

Red Teaming LLM on-premise : mét

📅 16 mai 2026 • ↻ Mis à jour le 17 mai 2026 • ⌚ 18 min de lecture • ≡ 3439 mots •

Découvrez la méthodologie complète du red teaming LLM on-premise : pr
2025, outils garak et PyRIT pour auditer vos modèles IA internes.



À RETENIR

A retenir - Red Teaming LLM on-premise

Le **red teaming LLM on-premise** est devenu une discipline incontournable po
langage en environnement d'entreprise. Les vecteurs d'attaque spécifiques
exfiltration de données -- nécessitent des outils et une méthodologie adaptée
L'OWASP LLM Top 10 2025 structure désormais ces risques et fournit un cad
sécurité doivent impérativement maîtriser garak, PyRIT et promptfoo pour au

Le déploiement massif de modèles de langage (LLM) en environnement on-premise
surface d'attaque radicalement nouvelle que les équipes de sécurité doivent impérati
Réponse sous 24h s'impose comme la réponse méthodologique structurée face à ces risques :

Devis gratuit →

des applications classiques, mais d'explorer les comportements émergents, les faiblesses inhérentes aux architectures transformer. En 2026, avec l'accélération des déploiements de modèles comme Qwen 2.5 en interne, comprendre comment un attaquant peut détourner, extraire ou manipuler les données d'un LLM est devenu une compétence fondamentale pour tout RSSI ou équipe red team. Cette session propose une vue d'ensemble complète, des outils éprouvés et des cas concrets pour auditer efficacement vos modèles. Nous aborderons jusqu'au reporting de remédiation. Nous couvrons les techniques d'attaque documentées, les outils de référence, la taxonomie des jailbreaks et les procédures de vérification post-audit.

Comprendre la surface d'attaque spécifique aux LLM on-premise

Contrairement aux applications web traditionnelles, un LLM déployé en entreprise expose une surface d'attaque qui s'étend bien au-delà des interfaces réseau classiques. Le modèle lui-même -- ainsi que les couches d'entraînement fine-tune -- constitue une cible à part entière. Les **vecteurs d'attaque** incluent les attaques sur l'entrée (prompt manipulation), les attaques sur le contexte (system prompt injection), et les attaques sur les sorties (data exfiltration via génération, model inversion).

Pour un LLM on-premise typique, la surface inclut l'API d'inférence (REST ou gRPC), les connecteurs de données (RAG, base documentaire, outils externes), et les pipelines de données. Chacune de ces couches mérite une analyse spécifique et systématique.

Les organisations qui déploient des LLM sans red teaming préalable s'exposent à des risques de fuite d'informations confidentielles présentes dans le contexte système, contournement des politiques de sécurité, et manipulation des agents IA pour exécuter des actions non autorisées. Une étude récente a révélé que 80% des LLM d'entreprise testés étaient vulnérables à au moins une forme de prompt injection.

La classification des modèles on-premise par niveau de risque est une première étape cruciale. Un modèle accessible via un réseau présente un profil de risque bien différent d'un agent autonome avec accès limité.

Cette cartographie conditionne la profondeur et les méthodes de red teaming.

Réponse sous 24h

Devis
gratuit



Réponse sous 24h

Devis
gratuit →