层体系

智能运营：
AI对抗AI



可视

可控

可管

网络围栏
【边界NGFW】

Agent围栏
【Agent安全网关】

主机围栏
【AIDR】

第一层

受控环境与最小权限



切断与公网的一切联系，白名单访问公网资源、内网模型和必要服务，严格遵循最小权限原则，防止系统层面的越权操作。

第二层

统一认证、授权、审计



所有流量必经Agent安全网关，实现统一的身份认证、权限控制和全链路审计，杜绝非法访问；模型围栏，过滤模型输入输出不合规内容及提示词注入攻击。

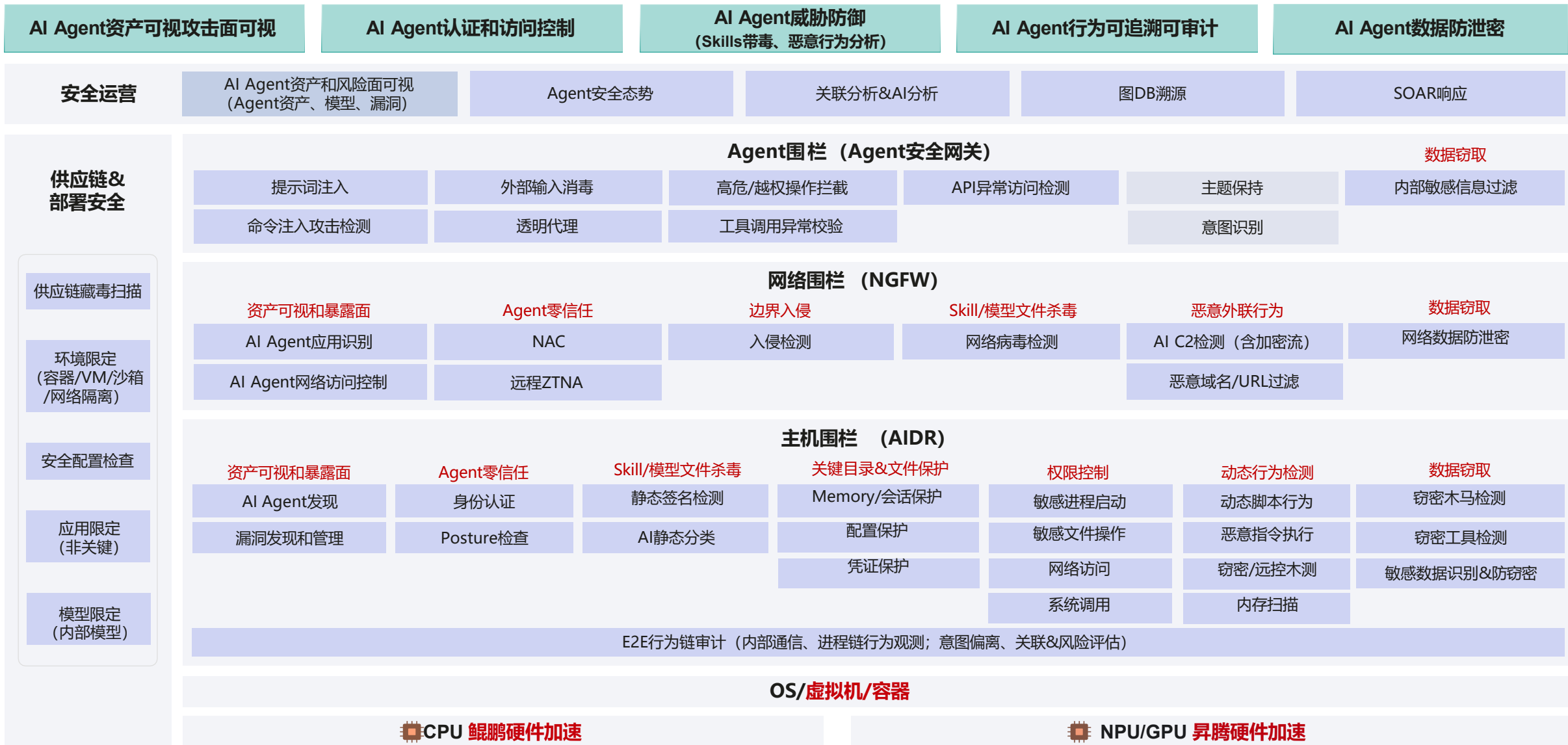
第三层

Agent异常行为监管



通过部署主机及服务器智算EDR，实时监控和审计AI Agentw在终端上的所有行为，及时发现异常。

华为AI Agent安全防护方案整体架构



安全防护：实现AI Agent可视、可控、可管

方案
价值



全面可视



全程可控



全链可管

AI 应用审计大屏



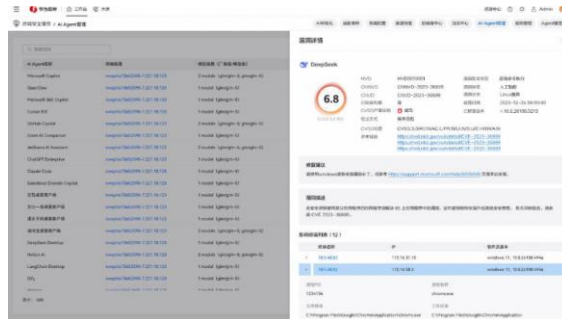
AI Agent资产和漏洞可视

AI资产可视

基本信息

AI Agent名称	Microsoft Copilot	安装包类型	Install.exe
安装路径	C:/Program Files	主程序	/usr/libexec/postfix/master
配置文件路径	/etc/docker/daemon.json	官网最新版本	v1.2.3
最近访问时间	2026-03-15 10:00:00	模型名称	ResNet50/Mistral
父进程	init	当前进程	process123
若干子进程	子进程1	进程状态	进行中
漏洞信息	3		

AI应用漏洞可视



光联不仅帮助客户“建好网络”，更持续帮助客户“用好网络、管优网络”。从建设到运营，全生命周期服务，成为客户值得信赖的长期合作伙伴



标准交付体系，保障项目高质量落地

光联通过标准化交付流程，专业团队保障，及资源保障能力，确保项目高质量落地，为长期稳定运行打下坚实基础。



光联骨干网：全球覆盖+安全智能，支撑高质量交付

依托全球骨干网与安全能力，为交付提供稳定、快速、安全的网络保障，确保项目高质量落地与稳定运行。



为交付体系提供坚实保障（五大维度）



智能运维体系，主动保障业务稳定运行

光联依托7x24专业运维体系，通过NOC网络运维、SOC安全运营、COC云托管运维协同联动，实现从监控、分析、响应到优化的闭环保障，为客户业务稳定运行保驾护航。



光联以AI运维能力 (AIOps) 与7x24持续运营能力为客户提供更智能、更高效、更稳定的运维保障，让网络、安全与云持续稳定支撑业务发展。

光联通过AIOps能力、标准化流程与7×24持续运维机制，构建智能运维闭环能力，持续保障客户网络、安全与云稳定运行。

NSOC网安融合运维中心



标准运维流程体系



THANK YOU

感谢聆听 期待合作



关注微信公众号



关注视频号