

Pentest IA : Méthodologie d'Audit des Systèmes Artificielle

30 April
2026Mis à jour le 30 April
202647 min de
lecture

Méthodologie pentest IA : OWASP ML Top 10, MITRE ATLAS, prompt injection, poisoning. Outils Garak, PromptFoo, ART. Conformité AI Act.

Le **pentest IA** et l'**audit de sécurité des systèmes d'intelligence artificielle** émergent comme des enjeux critiques face à la prolifération des modèles de langage (LLM), des systèmes de recommandation et des architectures RAG (Retrieval-Augmented Generation) dans les infrastructures cloud. L'application des méthodologies offensives traditionnelles aux systèmes d'IA — couvrant tout, de la **prompt injection** avancée à l'extraction de modèles, en passant par l'empoisonnement des modèles, les attaques adversariales sur le machine learning, les vulnérabilités des API de services de modèles (Hugging Face, registres de modèles). Les référentiels **OWASP ML Top 10** et le **NIST RMF** (Risk Management Framework) fournissent les cadres méthodologiques pour auditer ces systèmes, ainsi que des outils comme **Garak**, **PromptFoo**, **ART** (Adversarial Robustness Toolbox) et des frameworks d'intrusion spécifiques à l'IA. Avec l'entrée en application de l'**AI Act** européen et les réglementations associées, l'audit de sécurité IA n'est plus optionnel — il devient une obligation réglementaire à haut risque. Ce guide expert couvre l'ensemble de la méthodologie de pentest LLM.

fondamentaux théoriques aux techniques d'exploitation avancées, en intégrant les
matière de détection et de défense.

À RETENIR

Points clés de cet article :

Le **pentest IA** applique les méthodologies offensives aux systèmes d'intelli
classique, RAG, agents autonomes

Les référentiels **OWASP ML Top 10**, **MITRE ATLAS** et **NIST AI RMF** structure
d'évaluation

La **prompt injection** (directe et indirecte) reste la vulnérabilité la plus critique
production

L'**extraction de modèle** (model stealing) et l'**inversion de données d'entraîn**
intellectuelle et la confidentialité

Les outils **Garak**, **PromptFoo**, **ART** et **Counterfit** automatisent les tests de s

L'**AI Act** européen impose des évaluations de sécurité obligatoires pour les
partir d'août 2026

La **supply chain IA** (modèles pré-entraînés, datasets, frameworks ML) con
croissante et sous-évaluée

Cadres méthodologiques pour le pentest IA

L'évaluation de sécurité des systèmes d'IA nécessite des cadres méthodologiques
spécifiques au machine learning et aux modèles de langage, au-delà des vulnérab
Trois référentiels majeurs structurent aujourd'hui la discipline : l'OWASP ML Top 10
