

Guide Complet pour Débutants: Création d'un Chatbot Juridique avec Groq API, ChromaDB et MinLM

GHaythem Ghazouani

May 15, 2025

Contents

1	Introduction	3
2	Installation de l'environnement de développement	3
2.1	Installation de VSCode	3
2.2	Extensions VSCode recommandées	3
2.3	Installation de Python	3
3	Création du projet	4
3.1	Création du dossier du projet	4
3.2	Création de l'environnement virtuel	4
3.3	Installation des dépendances	4
3.4	Création de la structure du projet	4
4	Configuration de l'environnement	5
4.1	Configuration des variables d'environnement	5
5	Implémentation du chargeur de PDF	5
6	Implémentation du modèle de chatbot	7
7	Création des routes API	9
8	Création de l'application principale	11
9	Création du fichier principal	12
10	Exécution de l'application	12
10.1	Lancement du serveur	12
10.2	Accès à l'interface Swagger	12

11 Test du chatbot	13
11.1 Téléchargement d'un document PDF	13
11.2 Poser une question	13
12 Dépannage	13
12.1 Problèmes courants et solutions	13
13 Améliorations possibles	14
14 Conclusion	14

1 Introduction

Ce guide vous montrera, étape par étape, comment créer un chatbot juridique qui utilise l'API Groq avec LangChain, ChromaDB et le modèle MinLM pour répondre à des questions juridiques basées sur un document PDF. Nous utiliserons FastAPI pour créer une API REST et Swagger pour la tester. Ce guide est conçu pour les débutants et inclut toutes les commandes à exécuter dans VSCode.

2 Installation de l'environnement de développement

2.1 Installation de VSCode

Si vous n'avez pas encore VSCode, téléchargez-le depuis <https://code.visualstudio.com/> et installez-le.

2.2 Extensions VSCode recommandées

Ouvrez VSCode et installez les extensions suivantes pour faciliter le développement:

1. Python (Microsoft)
2. Pylance (Microsoft)
3. Python Indent (Kevin Rose)
4. autoDocstring (Nils Werner)

Pour installer une extension:

1. Cliquez sur l'icône des extensions dans la barre latérale (ou appuyez sur Ctrl+Shift+X)
2. Recherchez l'extension
3. Cliquez sur "Installer"

2.3 Installation de Python

Si vous n'avez pas Python installé:

1. Téléchargez Python 3.8 ou supérieur depuis <https://www.python.org/downloads/>
2. Lors de l'installation, cochez la case "Add Python to PATH"
3. Vérifiez l'installation en ouvrant un terminal dans VSCode (Terminal → New Terminal) et tapez:

```
1 python --version
```

3 Création du projet

3.1 Création du dossier du projet

Ouvrez VSCode et suivez ces étapes:

1. Cliquez sur "File" → "Open Folder"
2. Créez un nouveau dossier nommé "legal_chatbot" et ouvrez-le
3. Ouvrez un terminal dans VSCode (Terminal → New Terminal)

3.2 Création de l'environnement virtuel

Dans le terminal VSCode, exécutez les commandes suivantes:

```
1 # Création de l'environnement virtuel
2 python -m venv venv
3
4 # Activation de l'environnement virtuel
5 # Sur Windows
6 venv\Scripts\activate
7 # Sur macOS/Linux
8 source venv/bin/activate
```

Vous devriez voir "(venv)" apparaître au début de la ligne de commande, indiquant que l'environnement virtuel est activé.

3.3 Installation des dépendances

Toujours dans le terminal VSCode avec l'environnement virtuel activé, installez les dépendances nécessaires:

```
1 # Installation des packages requis
2 pip install fastapi uvicorn python-multipart python-dotenv
3 pip install langchain langchain-groq
4 pip install chromadb sentence-transformers
5 pip install pypdf
```

3.4 Création de la structure du projet

Exécutez les commandes suivantes pour créer la structure de base du projet:

```
1 # Sur Windows
2 mkdir app
3 mkdir app\routes
4 mkdir models
5 mkdir utils
6 mkdir data
7 type nul > app\__init__.py
8 type nul > app\main.py
9 type nul > app\routes\__init__.py
10 type nul > app\routes\chatbot.py
11 type nul > models\__init__.py
```

```

12 type nul > models\chatbot.py
13 type nul > utils\__init__.py
14 type nul > utils\pdf_loader.py
15 type nul > .env
16
17 # Sur macOS/Linux
18 mkdir -p app/routes
19 mkdir -p models
20 mkdir -p utils
21 mkdir -p data
22 touch app/__init__.py
23 touch app/main.py
24 touch app/routes/__init__.py
25 touch app/routes/chatbot.py
26 touch models/__init__.py
27 touch models/chatbot.py
28 touch utils/__init__.py
29 touch utils/pdf_loader.py
30 touch .env

```

4 Configuration de l'environnement

4.1 Configuration des variables d'environnement

Dans VSCode, ouvrez le fichier '.env' et ajoutez votre clé API Groq:

```
1 GROQ_API_KEY=votre_cle_api_groq
```

Pour obtenir une clé API Groq:

1. Créez un compte sur <https://console.groq.com/>
2. Accédez à la section API Keys
3. Créez une nouvelle clé API
4. Copiez la clé et collez-la dans le fichier .env

5 Implémentation du chargeur de PDF

Ouvrez le fichier 'utils/pdfloader.py' dans VSCode et ajoutez le code suivant :

```

1 # utils/pdf_loader.py
2 from langchain_community.document_loaders import PyPDFLoader
3 from langchain.text_splitter import RecursiveCharacterTextSplitter
4 from langchain_community.vectorstores import Chroma
5 from langchain_community.embeddings import HuggingFaceEmbeddings
6 import os
7 from dotenv import load_dotenv
8
9 # Charger les variables d'environnement
10 load_dotenv()
11
12 def load_and_process_pdf(pdf_path):

```

```

13 """
14 Charge un document PDF, le divise en chunks et crée une base de
    données vectorielle avec ChromaDB.
15
16 Args:
17     pdf_path (str): Chemin vers le fichier PDF à traiter
18
19 Returns:
20     Chroma: Base de données vectorielle contenant les chunks du
        document
21 """
22 print(f"Chargement du PDF: {pdf_path}")
23
24 # tape 1: Charger le PDF avec PyPDFLoader
25 loader = PyPDFLoader(pdf_path)
26 documents = loader.load()
27 print(f"PDF chargé avec succès: {len(documents)} pages")
28
29 # tape 2: Diviser le texte en chunks plus petits pour un
    meilleur traitement
30 text_splitter = RecursiveCharacterTextSplitter(
31     chunk_size=1000, # Taille de chaque chunk en caractères
32     chunk_overlap=200, # Chevauchement entre les chunks pour
        maintenir le contexte
33     length_function=len
34 )
35 chunks = text_splitter.split_documents(documents)
36 print(f"Document divisé en {len(chunks)} chunks")
37
38 # tape 3: Initialiser le modèle d'embedding MinLM
39 # Nous utilisons le modèle all-MiniLM-L6-v2 qui offre un bon é
    quilibre entre performance et rapidité
40 embeddings = HuggingFaceEmbeddings(
41     model_name="sentence-transformers/all-MiniLM-L6-v2",
42     model_kwargs={'device': 'cpu'} # Utiliser CPU pour la
        compatibilité
43 )
44 print("Modèle d'embedding MinLM initialisé")
45
46 # tape 4: Créer une base de données vectorielle avec ChromaDB
47 # ChromaDB stockera les embeddings et permettra des recherches
    sémantiques efficaces
48 db_directory = "./data/chroma_db"
49 vector_store = Chroma.from_documents(
50     documents=chunks,
51     embedding=embeddings,
52     persist_directory=db_directory # Stockage persistant des
        embeddings
53 )
54
55 # Persister la base de données pour une utilisation future
56 vector_store.persist()
57 print(f"Base de données vectorielle ChromaDB créée et persistée
    dans {db_directory}")
58
59 return vector_store

```

6 Implémentation du modèle de chatbot

Ouvrez le fichier 'models/chatbot.py' dans VSCode et ajoutez le code suivant:

```
1 # models/chatbot.py
2 from langchain.chains import ConversationalRetrievalChain
3 from langchain.memory import ConversationBufferMemory
4 from langchain_groq import ChatGroq
5 import os
6 from dotenv import load_dotenv
7
8 # Charger les variables d'environnement
9 load_dotenv()
10
11 class LegalChatbot:
12     def __init__(self, vector_store):
13         """
14         Initialise le chatbot juridique avec une base de données
15         vectorielle.
16
17         Args:
18             vector_store (Chroma): Base de données vectorielle
19             ChromaDB contenant les documents
20         """
21         print("Initialisation du chatbot juridique...")
22
23         # Stocker la base de données vectorielle
24         self.vector_store = vector_store
25
26         # Créer une mémoire pour stocker l'historique de
27         # conversation
28         self.memory = ConversationBufferMemory(
29             memory_key="chat_history",
30             return_messages=True
31         )
32
33         # Vérifier que la clé API Groq est disponible
34         groq_api_key = os.environ.get("GROQ_API_KEY")
35         if not groq_api_key:
36             raise ValueError("La clé API Groq n'est pas définie. Vérifiez votre fichier .env")
37
38         # Définir le prompt système pour guider le modèle
39         SYSTEM_PROMPT = """
40         Tu es un assistant juridique expert. Ta tâche est de répondre aux questions juridiques en te basant uniquement sur les documents fournis.
41
42         Voici quelques règles à suivre:
43         1. Réponds uniquement aux questions juridiques basées sur les documents fournis.
44         2. Si tu ne connais pas la réponse ou si l'information n'est pas dans les documents, indique-le clairement.
45         3. Cite les sections spécifiques des documents pour appuyer tes réponses.
46         4. Utilise un langage clair et précis, mais reste technique lorsque c'est nécessaire.
```

```

44     5. N'invente pas d'informations qui ne sont pas présentes
      dans les documents.
45
46     Format de réponse:
47     - Commence par une réponse directe à la question
48     - Développe avec des détails pertinents
49     - Cite les sources spécifiques (numéros de page, sections)
50     - Si nécessaire, mentionne les limitations de ta réponse
51     """
52
53     # Initialiser le modèle LLM avec Groq
54     self.llm = ChatGroq(
55         groq_api_key=groq_api_key,
56         model_name="llama3-70b-8192", # Modèle LLM puissant
      pour les réponses juridiques
57         temperature=0.2, # Température basse pour des réponses
      plus précises
58         system=SYSTEM_PROMPT # Prompt système pour guider le
      modèle
59     )
60     print("Modèle LLM Groq initialisé")
61
62     # Créer la chaîne de conversation avec récupération (RAG)
63     self.chain = ConversationalRetrievalChain.from_llm(
64         llm=self.llm,
65         retriever=self.vector_store.as_retriever(
66             search_kwargs={"k": 3} # Récupérer les 3 chunks
      les plus pertinents
67         ),
68         memory=self.memory,
69         return_source_documents=True # Retourner les documents
      sources pour les citations
70     )
71     print("Chaîne de conversation RAG initialisée")
72
73     async def get_response(self, query):
74         """
75         Obtient une réponse du chatbot pour une requête donnée.
76
77         Args:
78             query (str): La question posée par l'utilisateur
79
80         Returns:
81             dict: Dictionnaire contenant la réponse et les sources
82         """
83         print(f"Traitement de la requête: {query}")
84
85         # Invoquer la chaîne de conversation de manière asynchrone
86         result = await self.chain.ainvoke({"question": query})
87
88         # Extraire les sources pour la citation
89         sources = []
90         for doc in result.get("source_documents", []):
91             if hasattr(doc, "metadata") and "source" in doc.
      metadata:
92                 sources.append(doc.metadata["source"])

```



```

93         elif hasattr(doc, "metadata") and "page" in doc.
metadata:
94             sources.append(f"Page {doc.metadata['page']}")
95
96             # Retourner la réponse et les sources
97             response = {
98                 "answer": result["answer"],
99                 "sources": list(set(sources)) # limiter les doublons
100             }
101
102             print("Réponse générée avec succès")
103             return response

```

7 Création des routes API

Ouvrez le fichier ‘app/routes/chatbot.py’ dans VSCode et ajoutez le code suivant:

```

1 # app/routes/chatbot.py
2 from fastapi import APIRouter, UploadFile, File, HTTPException
3 from fastapi.responses import JSONResponse
4 import os
5 import tempfile
6 from pydantic import BaseModel
7 from typing import List, Optional
8 from models.chatbot import LegalChatbot
9 from utils.pdf_loader import load_and_process_pdf
10
11 # Créer un router FastAPI
12 router = APIRouter(tags=["chatbot"])
13
14 # Modèles de données pour la validation des entrées/sorties
15 class Query(BaseModel):
16     text: str
17
18     class Config:
19         schema_extra = {
20             "example": {
21                 "text": "Quelles sont les obligations légales d'un
employeur?"
22             }
23         }
24
25 class Response(BaseModel):
26     answer: str
27     sources: Optional[List[str]] = None
28
29 # Variable globale pour stocker l'instance du chatbot
30 chatbot_instance = None
31
32 @router.post("/upload-pdf", response_class=JSONResponse, summary="T
élécharger un document PDF")
33 async def upload_pdf(file: UploadFile = File(...)):
34     """
35     Télécharge un document PDF et initialise le chatbot juridique.

```

```

36
37 - **file**: Fichier PDF à télécharger (format .pdf uniquement)
38
39 Retourne un message de confirmation si le chargement est réussi
40
41 """
42 global chatbot_instance
43
44 # Vérifier que le fichier est un PDF
45 if not file.filename.endswith(".pdf"):
46     raise HTTPException(
47         status_code=400,
48         detail="Le fichier doit être au format PDF (.pdf)"
49     )
50
51 # Créer un fichier temporaire pour stocker le PDF
52 with tempfile.NamedTemporaryFile(delete=False, suffix=".pdf")
53 as temp_file:
54     # Lire le contenu du fichier téléchargé et l'écrire dans le
55     # fichier temporaire
56     content = await file.read()
57     temp_file.write(content)
58     temp_file_path = temp_file.name
59
60 try:
61     # Charger et traiter le PDF avec notre utilitaire
62     vector_store = load_and_process_pdf(temp_file_path)
63
64     # Initialiser le chatbot avec la base de données
65     vectorielle
66     chatbot_instance = LegalChatbot(vector_store)
67
68     return {
69         "message": "Document PDF chargé avec succès et chatbot
70         initialisé",
71         "filename": file.filename
72     }
73 except Exception as e:
74     # En cas d'erreur, renvoyer un message d'erreur détaillé
75     raise HTTPException(
76         status_code=500,
77         detail=f"Erreur lors du traitement du PDF: {str(e)}"
78     )
79 finally:
80     # Supprimer le fichier temporaire pour libérer de l'espace
81     if os.path.exists(temp_file_path):
82         os.unlink(temp_file_path)
83
84 @router.post("/ask", response_model=Response, summary="Poser une
85 question au chatbot")
86 async def ask_question(query: Query):
87     """
88     Pose une question au chatbot juridique et obtient une réponse.
89
90     - **query**: Objet contenant le texte de la question
91
92     Retourne la réponse du chatbot et les sources utilisées.

```

```

87     """
88     global chatbot_instance
89
90     # Vérifier que le chatbot a été initialisé
91     if chatbot_instance is None:
92         raise HTTPException(
93             status_code=400,
94             detail="Veuillez d'abord télécharger un document PDF
via /upload-pdf"
95         )
96
97     try:
98         # Obtenir une réponse du chatbot
99         response = await chatbot_instance.get_response(query.text)
100         return response
101     except Exception as e:
102         # En cas d'erreur, renvoyer un message d'erreur détaillé
103         raise HTTPException(
104             status_code=500,
105             detail=f"Erreur lors de la génération de la réponse: {
str(e)}"
106         )

```

8 Création de l'application principale

Ouvrez le fichier 'app/main.py' dans VSCode et ajoutez le code suivant:

```

1 # app/main.py
2 from fastapi import FastAPI
3 from fastapi.middleware.cors import CORSMiddleware
4 from app.routes.chatbot import router as chatbot_router
5
6 # Créer l'application FastAPI
7 app = FastAPI(
8     title="Chatbot Juridique API",
9     description="API pour un chatbot juridique utilisant Groq,
ChromaDB et MinLM",
10     version="1.0.0",
11     docs_url="/docs", # URL pour la documentation Swagger
12     redoc_url="/redoc" # URL pour la documentation ReDoc
13 )
14
15 # Configuration CORS pour permettre les requêtes cross-origin
16 app.add_middleware(
17     CORSMiddleware,
18     allow_origins=["*"], # Autoriser toutes les origines
19     allow_credentials=True,
20     allow_methods=["*"], # Autoriser toutes les méthodes
21     allow_headers=["*"], # Autoriser tous les headers
22 )
23
24 # Ajouter les routes du chatbot
25 app.include_router(chatbot_router, prefix="/api")
26
27 # Route racine
28 @app.get("/", tags=["root"])

```

```

29 async def root():
30     """
31     Route racine qui renvoie un message de bienvenue.
32     """
33     return {
34         "message": "Bienvenue sur l'API du Chatbot Juridique",
35         "documentation": "/docs",
36         "endpoints": {
37             "upload_pdf": "/api/upload-pdf",
38             "ask_question": "/api/ask"
39         }
40     }

```

9 Création du fichier principal

Créez un fichier ‘main.py’ à la racine du projet et ajoutez le code suivant:

```

1 # main.py
2 import uvicorn
3
4 if __name__ == "__main__":
5     # Démarrer le serveur uvicorn
6     uvicorn.run(
7         "app.main:app", # Chemin vers l'application FastAPI
8         host="0.0.0.0", # écouter sur toutes les interfaces
9         port=8000,      # Port 8000
10        reload=True      # Recharger automatiquement lors des
11        modifications
12    )

```

10 Exécution de l'application

10.1 Lancement du serveur

Dans le terminal VSCode, assurez-vous que l'environnement virtuel est activé, puis exécutez:

```

1 # Exécuter l'application
2 python main.py

```

Vous devriez voir une sortie similaire à:

```

1 INFO:      Started server process [12345]
2 INFO:      Waiting for application startup.
3 INFO:      Application startup complete.
4 INFO:      Uvicorn running on http://0.0.0.0:8000 (Press CTRL+C to
quit)

```

10.2 Accès à l'interface Swagger

Ouvrez votre navigateur et accédez à:

```

1 http://localhost:8000/docs

```

Vous verrez l'interface Swagger qui vous permettra de tester votre API.

11 Test du chatbot

11.1 Téléchargement d'un document PDF

Dans l'interface Swagger:

1. Cliquez sur l'endpoint `/api/upload-pdf`
2. Cliquez sur "Try it out"
3. Cliquez sur "Choose File" et sélectionnez un document PDF juridique
4. Cliquez sur "Execute"

11.2 Poser une question

Après avoir téléchargé un PDF:

1. Cliquez sur l'endpoint `/api/ask`
2. Cliquez sur "Try it out"
3. Dans le champ "Request body", entrez votre question:

```
1 {  
2   "text": "Quelles sont les principales obligations lé  
3   gales mentionnées dans le document?"  
4 }
```

4. Cliquez sur "Execute"

12 Dépannage

12.1 Problèmes courants et solutions

1. Erreur "ModuleNotFoundError"

Solution: Vérifiez que vous avez installé toutes les dépendances et que l'environnement virtuel est activé.

```
1 pip install -r requirements.txt  
2
```

2. Erreur "API key not found"

Solution: Vérifiez que votre fichier `.env` contient la clé API Groq correcte.

3. Erreur lors du chargement du PDF

Solution: Vérifiez que le PDF n'est pas corrompu et qu'il est lisible.

4. Erreur "Address already in use"

Solution: Changez le port dans le fichier main.py ou arrêtez le processus qui utilise le port 8000.

```
1 # Sur Windows
2 netstat -ano | findstr :8000
3 taskkill /PID <PID> /F
4
5 # Sur macOS/Linux
6 lsof -i :8000
7 kill -9 <PID>
8
```

13 Améliorations possibles

Voici quelques idées pour améliorer votre chatbot:

1. Interface utilisateur web

Créez une interface web simple avec HTML, CSS et JavaScript pour interagir avec votre API.

2. Support de plusieurs documents

Modifiez le code pour permettre le chargement de plusieurs documents PDF.

3. Historique des conversations

Ajoutez une base de données pour stocker l'historique des conversations.

4. Authentification

Ajoutez une authentification pour sécuriser votre API.

5. Déploiement

Déployez votre application sur un service comme Heroku, Vercel ou AWS.

14 Conclusion

Félicitations! Vous avez créé un chatbot juridique qui utilise l'API Groq avec LangChain, ChromaDB et le modèle MinLM pour répondre à des questions basées sur un document PDF. Ce chatbot peut être utilisé pour fournir des informations juridiques basées sur des documents spécifiques.

Ce projet vous a permis d'apprendre:

- Comment utiliser FastAPI pour créer une API REST
- Comment traiter des documents PDF avec PyPDFLoader
- Comment utiliser ChromaDB pour stocker et rechercher des embeddings

- Comment utiliser le modèle MinLM pour créer des embeddings
- Comment utiliser l'API Groq pour générer des réponses
- Comment implémenter la technique RAG (Retrieval Augmented Generation)