

# RAG ET ORACLE 23AI

Retrieval Augmented Generation avec Intelligence Artificielle Générative

Ce cours présente les concepts fondamentaux du Retrieval Augmented Generation (RAG) et son intégration avec Oracle 23ai et les services d'IA générative d'Oracle Cloud Infrastructure (OCI). Vous apprendrez à construire des applications intelligentes capables de générer des réponses contextualisées en s'appuyant sur des bases de données vectorielles.

## Table des matières

---

<b>1</b>	<b>Introduction aux Intégrations OCI Generative AI</b>	<b>5</b>
1.1	LangChain : Framework pour Applications LLM	5
1.2	Types de Modèles dans LangChain	5
1.3	Gestion des Prompts	5
1.3.1	Template de Prompt	5
1.3.2	Template de Prompt de Chat	6
1.4	Création de Chaînes	6
1.5	Flux de Traitement avec LangChain	6
1.6	Gestion de la Mémoire	6
1.7	Oracle 23ai comme Magasin de Vecteurs	6
1.8	Intégrations OCI Generative AI	6
<b>2</b>	<b>Retrieval Augmented Generation (RAG)</b>	<b>7</b>
2.1	Concept et Problématique	7
2.2	Avantages du RAG	7
2.2.1	Réduction des Biais et Erreurs	7
2.2.2	Dépassement des Limitations du Modèle	7
2.2.3	Extensibilité	7
2.3	Architecture du Pipeline RAG	7
2.3.1	Phase d'Ingestion	7
2.3.2	Phase de Récupération	8
2.3.3	Phase de Génération	8
2.4	Efficacité du RAG	8
<b>3</b>	<b>Traitement des Documents</b>	<b>8</b>
3.1	Chargement des Documents	8
3.1.1	Variété des Sources et Formats	8
3.1.2	Classes de Chargement LangChain	8
3.2	Segmentation des Documents	8
3.2.1	Considérations pour la Taille des Chunks	9
3.2.2	Continuité Contextuelle avec Chunk Overlap	9
3.2.3	Segmentation Sémantique	9
3.3	Exemple de Code pour le Traitement PDF	10
<b>4</b>	<b>Embeddings et Stockage des Documents</b>	<b>10</b>
4.1	Comprendre les Embeddings	10
4.1.1	Analogie Conceptuelle	10
4.1.2	Représentation Visuelle	11
4.2	Génération des Embeddings	11
4.2.1	Modèles d'Embeddings	11
4.2.2	Approches de Génération	11
4.3	Stockage dans Oracle 23ai	11
4.3.1	Type de Données Vector	11
4.3.2	Opérations de Base de Données Standards	11
4.4	Exemple d'Implémentation	12

<b>5</b>	<b>Récupération et Génération</b>	<b>12</b>
5.1	Phase de Récupération . . . . .	13
5.1.1	Processus de Recherche . . . . .	13
5.1.2	Mesures de Similarité . . . . .	13
5.2	Optimisation des Performances . . . . .	13
5.2.1	Problématique de l'Échelle . . . . .	13
5.2.2	Solution : Indexation . . . . .	13
5.2.3	Techniques d'Indexation . . . . .	13
5.3	Phase de Génération . . . . .	14
5.3.1	Contextualisation . . . . .	14
5.3.2	Exemple d'Implémentation . . . . .	15
<b>6</b>	<b>RAG Conversationnel</b>	<b>15</b>
6.1	Concept du RAG Conversationnel . . . . .	16
6.1.1	Différence avec RAG Standard . . . . .	16
6.1.2	Importance du Contexte Conversationnel . . . . .	16
6.2	Gestion de la Mémoire . . . . .	16
6.2.1	Concept de Mémoire . . . . .	16
6.2.2	Types de Mémoire LangChain . . . . .	16
6.3	Architecture du RAG Conversationnel . . . . .	17
6.4	Implémentation Pratique . . . . .	17
<b>7</b>	<b>Considérations Techniques et Bonnes Pratiques</b>	<b>17</b>
7.1	Optimisation des Performances . . . . .	17
7.1.1	Gestion de la Taille des Chunks . . . . .	18
7.1.2	Stratégies de Recherche . . . . .	18
7.2	Qualité des Embeddings . . . . .	18
7.2.1	Choix du Modèle d'Embeddings . . . . .	18
7.2.2	Cohérence des Modèles . . . . .	18
7.3	Surveillance et Évaluation . . . . .	18
7.3.1	Métriques de Performance . . . . .	18
7.3.2	Tests et Validation . . . . .	19
<b>8</b>	<b>Architecture et Intégration avec Oracle 23ai</b>	<b>19</b>
8.1	Architecture Globale . . . . .	19
8.1.1	Composants Principaux . . . . .	19
8.1.2	Flux de Données . . . . .	19
8.2	Avantages d'Oracle 23ai . . . . .	19
8.2.1	Performance Native . . . . .	19
8.2.2	Intégration SQL . . . . .	19
8.2.3	Sécurité et Gouvernance . . . . .	20
<b>9</b>	<b>Cas d'Usage et Applications</b>	<b>20</b>
9.1	Applications Métiers . . . . .	20
9.1.1	Support Client Intelligent . . . . .	20
9.1.2	Analyse Documentaire . . . . .	20
9.1.3	Formation et Éducation . . . . .	20
9.2	Défis et Limitations . . . . .	21
9.2.1	Défis Techniques . . . . .	21
9.2.2	Considérations Éthiques . . . . .	21

<b>10 Évolutions et Perspectives</b>	<b>21</b>
10.1 Tendances Technologiques . . . . .	21
10.1.1 Amélioration des Modèles . . . . .	21
10.1.2 Optimisations Techniques . . . . .	21
10.2 Intégration Avancée avec Oracle . . . . .	21
10.2.1 Fonctionnalités Émergentes . . . . .	21
10.2.2 Écosystème Oracle . . . . .	22
<b>11 Conclusion</b>	<b>22</b>
11.1 Points Clés à Retenir . . . . .	22
11.2 Avantages de l'Approche Oracle . . . . .	22
11.3 Perspectives d'Avenir . . . . .	22
11.4 Recommandations Pratiques . . . . .	23

## 1 Introduction aux Intégrations OCI Generative AI

Oracle Cloud Infrastructure (OCI) Generative AI offre une variété de fonctionnalités pour générer, résumer et créer des embeddings de données. Cette section explore comment ces services s'intègrent avec d'autres frameworks et services OCI utiles, permettant de construire des applications avancées.

### 1.1 LangChain : Framework pour Applications LLM

LangChain est un framework destiné au développement d'applications alimentées par des modèles de langage. Il permet de créer des applications conscientes du contexte qui s'appuient sur un modèle de langage pour répondre en fonction du contexte fourni.

LangChain offre une multitude de composants qui facilitent la construction d'applications basées sur les LLM avec un effort minimal. Les principaux composants incluent :

- **Modèles de langage large (LLM)** : Le cœur de toute application
- **Prompts** : Instructions et contexte pour le modèle
- **Mémoire** : Stockage des conversations passées
- **Chaînes** : Orchestration des composants
- **Magasins de vecteurs** : Stockage des embeddings
- **Chargeurs de documents** : Import de différents formats
- **Diviseurs de texte** : Segmentation intelligente du contenu

Un avantage majeur de LangChain est l'interchangeabilité de ses composants. Par exemple, il est possible de passer d'un LLM à un autre avec des modifications minimales du code.

### 1.2 Types de Modèles dans LangChain

LangChain intègre deux types principaux de modèles, définis par leurs types d'entrée et de sortie :

#### Point important

**LLM** dans LangChain fait référence à un modèle de complétion de texte pur. Ces modèles prennent une chaîne de caractères comme prompt en entrée et produisent une chaîne de caractères en sortie.

**Modèles de chat** sont souvent basés sur des LLM mais sont spécifiquement ajustés pour les conversations. Ils prennent une liste de messages de chat en entrée et retournent un message IA en sortie.

### 1.3 Gestion des Prompts

Dans LangChain, les prompts peuvent être créés en utilisant deux types de classes principales :

#### 1.3.1 Template de Prompt

Créé à partir d'une chaîne Python formatée qui combine du texte fixe et des espaces réservés remplis lors de l'exécution. Généralement utilisé avec les modèles de génération, mais compatible avec les modèles de chat.

### 1.3.2 Template de Prompt de Chat

Composé d'une liste de messages de chat, chacun ayant un rôle et un contenu. Utilisé spécifiquement avec les modèles de chat.

## 1.4 Création de Chaînes

LangChain propose des frameworks pour créer des chaînes de composants. Deux approches principales :

1. **LangChain Expression Language (LCEL)** : Approche déclarative et préférée
2. **Classes Python LangChain** : Comme LLM Chain, approche plus traditionnelle

## 1.5 Flux de Traitement avec LangChain

### Exemple

Le processus typique suit cette séquence :

1. L'utilisateur soumet une requête
2. Le système utilise des prompts pour rassembler le contexte additionnel
3. Les prompts combinent le texte fixe et les variables capturées
4. La chaîne orchestre les opérations : entrée → prompt → LLM → réponse

## 1.6 Gestion de la Mémoire

La mémoire stocke les conversations avec le chatbot à un moment donné. Le processus fonctionne ainsi :

- La chaîne récupère la conversation (série de messages) depuis la mémoire
- Elle transmet cette conversation au LLM avec la nouvelle question
- Une fois la réponse générée, elle sauvegarde la question et la réponse en mémoire

LangChain propose divers types de mémoire selon ce qui est retourné : résumé du contenu, contenu actuel, ou même des entités extraites comme les noms.

## 1.7 Oracle 23ai comme Magasin de Vecteurs

Oracle 23ai peut être utilisé comme magasin de vecteurs, et LangChain fournit des classes Python pour stocker et rechercher des embeddings dans Oracle 23ai.

## 1.8 Intégrations OCI Generative AI

OCI Generative AI s'intègre avec Oracle 23ai de plusieurs façons :

1. **Génération d'embeddings externes** : Accès au service OCI Generative AI via DB UTILS et API REST
2. **SELECT AI d'Oracle 23ai** : Utilise OCI Generative AI pour générer du SQL à partir de langage naturel
3. **Classes LangChain** : Intégration directe via les classes LangChain

Les applications utilisant à la fois OCI Generative AI et Oracle 23ai peuvent également être créées en utilisant le SDK Python.

## 2 Retrieval Augmented Generation (RAG)

### 2.1 Concept et Problématique

Le Retrieval Augmented Generation (RAG) est une approche qui améliore les modèles de langage traditionnels en récupérant des informations à jour à partir de sources externes et en fournissant ces informations additionnelles et spécifiques au LLM, accompagnées de la requête utilisateur.

Les modèles de langage traditionnels génèrent des réponses basées uniquement sur leurs données d'entraînement, qui peuvent devenir obsolètes. RAG résout ce problème en récupérant des informations actualisées et en améliorant le contexte fourni au LLM pour générer des réponses plus pertinentes.

### 2.2 Avantages du RAG

#### 2.2.1 Réduction des Biais et Erreurs

Les LLM standard peuvent parfois perpétuer des biais ou des erreurs présents dans leurs données d'entraînement. RAG peut atténuer cela en puisant dans une variété de perspectives et de sources, conduisant à des réponses plus équilibrées et précises.

#### 2.2.2 Dépassement des Limitations du Modèle

RAG peut surmonter les limitations du modèle telles que les limites de tokens, puisque nous ne fournissons que les top K résultats de recherche aux LLM au lieu de l'ensemble du document.

#### 2.2.3 Extensibilité

RAG permet aux modèles de traiter un éventail plus large de requêtes sans nécessiter des ensembles de données d'entraînement exponentiellement plus importants.

### 2.3 Architecture du Pipeline RAG

Le pipeline RAG de base comprend trois phases principales :

#### 2.3.1 Phase d'Ingestion

##### Point important

Cette première phase est cruciale pour la qualité du système RAG final. Elle comprend :

1. **Chargement des documents** : Import du corpus de texte original
2. **Segmentation** : Division en morceaux plus petits et gérables (chunks)
3. **Création d'embeddings** : Transformation en représentations mathématiques
4. **Indexation** : Stockage dans une base de données pour récupération efficace

### 2.3.2 Phase de Récupération

Cette phase utilise les données indexées pour trouver des informations pertinentes :

1. **Requête utilisateur** : L'utilisateur soumet une question
2. **Recherche** : Le système cherche dans les embeddings stockés
3. **Sélection** : Choix des top K résultats les plus pertinents

### 2.3.3 Phase de Génération

Phase finale où le système génère une réponse basée sur les informations récupérées :

1. **Contextualisation** : Les chunks sélectionnés sont fournis au modèle génératif
2. **Génération** : Le modèle utilise le contexte pour produire une réponse cohérente et pertinente

## 2.4 Efficacité du RAG

L'architecture RAG est particulièrement efficace dans les scénarios où les modèles génératifs doivent être complétés avec des informations spécifiques qui peuvent ne pas être présentes dans les données d'entraînement.

## 3 Traitement des Documents

---

Cette section détaille la phase d'ingestion du pipeline RAG, en commençant par le chargement et la segmentation des documents.

### 3.1 Chargement des Documents

#### 3.1.1 Variété des Sources et Formats

Les documents peuvent provenir de diverses sources et avoir multiple formats :

- **PDF** : Documents formatés, rapports, manuels
- **CSV** : Données tabulaires, feuilles de calcul
- **HTML** : Pages web, documentation en ligne
- **JSON** : Données structurées, APIs
- **Autres formats** : TXT, DOCX, XML, etc.

#### 3.1.2 Classes de Chargement LangChain

La plupart des frameworks LLM, y compris LangChain, offrent des classes spécialisées pour charger différents types de documents. Ces classes de chargement supportent également :

- Chargement d'un document unique
- Chargement de tous les documents d'un répertoire donné
- Traitement par lots pour de gros volumes

### 3.2 Segmentation des Documents

Une fois les documents chargés, l'étape suivante consiste à les diviser en morceaux plus petits, également appelés chunks.



### 3.2.1 Considérations pour la Taille des Chunks

#### Point important

La taille des chunks est un paramètre critique qui impacte directement la qualité du système RAG :

**Contraintes techniques** : La plupart des LLM ont des contraintes de taille d'entrée maximale, limitant la taille des chunks par la fenêtre contextuelle du LLM.

**Équilibre sémantique** :

- Chunks trop petits : Peuvent ne pas être sémantiquement utiles
- Chunks trop grands : Peuvent ne pas être sémantiquement spécifiques

### 3.2.2 Continuité Contextuelle avec Chunk Overlap

Pour maintenir la continuité du contexte d'un chunk à l'autre, nous incluons une partie du chunk précédent dans le chunk suivant. Cette technique est appelée **chunk overlap**.

#### Exemple

Si nous avons un chunk de 1000 caractères avec un overlap de 200 caractères :

- Chunk 1 : caractères 1-1000
- Chunk 2 : caractères 801-1800 (overlap de 200)
- Chunk 3 : caractères 1601-2600 (overlap de 200)

### 3.2.3 Segmentation Sémantique

La segmentation cherche à préserver le sens en utilisant la structure naturelle du texte :

1. **Séparateurs de paragraphe** : Division primaire
2. **Séparateurs de phrase** : Division secondaire
3. **Séparateurs de mot** : Division tertiaire

L'objectif est d'obtenir des chunks de la taille désirée tout en conservant des chunks sémantiquement riches.

### 3.3 Exemple de Code pour le Traitement PDF

#### Exemple

Voici un exemple de code pour charger et diviser un document PDF :

```
# Chargement du PDF
pdf_reader = PDFReader()
reader_object = pdf_reader.load("document.pdf")

# Extraction du texte
text = ""
for page in reader_object.pages:
    text += page.extract_text()

# Création du diviseur de texte
text_splitter = TextSplitter(
    chunk_size=1000,
    chunk_overlap=200
)

# Division du texte en chunks
chunks = text_splitter.split_text(text)
```

## 4 Embeddings et Stockage des Documents

Cette section couvre la transformation des chunks en embeddings et leur stockage dans Oracle 23ai.

### 4.1 Comprendre les Embeddings

#### 4.1.1 Analogie Conceptuelle

Les embeddings sont des représentations mathématiques qui permettent aux machines de comprendre la similarité sémantique entre mots, phrases ou documents. Dans un espace multidimensionnel, les embeddings d'éléments similaires sont proches les uns des autres.

#### Exemple

Considérons trois groupes de mots : animaux, fruits, et lieux. Si nous donnons le mot "tigre" :

- Les humains le placent naturellement dans le groupe "animaux"
- Les machines utilisent les embeddings pour faire cette même association
- L'embedding de "tigre" sera plus proche des embeddings du groupe "animaux"

### 4.1.2 Représentation Visuelle

Dans une représentation bidimensionnelle simplifiée, nous pouvons observer comment les mots sémantiquement similaires se regroupent naturellement. Les embeddings capturent cette similarité sémantique dans un espace multidimensionnel (généralement 768, 1024, ou plus de dimensions).

## 4.2 Génération des Embeddings

### 4.2.1 Modèles d'Embeddings

Les embeddings sont créés en utilisant des modèles d'embeddings entraînés. Oracle 23ai supporte l'utilisation de modèles d'embeddings à la fois à l'intérieur et à l'extérieur de la base de données.

### 4.2.2 Approches de Génération

#### Point important

**Génération externe** : Pour générer des embeddings en dehors de la base de données, des modèles d'embeddings tiers peuvent être utilisés via les services OCI Generative AI.

**Génération interne** : Pour conserver les données entièrement dans la base de données, nous pouvons importer des modèles d'embeddings au format ONNX dans Oracle 23ai et créer les embeddings directement dans la base de données.

## 4.3 Stockage dans Oracle 23ai

### 4.3.1 Type de Données Vector

Oracle 23ai a introduit un nouveau type de données **vector** pour stocker les embeddings dans une colonne de base de données. Cette colonne peut être créée aux côtés d'autres types de données traditionnels.

#### Exemple

Création d'une table avec colonne vector :

```
CREATE TABLE documents (  
  id NUMBER PRIMARY KEY,  
  content CLOB,  
  embeddings VECTOR(1024),  
  metadata JSON  
);
```

### 4.3.2 Opérations de Base de Données Standards

Les embeddings peuvent être manipulés en utilisant les instructions SQL habituelles :

- **INSERT** : Ajouter de nouveaux embeddings
- **UPDATE** : Modifier des embeddings existants
- **SELECT** : Récupérer des embeddings avec recherche de similarité

## 4.4 Exemple d'Implémentation

### Exemple

Code pour créer des embeddings et les stocker :

```
# Connexion à la base de données
connection = oracle_db_connect(
    username="user",
    password="password",
    dsn="database_dsn"
)

# Transformation des chunks en documents
documents = []
for i, chunk in enumerate(chunks):
    doc_dict = {
        'id': i,
        'link': source_link,
        'text': chunk
    }
    document = chunks_to_docs_wrapper(doc_dict)
    documents.append(document)

# Création du modèle d'embeddings
embedding_model = OCIGenAIEmbeddings(
    model_name="cohere.embed-multilingual-v3.0",
    service_endpoint="https://inference.generativeai...",
    compartment_id="compartment_id",
    auth_type="API_KEY"
)

# Création du vector store
vector_store = OracleVS.from_documents(
    documents=documents,
    embedding=embedding_model,
    connection=connection,
    table_name="document_embeddings",
    distance_strategy="COSINE"
)
```

Le vector store est maintenant prêt pour effectuer des recherches de documents correspondant aux requêtes utilisateur.

## 5 Récupération et Génération

Cette section détaille les phases de récupération et de génération du pipeline RAG.

## 5.1 Phase de Récupération

### 5.1.1 Processus de Recherche

Lorsqu'un utilisateur soumet une requête :

1. **Encodage de la requête** : La requête est transformée en embedding en utilisant le même modèle d'embeddings utilisé pour encoder les chunks
2. **Recherche vectorielle** : Une recherche est effectuée dans la base de données où les chunks embeddings sont stockés
3. **Sélection des résultats** : Le système retourne les chunks les plus similaires à la requête

### 5.1.2 Mesures de Similarité

#### Point important

La recherche vectorielle utilise typiquement deux mesures de similarité principales :

**Produit scalaire (Dot Product)** : Mesure la magnitude de la projection d'un vecteur sur l'autre. Considère à la fois la magnitude et l'angle entre les vecteurs pour calculer la similarité.

**Distance cosinus (Cosine Distance)** : Ne considère que l'angle entre les vecteurs, pas la magnitude, pour calculer la similarité.

Dans le contexte du NLP :

- Plus de magnitude = contenu sémantiquement plus riche
- Moins d'angle = plus de similarité

## 5.2 Optimisation des Performances

### 5.2.1 Problématique de l'Échelle

Comparer l'embedding de la requête utilisateur avec l'embedding de chaque chunk est acceptable pour un petit nombre d'embeddings, mais devient problématique à grande échelle.

### 5.2.2 Solution : Indexation

Les index sont comme des tables des matières qui permettent de trouver facilement les embeddings. Ce sont des structures de données spécialisées conçues pour les recherches de similarité.

### 5.2.3 Techniques d'Indexation

Diverses techniques comme le clustering, le partitionnement et les graphes de voisinage sont utilisées pour grouper les embeddings, ce qui aide à réduire l'espace de recherche.

**Hierarchical Navigable Small-World Graph (HNSW)** : Forme de graphe de voisinage en mémoire pour la recherche vectorielle. Index très efficace pour la recherche de similarité approximative.

**Inverted File Flat (IVF)** : Forme d'index vectoriel basé sur le partitionnement de voisinage. Index basé sur les partitions qui atteint l'efficacité de recherche en réduisant la zone de recherche grâce aux partitions ou clusters de voisinage.

## 5.3 Phase de Génération

### 5.3.1 Contextualisation

Une fois le contexte récupéré sous forme de chunks pertinents, nous pouvons l'envoyer au LLM avec la requête. Le LLM considère le contexte et la requête utilisateur pour générer une réponse pertinente et spécifique.

### 5.3.2 Exemple d'Implémentation

#### Exemple

Code pour la récupération et génération :

```
# Imports nécessaires
from langchain.chains import RetrievalQA
from langchain.llms import ChatOCIGenAI
from langchain.vectorstores import OracleVS

# Création du vector store
vector_store = OracleVS(
    embedding=embedding_model,
    connection=connection,
    table_name="document_embeddings",
    distance_strategy="COSINE"
)

# Création du retriever
retriever = vector_store.as_retriever(
    search_type="similarity",
    search_kwargs={"k": 3} # Top 3 résultats
)

# Création du LLM
llm = ChatOCIGenAI(
    model_id="cohere.command-r-plus",
    service_endpoint="https://inference.generativeai...",
    compartment_id="compartment_id",
    auth_type="API_KEY"
)

# Création de la chaîne RAG
qa_chain = RetrievalQA.from_chain_type(
    llm=llm,
    retriever=retriever,
    return_source_documents=True
)

# Invocation avec une question
response = qa_chain.invoke({"query": "Votre question ici"})
answer = response['result']
sources = response['source_documents']
```

## 6 RAG Conversationnel

Cette section explore l'implémentation du RAG pour créer des chatbots conversationnels.

## 6.1 Concept du RAG Conversationnel

### 6.1.1 Différence avec RAG Standard

Le RAG conversationnel étend le RAG standard en maintenant un historique des interactions pour créer une expérience de chat continue.

Un chat est une série de questions et réponses où :

1. L'utilisateur pose une question
2. Le LLM fournit une réponse
3. L'utilisateur pose une question de suivi
4. Le processus se répète

### 6.1.2 Importance du Contexte Conversationnel

#### Exemple

Considérons cette séquence :

- Question 1 : "Parlez-moi de Las Vegas"
- Réponse 1 : [Informations sur Las Vegas]
- Question 2 : "Parlez-moi de sa température typique tout au long de l'année"

La seconde question fait référence à "sa" (Las Vegas), nécessitant le contexte de la conversation précédente.

## 6.2 Gestion de la Mémoire

### 6.2.1 Concept de Mémoire

#### Point important

Pour maintenir une liste des questions posées et des réponses données, le concept de **mémoire** est utilisé :

- Les contenus de la mémoire sont mis à jour chaque fois qu'une nouvelle question est posée et qu'une réponse est générée
- Les contenus de la mémoire sont transmis comme contexte additionnel au LLM
- Le LLM répond à la question suivante en considérant les documents récupérés ET l'historique de conversation

### 6.2.2 Types de Mémoire LangChain

LangChain propose une variété de types de mémoire selon ce qui est retourné de la mémoire :

- **Mémoire complète** : Retourne tout l'historique de conversation
- **Mémoire résumée** : Retourne un résumé des contenus plutôt que le contenu actuel
- **Mémoire d'entités** : Retourne des entités extraites comme les noms, lieux, etc.



### 6.3 Architecture du RAG Conversationnel

Le système RAG conversationnel fonctionne ainsi :

1. **Récupération de l'historique** : La chaîne récupère la conversation depuis la mémoire
2. **Recherche vectorielle** : Recherche de documents pertinents pour la nouvelle question
3. **Contextualisation** : Combinaison de l'historique + documents récupérés
4. **Génération** : Le LLM génère une réponse contextuelle
5. **Mise à jour mémoire** : Sauvegarde de la nouvelle question et réponse

### 6.4 Implémentation Pratique

#### Exemple

Exemple d'implémentation d'un chatbot RAG conversationnel :

```
from langchain.memory import ConversationBufferMemory
from langchain.chains import ConversationalRetrievalChain

# Création de la mémoire
memory = ConversationBufferMemory(
    memory_key="chat_history",
    return_messages=True,
    output_key='answer'
)

# Création de la chaîne conversationnelle
conversation_chain = ConversationalRetrievalChain.from_llm(
    llm=llm,
    retriever=retriever,
    memory=memory,
    return_source_documents=True
)

# Boucle de conversation
while True:
    user_question = input("Votre question: ")
    if user_question.lower() in ['quit', 'exit']:
        break

    response = conversation_chain({"question": user_question})
    print(f"Réponse: {response['answer']}")
```

## 7 Considérations Techniques et Bonnes Pratiques

### 7.1 Optimisation des Performances

### 7.1.1 Gestion de la Taille des Chunks

#### Point important

Le choix de la taille des chunks impacte directement :

- **Précision** : Chunks plus petits = récupération plus précise
- **Contexte** : Chunks plus grands = plus de contexte par chunk
- **Performance** : Équilibre entre vitesse et qualité

Recommandations pratiques :

- Documents techniques : 500-1000 caractères
- Documents narratifs : 1000-2000 caractères
- Overlap : 10-20% de la taille du chunk

### 7.1.2 Stratégies de Recherche

**Recherche de similarité simple** : Retourne les K documents les plus similaires.

**Recherche avec seuil** : Retourne tous les documents au-dessus d'un seuil de similarité.

**Recherche hybride** : Combine recherche vectorielle et recherche par mots-clés.

## 7.2 Qualité des Embeddings

### 7.2.1 Choix du Modèle d'Embeddings

Le choix du modèle d'embeddings est crucial pour la performance du système RAG :

- **Modèles multilingues** : Pour traiter plusieurs langues
- **Modèles spécialisés** : Pour des domaines techniques spécifiques
- **Dimension des vecteurs** : Plus de dimensions = plus de nuances sémantiques

### 7.2.2 Cohérence des Modèles

#### Point important

Il est essentiel d'utiliser le même modèle d'embeddings pour :

- L'encodage des documents lors de l'ingestion
- L'encodage des requêtes lors de la recherche

Un changement de modèle nécessite une réindexation complète de la base de données vectorielle.

## 7.3 Surveillance et Évaluation

### 7.3.1 Métriques de Performance

**Précision de récupération** : Pourcentage de documents pertinents dans les résultats.

**Rappel** : Pourcentage de documents pertinents effectivement récupérés.

**Latence** : Temps de réponse du système complet.

**Satisfaction utilisateur** : Évaluation qualitative des réponses générées.

### 7.3.2 Tests et Validation

- **Jeux de test** : Création de questions-réponses de référence
- **Tests A/B** : Comparaison de différentes configurations
- **Retour utilisateur** : Collecte de feedback pour amélioration continue

## 8 Architecture et Intégration avec Oracle 23ai

### 8.1 Architecture Globale

#### 8.1.1 Composants Principaux

L'architecture RAG avec Oracle 23ai comprend plusieurs couches :

1. **Couche de données** : Oracle 23ai avec support vectoriel
2. **Couche d'embeddings** : OCI Generative AI ou modèles ONNX
3. **Couche d'orchestration** : LangChain pour la gestion des workflows
4. **Couche d'interface** : Applications utilisateur (web, mobile, API)

#### 8.1.2 Flux de Données

##### Exemple

Flux complet d'une requête RAG :

1. Réception de la requête utilisateur
2. Génération d'embedding pour la requête
3. Recherche vectorielle dans Oracle 23ai
4. Récupération des documents pertinents
5. Contextualisation avec l'historique de conversation
6. Génération de la réponse via OCI Generative AI
7. Retour de la réponse à l'utilisateur
8. Mise à jour de l'historique de conversation

### 8.2 Avantages d'Oracle 23ai

#### 8.2.1 Performance Native

Oracle 23ai offre des performances optimisées pour les opérations vectorielles :

- **Index HNSW natif** : Recherche vectorielle ultra-rapide
- **Parallélisation** : Traitement concurrent des requêtes
- **Cache intelligent** : Mise en cache des embeddings fréquemment utilisés

#### 8.2.2 Intégration SQL

La possibilité d'utiliser SQL pour les requêtes vectorielles simplifie l'intégration :

### Exemple

```
-- Recherche vectorielle avec SQL
SELECT content, VECTOR_DISTANCE(embedding, :query_vector) as similarity
FROM documents
ORDER BY similarity
FETCH FIRST 5 ROWS ONLY;
```

### 8.2.3 Sécurité et Gouvernance

Oracle 23ai apporte les avantages entreprise :

- **Chiffrement** : Données vectorielles chiffrées au repos et en transit
- **Contrôle d'accès** : Gestion fine des permissions
- **Audit** : Traçabilité complète des opérations
- **Sauvegarde** : Protection des données vectorielles

## 9 Cas d'Usage et Applications

### 9.1 Applications Métiers

#### 9.1.1 Support Client Intelligent

### Exemple

Un système de support client utilisant RAG peut :

- Accéder instantanément à la documentation produit
- Fournir des réponses personnalisées basées sur l'historique client
- Maintenir la cohérence des réponses across différents agents
- S'améliorer continuellement avec les nouvelles interactions

#### 9.1.2 Analyse Documentaire

Les systèmes RAG excellent dans l'analyse de gros volumes documentaires :

- **Recherche juridique** : Analyse de contrats et réglementations
- **Recherche médicale** : Consultation de bases de données médicales
- **Intelligence économique** : Analyse de rapports et études de marché

#### 9.1.3 Formation et Éducation

**Tuteurs virtuels** : Systèmes adaptatifs qui s'ajustent au niveau de l'apprenant.

**Bases de connaissances** : Accès intelligent aux ressources pédagogiques.

**Évaluation automatique** : Génération de questions et correction intelligente.

## 9.2 Défis et Limitations

### 9.2.1 Défis Techniques

#### Point important

**Latence** : Le processus de récupération + génération peut introduire des délais.

**Coût computationnel** : Les opérations d'embedding et de recherche vectorielle sont coûteuses.

**Maintenance** : Nécessité de maintenir les embeddings à jour avec les nouvelles données.

**Qualité des sources** : La qualité des réponses dépend directement de la qualité des documents source.

### 9.2.2 Considérations Éthiques

- **Biais dans les données** : Les documents source peuvent contenir des biais
- **Confidentialité** : Gestion appropriée des données sensibles
- **Transparence** : Capacité à expliquer les sources des réponses
- **Responsabilité** : Clarté sur la responsabilité des réponses générées

## 10 Évolutions et Perspectives

### 10.1 Tendances Technologiques

#### 10.1.1 Amélioration des Modèles

**Modèles multimodaux** : Intégration de texte, images, audio et vidéo dans le même système RAG.

**Embeddings contextuels** : Embeddings qui s'adaptent au contexte de la requête.

**Apprentissage continu** : Systèmes qui s'améliorent automatiquement avec l'usage.

#### 10.1.2 Optimisations Techniques

- **Quantification** : Réduction de la taille des embeddings sans perte de qualité
- **Compression** : Techniques avancées de compression vectorielle
- **Edge computing** : Déploiement de RAG sur des dispositifs locaux

### 10.2 Intégration Avancée avec Oracle

#### 10.2.1 Fonctionnalités Émergentes

**SELECT AI évolué** : Génération automatique de requêtes SQL complexes à partir de langage naturel.

**RAG fédéré** : Recherche simultanée dans multiple bases de données.

**Optimisation automatique** : Ajustement automatique des paramètres RAG selon l'usage.

### 10.2.2 Écosystème Oracle

L'intégration avec l'écosystème Oracle s'étend :

- **Oracle APEX** : Création rapide d'interfaces RAG
- **Oracle Analytics** : Visualisation des performances RAG
- **Oracle Integration** : Connexion avec systèmes externes

## 11 Conclusion

Le Retrieval Augmented Generation représente une évolution majeure dans l'utilisation des modèles de langage large, permettant de créer des applications intelligentes qui combinent la puissance génératrice des LLM avec l'accès à des bases de connaissances spécifiques et actualisées.

### 11.1 Points Clés à Retenir

Les éléments essentiels du RAG incluent :

- **Architecture en trois phases** : Ingestion, récupération, génération
- **Importance des embeddings** : Représentation sémantique des documents
- **Optimisation des performances** : Index vectoriels et stratégies de recherche
- **Contextualisation intelligente** : Utilisation du contexte conversationnel

### 11.2 Avantages de l'Approche Oracle

L'intégration avec Oracle 23ai et OCI Generative AI offre :

- **Performance enterprise** : Scalabilité et fiabilité
- **Sécurité renforcée** : Protection des données sensibles
- **Intégration native** : Facilité de déploiement dans les environnements existants
- **Flexibilité** : Support de multiples modèles et stratégies

### 11.3 Perspectives d'Avenir

Le domaine du RAG continue d'évoluer rapidement avec des innovations dans :

#### Point important

**Intelligence artificielle multimodale** : Intégration de différents types de médias dans le même système RAG.

**Personnalisation avancée** : Adaptation dynamique aux préférences et au contexte utilisateur.

**Efficacité énergétique** : Optimisation des ressources computationnelles nécessaires.

**Démocratisation** : Outils permettant à des non-experts de créer des systèmes RAG.

## 11.4 Recommandations Pratiques

Pour implémenter avec succès un système RAG :

1. **Commencer petit** : Débuter avec un cas d'usage spécifique et bien défini
2. **Prioriser la qualité des données** : Investir dans la curation et la structuration des documents source
3. **Tester itérativement** : Évaluer régulièrement les performances et ajuster les paramètres
4. **Planifier l'évolutivité** : Concevoir l'architecture pour supporter la croissance
5. **Former les utilisateurs** : Accompagner l'adoption par la formation et le support

Le RAG avec Oracle 23ai représente une opportunité unique de créer des applications intelligentes qui transforment la façon dont les organisations accèdent et utilisent leurs connaissances. En maîtrisant ces concepts et techniques, les développeurs et les data scientists peuvent construire des solutions qui apportent une valeur réelle aux utilisateurs finaux.

**Ce document a couvert les aspects fondamentaux du RAG avec Oracle 23ai. Pour approfondir ces concepts, nous recommandons la pratique avec des projets concrets et l'exploration continue des dernières évolutions technologiques dans ce domaine en rapide expansion.**