

Mini RAG Local

Retrieval-Augmented Generation sur vos documents PDF

Gemini Embedding 2 · Qdrant Vector DB · Gemini 2.5 Flash · PyMuPDF

2 scripts Python autonomes	Dual Texte + Image	100 % local & offline	1 024 dimensions
-----------------------------------	---------------------------	----------------------------------	-------------------------

Table des matières

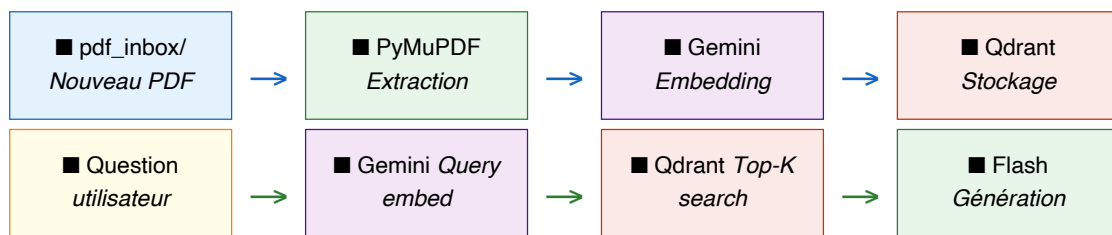
1. Vue d'ensemble et architecture
2. Pourquoi Gemini Embedding 2 ?
3. Installation complète (Conda + Qdrant)
4. Source — pdfInject.py
5. Source — rag.py
6. Exemples d'utilisation
7. Paramètres et configuration
8. Limites et pistes d'amélioration

■ 1 · Vue d'ensemble et architecture

Ce projet implémente un système **RAG (Retrieval-Augmented Generation)** entièrement local. Il vous permet d'interroger en langage naturel une collection de PDF stockés sur votre machine, sans envoyer vos documents à un cloud externe. Seuls les appels d'embedding et de génération transitent par l'API Google.

Composant	Rôle	Technologie	Local ?
Extraction PDF	Lire texte et raster chaque page	PyMuPDF (fitz)	■ Oui
Embedding	Convertir contenu en vecteurs	Gemini Embedding 2	API Google
Stockage	Index vectoriel + métadonnées	Qdrant (binaire local)	■ Oui
Génération	Réponse contextualisée NLP	Gemini 2.5 Flash	API Google
Orchestration	Pipeline complet en Python	pdfInject + rag.py	■ Oui

Flux de données



Haut : pipeline d'indexation (pdfInject.py) · Bas : pipeline de requête (rag.py)

■ 2 · Pourquoi Gemini Embedding 2 ?

Le choix du modèle d'embedding est crucial pour la qualité d'un système RAG. **Gemini Embedding 2** (gemini-embedding-2-preview / gemini-embedding-2 GA mars 2026) se distingue par plusieurs atouts majeurs, notamment sa capacité **multimodale native** qui permet d'utiliser *le même modèle* pour vectoriser du texte et des images dans un espace vectoriel commun.

Multimodalité native



Texte et images partagent le même espace vectoriel 1 024D. Une requête textuelle peut retrouver une page PDF riche en graphiques ou tableaux sans OCR supplémentaire. C'est l'atout principal pour ce projet.

Dimensionnalité configurable



L'API accepte output_dimensionality de 64 à 3 072. On choisit 1 024 pour un bon compromis précision / empreinte mémoire Qdrant.

Task types spécialisés



RETRIEVAL_DOCUMENT pour l'indexation, RETRIEVAL_QUERY pour les questions. Cela améliore la similarité asymétrique chunk ↔ requête courte.

MTEB State-of-the-art



Gemini Embedding 2 atteint des scores top-tier sur les benchmarks MTEB (Massive Text Embedding Benchmark) en langues multiples dont le français.

Multilingue



Supporte plus de 100 langues dans le même espace vectoriel : une requête en français retrouve un chunk en anglais ou en espagnol de façon transparente.

Intégration API unifiée



Un seul appel gemini.models.embed_content() gère texte et images. Pas de SDK distinct, pas de modèle de vision séparé à maintenir.

Comparaison avec d'autres modèles d'embedding populaires

Modèle	Dimensions	Multimodal	Multilingue	API unifiée
Gemini Embedding 2	64–3 072 (config.)	■ Texte + Image	■ 100+ langues	■ Oui
OpenAI text-embedding-3	256–3 072 (config.)	■ Texte seul	■ Partiel	■ Séparé
Cohere Embed v3	1 024 fixe	■ Texte + Image	■ Oui	■ Séparé
all-MiniLM-L6-v2	384 fixe	■ Texte seul	■ Limité	■ Local
nomic-embed-text	768 fixe	■ Texte seul	■ Limité	■ Local

■ **Point clé** : Le fait d'utiliser **le même espace vectoriel** pour le texte et les images signifie qu'une question comme « *Montre-moi le graphique de ventes* » peut retrouver une page PDF à dominante visuelle même si elle contient peu de texte extrait.

Avantage décisif pour les PDF à contenu mixte

La grande majorité des PDF professionnels ne sont pas de simples blocs de texte. Ils contiennent des **tableaux financiers**, des **graphiques**, des **schémas techniques**, des **captures d'écran**, des **diagrammes**, ou encore des pages entièrement scannées. Dans ces cas, l'extraction de texte brut par PyMuPDF est partielle ou inexistante — le contenu sémantique réside dans la *structure visuelle*, pas dans les caractères.

Type de contenu PDF	Texte extrait (PyMuPDF)	Avec Gemini Embedding 2
Page entièrement scannée	■ Rien	■ Image vectorisée et retrouvable
Tableau avec bordures	■ Texte décousu	■ Mise en page et valeurs comprises
Graphique / chart	■ Légendes seulement	■ Contenu visuel vectorisé
Schéma technique / plan	■ Titres uniquement	■ Structure visuelle indexée
PDF natif texte	■ Extraction complète	■ Texte + contexte visuel

■ **En pratique** : grâce à l'injection duale (texte *et* image par page), un PDF de rapport financier avec des tableaux complexes sera retrouvé même si PyMuPDF n'a extrait que des chiffres isolés. La vectorisation de la page rendue en PNG capture la **sémantique visuelle** — disposition, couleurs, formes — dans le même espace que votre question textuelle. C'est un avantage considérable sur les RAG classiques qui se limitent au texte extrait.

■ 3 · Installation complète

3.1 — Clé API Google (Gemini)

Le projet utilise l'API Google Gemini pour les embeddings et la génération. Vous bénéficiez de **300 \$ de crédits gratuits** avec un compte Google Cloud. La clé API est à définir comme variable d'environnement :

```
1 # Créer votre clé sur : https://aistudio.google.com/app/apikey
2 export GOOGLE_API_KEY='votre_cle_api_ici'
3
4 # Pour la rendre permanente (bash/zsh) :
5 echo 'export GOOGLE_API_KEY="votre_cle"' >> ~/.zshrc
6 source ~/.zshrc
```

3.2 — Environnement Python avec Conda

L'utilisation de Conda garantit un environnement isolé et reproductible, sans conflit avec d'autres projets Python.

```
1 # 1. Créer l'environnement Python 3.11
2 conda create -n rag_env python=3.11
3
4 # 2. Activer l'environnement
5 conda activate rag_env
6
7 # 3. Installer les dépendances
8 pip install pymupdf pillow qdrant-client google-genai
9
10 # 4. Vérifier l'installation
11 python -c "import fitz, qdrant_client, google.genai; print('OK')"
```

3.3 — Installation de Qdrant (binaire natif)

Qdrant est téléchargé et lancé directement comme binaire — pas besoin de Docker. Choisissez la version correspondant à votre architecture :


```
1 # Terminal 1 – Démarrer Qdrant
2 ./qdrant

3 # Terminal 2 – Démarrer le watcher d'indexation
4 conda activate rag_env
5 python pdfInject.py

6 # Terminal 3 – Interface de recherche
7 conda activate rag_env
8 python rag.py

9 # Utilisation : déposer des PDF dans pdf_inbox/
10 cp ~/Documents/rapport_annuel.pdf pdf_inbox/
11 cp ~/Documents/contrat.pdf pdf_inbox/
12 # -> pdfInject.py les détectera et indexera automatiquement
```


■ 6 · Exemples d'utilisation

■ Question

Quel est le chiffre d'affaires total du dernier exercice ?

■ Réponse

D'après le rapport annuel 2023, le chiffre d'affaires consolidé s'élève à 142,3 M€, en hausse de 8,4 % par rapport à l'exercice précédent (131,3 M€). La croissance est portée principalement par le segment Europe (+12 %).

■ Sources : rapport_2023.pdf — page 4 (score=0.91), rapport_2023.pdf — chunk 3 (score=0.88)

■ Question

Quelles sont les clauses de résiliation dans le contrat de prestation ?

■ Réponse

Le contrat prévoit deux modalités de résiliation : (1) résiliation à l'amiable avec un préavis de 30 jours par lettre recommandée, (2) résiliation pour faute grave sans préavis. Dans les deux cas, les travaux réalisés sont dus à la date d'effet de la résiliation.

■ Sources : contrat_prestation.pdf — chunk 7 (score=0.94), contrat_prestation.pdf — chunk 8 (score=0.87)

■ Question

Montre-moi les graphiques du tableau de bord Q3

■ Réponse

Je ne dispose pas du contenu textuel des graphiques, mais j'ai identifié 3 pages à fort contenu visuel correspondant à votre requête. Les pages 12, 15 et 18 semblent contenir les tableaux de bord Q3 selon leur similarité vectorielle.

■ Sources : dashboard_Q3.pdf — page 12 (score=0.89), page 15 (score=0.85), page 18 (score=0.82)

■ Question

Résume les points clés de la politique de confidentialité

■ Réponse

La politique de confidentialité couvre quatre points principaux : (1) collecte de données limitée au strict nécessaire, (2) durée de conservation de 3 ans maximum, (3) droit d'accès et de rectification sur simple demande écrite, (4) aucune transmission à des tiers sans consentement explicite.

■ Sources : politique_confidentialite.pdf — chunk 1 (score=0.96), chunk 2 (score=0.91)

Session interactive complète (terminal)

■ 7 · Paramètres et configuration

Paramètre	Défaut	Impact / Conseil
CHUNK_SIZE	1 000 car.	Augmenter (1 500–2 000) pour des documents denses. Diminuer (500) pour une
CHUNK_OVERLAP	200 car.	20 % de CHUNK_SIZE est une bonne règle. Évite de couper une phrase clé.
VECTOR_DIM	1 024	512 suffit pour des petites collections. 3 072 pour une précision maximale (+ coût)
TOP_K	5	Augmenter à 8–10 pour les questions complexes. Diminuer à 3 pour la vitesse.
WATCH_FOLDER	pdf_inbox/	Peut pointer vers n'importe quel dossier, y compris un dossier réseau monté.
EMBED_MODEL	gemini-embedding-2-preview	Remplacer par 'gemini-embedding-2' (GA) dès disponibilité générale.
GENERATE_MODEL	gemini-2.5-flash	Remplacer par 'gemini-2.5-pro' pour des réponses plus riches (+ latence).
sleep(30)	30 secondes	Réduire à 5–10 s pour une indexation plus réactive (+ charge CPU/API).

■ 8 · Limites et pistes d'amélioration

◆ Images référencées, non injectées

Les hits image pointent vers la page source mais n'envoient pas le PNG à Gemini 2.5 Flash. Évolution : passer l'image directement à Gemini 2.5 Flash (vision multimodale) pour des réponses sur le contenu visuel.

◆ Polling vs. événements système

Le sleep(30) consomme des cycles inutilement. Sur Linux, inotifywait (inotify-tools) ou la lib Python watchdog permettrait une réactivité instantanée sur nouveau fichier.

◆ Collection recrée au démarrage

Supprimer et recréer la collection à chaque lancement détruit l'index existant. En prod : supprimer ce bloc et utiliser upsert pour une indexation incrémentale vraiment persistante.

◆ Pas d'interface graphique

Une UI Streamlit ou Gradio améliorerait l'ergonomie : upload drag-and-drop, affichage des scores, historique des échanges, prévisualisation de la page source.

◆ Pas de re-ranking

Après la recherche vectorielle (recall), un re-ranking cross-encoder (ex: Cohere Rerank) améliorerait la précision des résultats avant injection dans le prompt.

◆ Un seul type de chunk

Tenter un chunking sémantique (sentences, paragraphes) plutôt que caractères fixes améliorerait la cohérence des chunks et donc la qualité des réponses générées.

