

Guide RAG - Retrieval-Augmented Generation

Introduction au RAG

Le RAG est une technique qui combine la recherche d'information et la génération de texte pour créer des systèmes d'IA plus précis.

Architecture

1. Ingestion des documents : Les documents sont découpés en chunks
2. Vectorisation : Chaque chunk est transformé en vecteur d'embedding
3. Stockage : Les vecteurs sont stockés dans une base vectorielle
4. Recherche : Les requêtes sont vectorisées et comparées aux documents
5. Génération : Un LLM génère une réponse basée sur les documents trouvés

Avantages du RAG

- Réponses basées sur des sources fiables et à jour
- Réduction des hallucinations du modèle
- Possibilité de citer les sources utilisées
- Adaptation facile à de nouveaux domaines