



Indirect Prompt Injection 2026 : Empoisonner u LLM



16 mai 2026



Mis à jour le 17 mai 2026



20 min de lecture



3762 mots



L'indirect prompt injection (IPI) injecte une instruction adversariale dans un document du corpus RAG. ASR > 60% sur LLM 2026 sans defense.

À RETENIR

A retenir — Indirect Prompt Injection

IPI (Greshake et al., 2023) injecte une instruction adversariale dans un doc corpus RAG. ASR **74%** sur GPT-5 en 2026.

Vecteurs courants : commentaires HTML, metadata PDF, alt-text image, co base64, white-on-white CSS, ZWS Unicode steganographique.

Impact business : exfiltration de PII, manipulation de reponses chatbot client, propagation par *worm prompts* (Comprompter, 2024).

In projet cybersécurité
Reponse sous 24h

Devis gratuit



Defenses 2026 : **Spotlighting** (Hines et al., Microsoft 2024), **StruQ** (Chen et al., 2024), Constitutional Classifiers, isolation prompts/data.

OWASP LLM01 classe IPI comme attaque prioritaire. AI Act annexe IV exige la documentation des tests IPI.

L'**indirect prompt injection** (IPI) est l'attaque la plus pernicieuse de l'écosystème LLM. Contrairement aux jailbreaks frontaux ([GCG Adversarial Suffix](#), [Multi-Turn Jailbreak](#), [Crescendo](#)), l'IPI exploite une propriété structurelle des architectures [RAG \(Retrieval Augmented Generation\)](#) : un LLM ne distingue pas, dans son contexte, ce qui est issu de l'utilisateur et ce qui est donné en retrieval. Un attaquant qui contrôle un document en retrieval RAG d'une entreprise peut injecter des instructions adversariales que le LLM exécute s'il s'agit du prompt utilisateur. Cet article présente le code Python d'attaque, les types de payload (HTML, metadata, steganographie), les défenses 2026 (Spotlighting, Suffix Mapping conforme AI Act / [OWASP LLM Top 10](#)). Pour les architectes RAG 2026, l'IPI passe du statut de risque émergent à celui de menace opérationnelle confirmée, exigeant une approche de défense-in-depth coordonnée avec les régulateurs (CNIL, ENISA, NIS2).

1. Genèse et état de l'art

Le concept d'IPI est formalisé par Kai Greshake et al. en avril 2023 dans *Not what you signed up for: Compromising Real-World LLM-Integrated Applications with Indirect Prompt Injection*. La demo est saisissante : un site web piège contient une instruction cachée. Lorsqu'un assistant IA navigue le site (Bing Chat à l'époque), prend le contrôle de la conversation utilisateur.

Réponse sous 24h

Devis
gratuit



Réponse sous 24h

Devis
gratuit

