

Agent IA Jailbreak 2026 : MCP & Tool Injection



16 mai
2026



Mis à jour le 17 mai
2026



20 min de
lecture



3420
mots



Les agents IA branches sur outils (MCP, fonction calling) sont vulnérables au plan hijacking. ASR > 70% sur agents non hardenés.

À RETENIR

A retenir — Agent IA Jailbreak via MCP & Tool Injection

Tool injection : un IPI dans une donnée retrieved force l'agent à appeler de (envoi mail, exfil fichiers, transfert fonds).

MCP (Model Context Protocol) standardise depuis novembre 2024 facilite multiplie la surface d'attaque (servers tiers non audités).

Plan hijacking : manipulation du ReAct loop pour devier le plan d'exécution l'apparence de cohérence.

In projet cyber security
Réponse sous 24h
Cas réel 2025 : agent SOC compromis via pro de 12000 credentials.

Devis
gratuit



un email P

Defenses 2026 : *tool allowlists, human-in-the-loop pour actions critiques, MCP, output filtering.*

L'**agent ai jailbreak** est, en 2026, le risque le plus critique de l'écosystème LLM en combine trois ingrédients dangereux : (1) un LLM avec ses vulnérabilités d'alignement (Suffix, Multi-Turn Jailbreaks Crescendo), (2) un accès à des outils (file system, na (3) un loop d'exécution autonome (ReAct, Plan-and-Execute) qui orchestre les act combine Indirect Prompt Injection RAG sur le contexte avec une tool injection peut des actions catastrophiques : envoi de mails frauduleux, exfiltration de données, t article présente les attaques (tool injection via MCP, plan hijacking, ReAct loop ma Python exploit, et les défenses (allowlists, HITL, capability scopes). Pour les RSSI, devenu la priorité numéro un en 2026, dépassant les attaques sur LLM unitaires en business. La conformité AI Act article 14 (human oversight) impose explicitement c

1. Genèse et état de l'art

La littérature sur les agents IA émerge en 2022 avec ReAct (Yao et al., 2022) puis (2023). Pour la sécurité, les premiers travaux significatifs :

Greshake et al. (2023) — *Not what you've signed up for*, démonstration d'IPI sur ChatGPT plugins).

Cohen et al. (2024) — *CompromptMized*, premier ver IA via agents email.

Anthropic (2024) — publication du Model Context Protocol (MCP) en novembre l'interface agent-outil.

Bagdasaryan et al. (2024) — *Adversarial Illusions* embeddings, multimodaux.

Réponse sous 24h

Devis
gratuit



Réponse sous 24h

Devis
gratuit →