

# Attaques RAG : empoisonnement vectoriel

📅 16 mai 2026 • 🔄 Mis à jour le 17 mai 2026 • ⌚ 16 min de lecture • ☰ 3107 mots •

Maîtrisez les attaques contre les systèmes RAG : empoisonnement vectoriel. Défenses et hardening pour sécuriser vos bases vectorielles.



**À RETENIR**

## A retenir - Attaques RAG et empoisonnement vectoriel

Les **attaques contre les systèmes RAG** représentent l'une des menaces les plus critiques pour une entreprise. Un attaquant capable d'injecter des documents malveillants dans les réponses du système, exfiltrer des données sensibles et contourner les garde-fous de sécurité. Ces attaques sont discrètes, persistantes et difficiles à détecter sans un monitoring spécialisé. La défense repose sur des techniques avancées telles que le monitoring de la similarité cosinus et l'isolation des sources de données.

Les architectures **RAG (Retrieval-Augmented Generation)** sont devenues le standard pour les applications d'IA générative, permettant d'ancrer les réponses du modèle dans une base de connaissances pertinente. Cependant, cette puissance vient avec une surface d'attaque non négligeable et souvent sous-estimée.

Réponse sous 24h

**Devis gratuit** →

**vectoriel** exploitent la confiance implicite accordée aux documents récupérés par d'injecter des chunks malveillants dans la base vectorielle peut influencer, manipuler le système, transformant votre assistant IA métier en vecteur d'exfiltration ou de données. Parmi les solutions RAG sur Pinecone, Qdrant, Weaviate et ChromaDB, la sécurisation de la base de données est une priorité de sécurité que tout RSSI doit intégrer dans sa stratégie de protection des systèmes d'information. Les attaques documentées, les mécanismes sous-jacents, et les défenses concrètes pour protéger les vecteurs les plus fréquemment exploités par les attaquants en 2026.

---

## Architecture RAG et surface d'attaque -- vue d'ensemble

---

Un système RAG typique comprend quatre composants : le pipeline d'ingestion (crawling, API tierces), la base **vectorielle** (stockage des embeddings), le retrieval engine (recherche de similarité) et le LLM, s'appuyant sur les chunks récupérés. Chacun de ces composants est potentiellement vulnérable.

La surface d'attaque se structure ainsi :

**Pipeline d'ingestion** : injection de documents malveillants via les sources de données (crawling, API tierces)

**Base vectorielle** : accès direct à l'API de la vector DB pour insertion de vecteurs malveillants (accès insuffisante)

**Retrieval** : manipulation de la requête d'embedding pour biaiser les résultats de recherche

**Contexte LLM** : injection d'instructions via les chunks récupérés (indirect poisoning)

Pour aller plus loin sur les attaques d'injection indirecte dans les RAG, consultez notre article sur [indirect poisoning](#). Les aspects de red teaming LLM complet incluant le testing RAG sont couverts dans notre [première](#).

---

---

Réponse sous 24h

Devis  
gratuit →