

Get 50% off your ticket to MongoDB.local NYC on May 2. Use code Web50! >>



MongoDB Developer



[MONGODB DEVELOPER CENTER](#) > [DEVELOPER TOPICS](#) > [PRODUCTS](#) >

[MONGODB](#) > [TUTORIALS](#)

# Revolutionizing AI Interaction: Integrating Mistral AI and MongoDB for a Custom LLM GenAI Application



Han Heloir

11 min read • Published Feb 08, 2024 • Updated Feb 13, 2024

MongoDB



Rate this tutorial

Large language models (LLMs) are known for their ability to converse with us in an almost human-like manner. Yet, the complexity of their inner workings often remains covered in mystery, sparking intrigue. This intrigue intensifies when we factor in the privacy challenges associated with AI technologies.

In addition to privacy concerns, cost is another significant challenge. Deploying a large language model is crucial for AI applications, and there are two primary options: self-hosted or API-based models. With API-based LLMs, the model is hosted by a service provider, and costs accrue with each API request. In contrast, a self-hosted LLM runs on your own infrastructure, giving you complete control over costs. The bulk of expenses for a self-hosted LLM pertains to the necessary hardware.

Another aspect to consider is the availability of LLM models. With API-based models, during times of high demand, model availability can be compromised. In contrast, managing your own LLM ensures control over availability. You will be able to make sure all your queries to your self-managed LLM can be handled properly and under your control.

**Mistral AI**, a French startup, has introduced innovative solutions with the Mistral 7B model, Mistral Mixture of Experts, and Mistral Platform, all standing for a spirit of openness. This article explores how Mistral AI, in collaboration with MongoDB, a developer data platform that unifies operational, analytical, and vector search data services, is revolutionizing our interaction with AI. We will delve into the integration of Mistral AI with MongoDB Atlas and discuss its impact on privacy, cost efficiency, and AI accessibility.

## Mistral AI: a game-changer

Mistral AI has emerged as a pivotal player in the open-source AI community, setting new standards in AI innovation. Let's break down what makes Mistral AI so transformative.

## A beacon of openness: Mistral AI's philosophy

Mistral AI's commitment to openness is at the core of its philosophy. This commitment extends beyond just providing open-source code; it's about advocating for transparent and adaptable AI models. By prioritizing transparency, Mistral AI empowers users to truly own and shape the future of AI. This approach is fundamental to ensuring AI remains a positive, accessible force for everyone.

## Unprecedented performance with Mistral 8x7B

Mistral AI has taken a monumental leap forward with the release of Mixtral 8x7B, an innovative sparse mixture of experts model (SMoE) with open weights. An SMoE is a neural network architecture that boosts traditional model efficiency and scalability. It utilizes specialized “expert” sub-networks to handle different input segments. Mixtral incorporates eight of these expert sub-networks.

Licensed under Apache 2.0, Mixtral sets a new benchmark in the AI landscape. Here's a closer look at what makes Mixtral 8x7B a groundbreaking advancement.

## High-performance with sparse architectures

Mixtral 8x7B stands out for its efficient utilization of parameters and high-quality performance. Despite its total parameter count of 46.7 billion, it operates using only 12.9 billion parameters per token. This unique architecture allows Mixtral to maintain the speed and cost efficiency of a 12.9 billion parameter model while offering the capabilities of a much larger model.

## Superior performance, versatility, and cost-performance optimization

Mixtral rivals leading models like Llama 2 70B and GPT-3.5, excelling in handling large contexts, multilingual processing, code generation, and instruction-following. The Mixtral 8x7B model combines cost efficiency with high performance, using a sparse mixture of experts network for optimized resource usage, offering premium outputs at lower costs compared to similar models.

## Mistral “La plateforme”

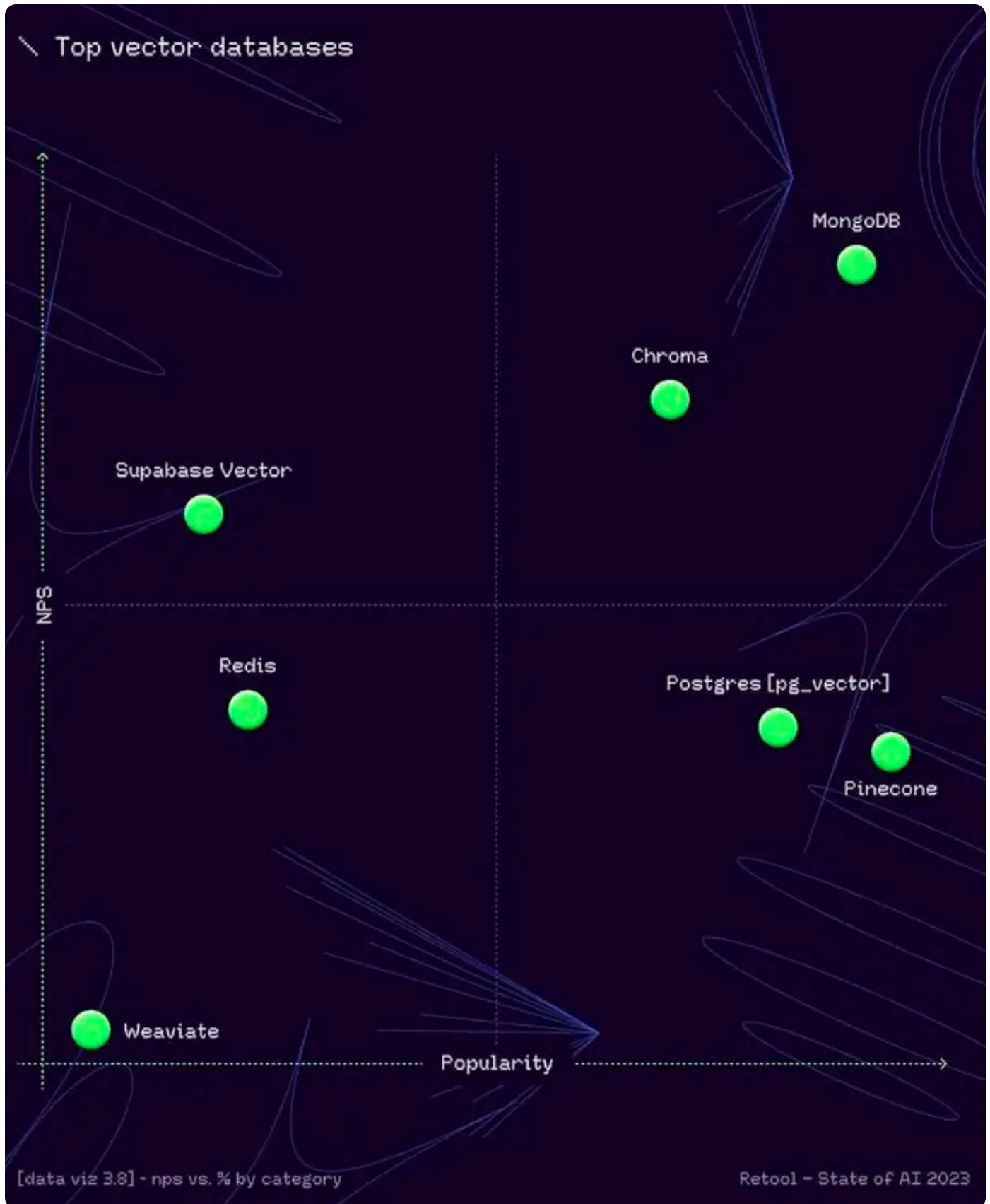
Mistral AI's beta platform offers developers generative models focusing on simplicity: Mistral-tiny for cost-effective, English-only text generation (7.6 MT-Bench score), Mistral-small for multilingual support including coding (8.3 score), and Mistral-medium for high-quality, multilingual output (8.6 score). These user-friendly, accurately fine-tuned models facilitate efficient AI deployment, as demonstrated in our article using the Mistral-tiny and the platform's embedding model.

## Why MongoDB Atlas as a vector store?

[MongoDB Atlas](#) is a unique, fully-managed platform integrating enterprise data, vector search, and analytics, allowing the creation of tailored AI applications. It goes beyond standard vector search with a comprehensive ecosystem, including models like Mistral, setting it apart in terms of unification, scalability, and security.

MongoDB Atlas unifies operational, analytical, and vector search data services to streamline the building of generative AI-enriched apps. From proof-of-concept to production, MongoDB Atlas empowers developers with scalability, security, and performance for their mission-critical production applications.

According to the [Retool AI report](#), MongoDB takes the lead, earning its place as the top-ranked vector database.



- Vector store easily works together with current MongoDB databases, making it a simple addition for groups already using MongoDB for managing their

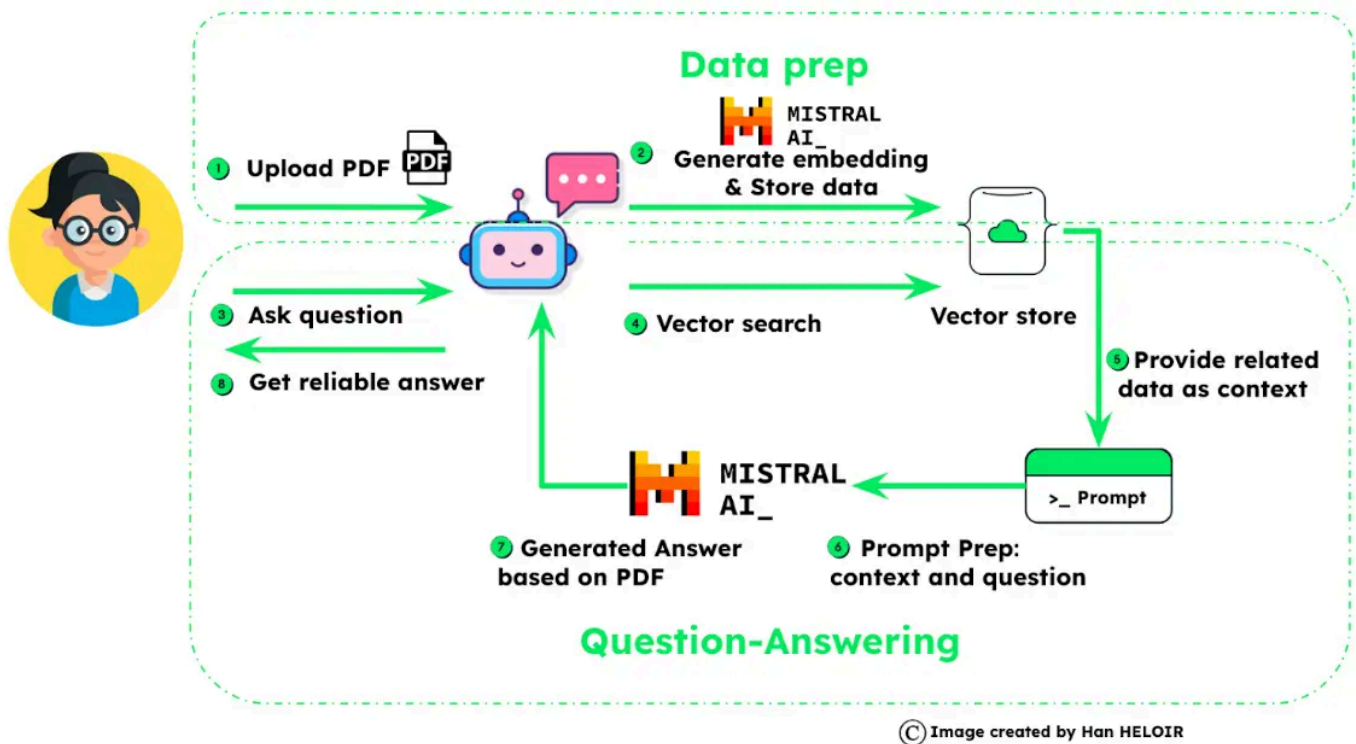
data. This means they can start using vector storage without needing to make big changes to their systems.

- MongoDB Atlas is purpose-built to handle large-scale, operation-critical applications, showcasing its robustness and reliability. This is especially important in applications where it's critical to have accurate and accessible data.
- Data in MongoDB Atlas is stored in JSON format, making it an ideal choice for managing a variety of data types and structures. This is particularly useful for AI applications, where the data type can range from embeddings and text to integers, floating-point values, GeoJSON, and more.
- MongoDB Atlas is designed for enterprise use, featuring top-tier security, the ability to operate across multiple cloud services, and is fully managed. This ensures organizations can trust it for secure, reliable, and efficient operations.

With MongoDB Atlas, organizations can confidently store and retrieve embeddings alongside your existing data, unlocking the full potential of AI for their applications.

## Overview and implementation of your custom LLM GenAI app

Creating a self-hosted LLM GenAI application integrates the power of open-source AI with the robustness of an enterprise-grade vector store like MongoDB. Below is a detailed step-by-step guide to implementing this innovative system:



## 1. Data acquisition and chunk

The first step is gathering data relevant to your application's domain, including text documents, web pages, and importantly, operational data already stored in MongoDB Atlas. Leveraging Atlas's operational data adds a layer of depth, ensuring your AI application is powered by comprehensive, real-time data, which is crucial for contextually enriched AI responses.

Then, we divide the data into smaller, more manageable chunks. This division is crucial for efficient data processing, guaranteeing the AI model interacts with data that is both precise and reflective of your business's operational context.

### 2.1 Generating embeddings

Utilize **Mistral AI embedding endpoint** to transform your segmented text data into embeddings. These embeddings are numerical representations that capture the essence of your text, making it understandable and usable by AI models.

### 2.2 Storing embeddings in MongoDB vector store

Once you have your embeddings, store them in MongoDB's vector store. MongoDB Atlas, with its advanced search capabilities, allows for the efficient storing and managing of these embeddings, ensuring that they are easily accessible when needed.

## 2.3 Querying your data

Use **MongoDB's vector search** capability to query your stored data. You only need to create a vector search index on the embedding field in your document. This powerful feature enables you to perform complex searches and retrieve the most relevant pieces of information based on your query parameters.

## 3. & 4. Embedding questions and retrieving similar chunks

When a user poses a question, generate an embedding for this query. Then, using MongoDB's search functionality, retrieve data chunks that are most similar to this query embedding. This step is crucial for finding the most relevant information to answer the user's question.

## 5. Contextualized prompt creation

Combine the retrieved segments and the original user query to create a comprehensive prompt. This prompt will provide a context to the AI model, ensuring that the responses generated are relevant and accurate.

## 6. & 7. Customized answer generation from Mistral AI

Feed the contextualized prompt into the Mistral AI 7B LLM. The model will then generate a customized answer based on the provided context. This step leverages the advanced capabilities of Mistral AI to provide specific, accurate, and relevant answers to user queries.



# Implementing a custom LLM GenAI app with Mistral AI and MongoDB Atlas

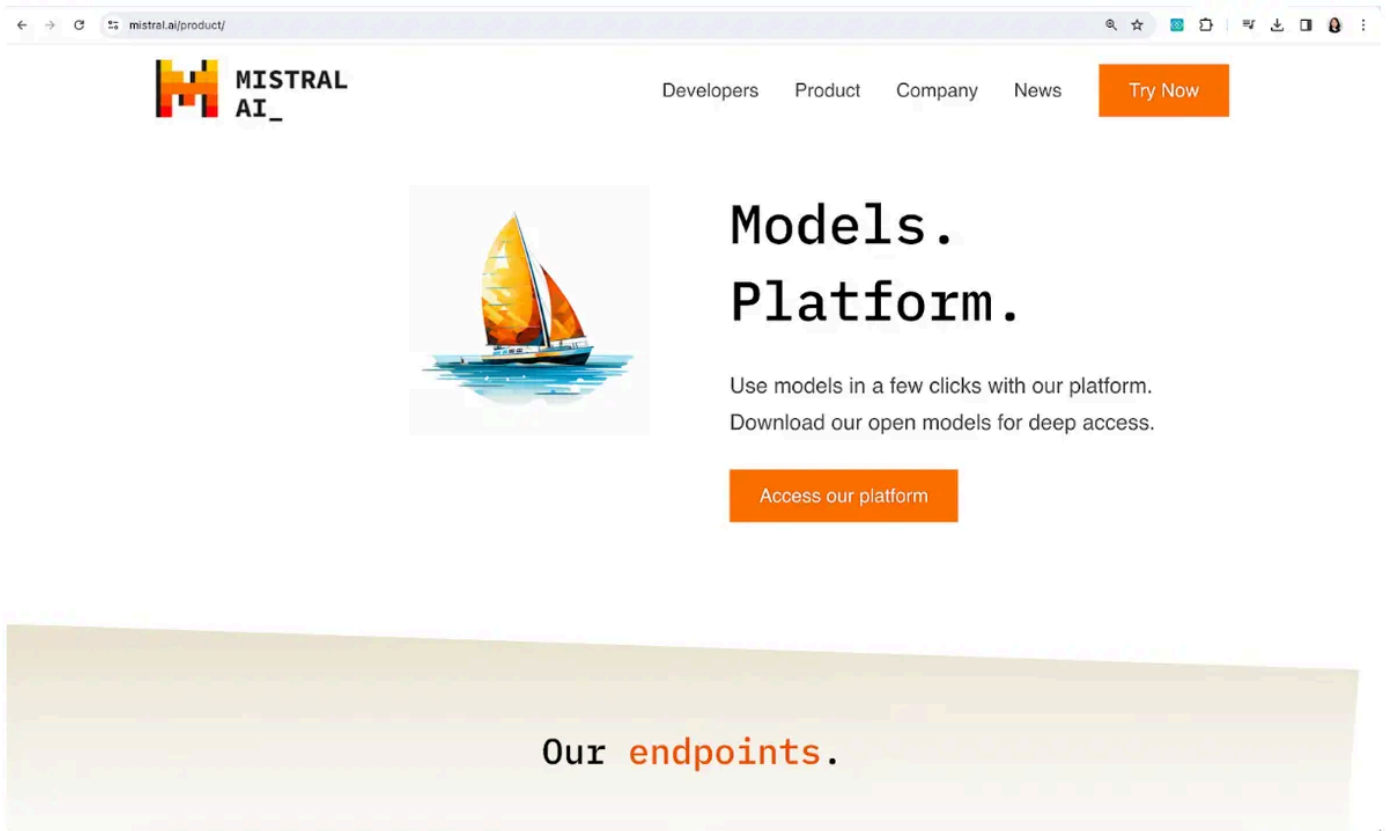
Now that we have a comprehensive understanding of Mistral AI and MongoDB Atlas and the overview of your next custom GenAI app, let's dive into implementing a custom large language model GenAI app. This app will allow you to have your own personalized AI assistant, powered by the Mistral AI and supported by the efficient data management of MongoDB Atlas.

In this section, we'll explain the prerequisites and four parts of the code:

- Needed libraries
- Data preparation process
- Question and answer process
- User interface through Gradio

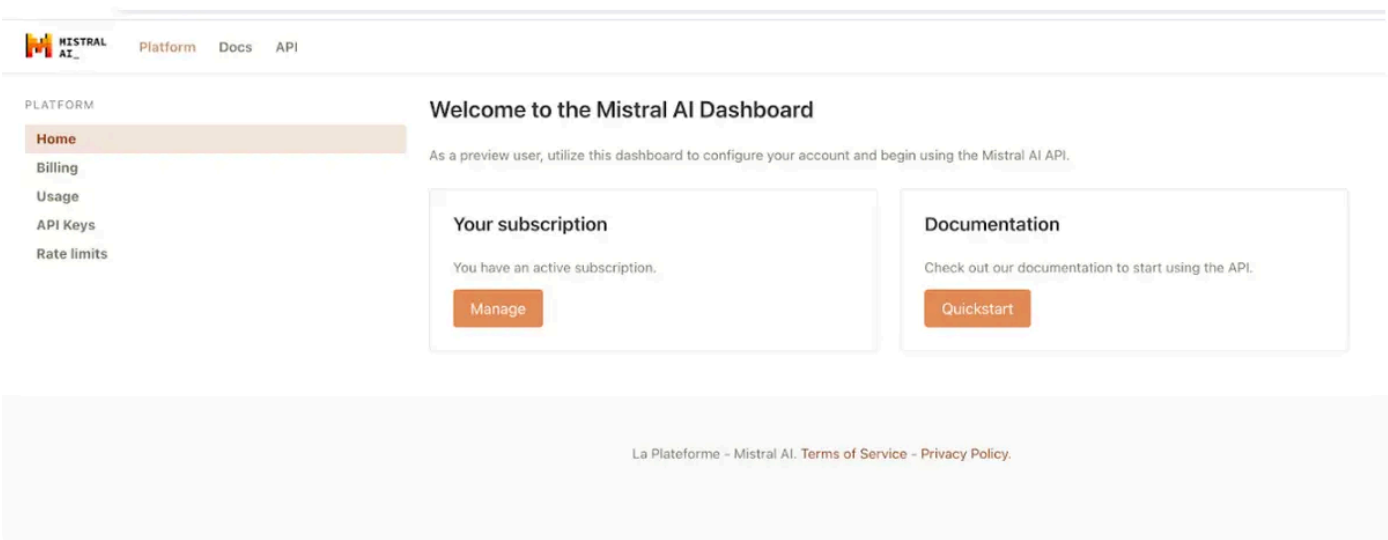
## O. Prerequisites

As explained above, in this article, we are going to leverage the Mistral AI model through Mistral "La plateforme." To get access, you should first [create an account on Mistral AI](#). You may need to wait a few hours (or one day) before your account is activated.



Once your account is activated, you can add your subscription. Follow the instructions step by step on the Mistral AI platform.

Once you have set up your subscription, you can then generate your API key for future usage.



**PLATFORM**

- Home
- Billing
- Usage
- API Keys**
- Rate limits

### Your API keys

API keys are used to query the API. You can create up to 10 API keys, and delete them when you don't need them anymore.

Active Expired + Create new key

Name	Key	Expiration date
No name	*****lnZH	Never
No name	*****JgEB	Never

La Plateforme - Mistral AI. [Terms of Service](#) - [Privacy Policy](#).

Besides using the Mistral AI “La plateforme,” you have another option to implement the Mistral AI model on a machine featuring Nvidia V100, V100S, or A100 GPUs (not an exhaustive list). If you want to deploy a self-hosted large language model on a public or private cloud, you can refer to my previous article on [how to deploy Mistral AI within 10 minutes](#).

## 1. Import needed libraries

This section shows the versions of the required libraries. Personally, I run my code in VScode. So you need to install the following libraries beforehand. Here is the version at the moment I’m running the following code.

```

1  mistralai
2  pymongo                                4.3.3
3  gradio                                 4.10.0
4  gradio_client                          0.7.3
5  langchain                              0.0.348
6  langchain-core                         0.0.12
7  pandas                                 2.0.3

```



These include libraries for data processing, web scraping, AI models, and database interactions.

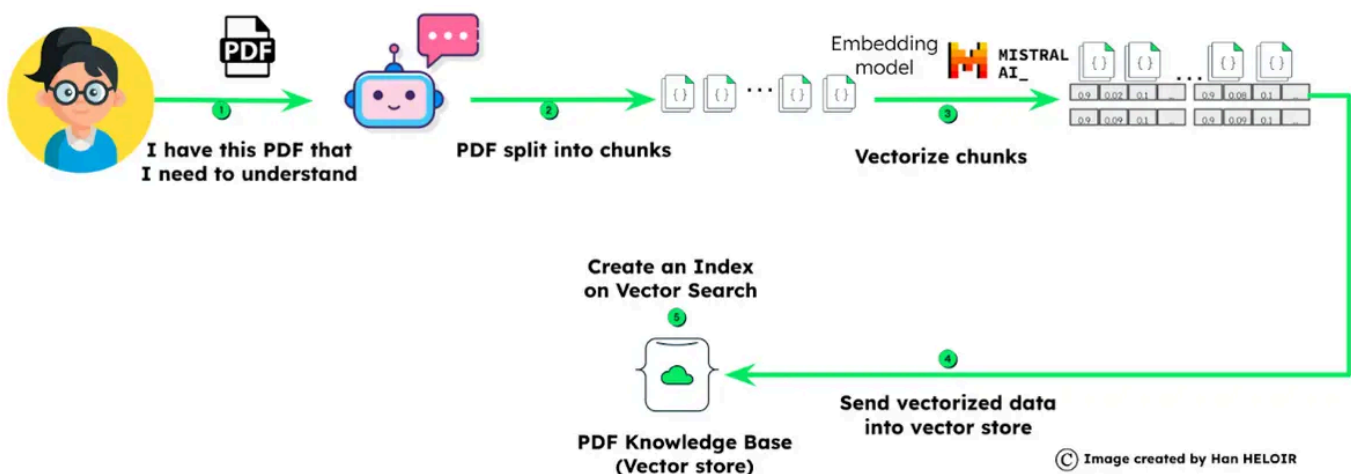
```

1  import gradio as gr
2  import os
3  import pymongo
4  import pandas as pd
5  from mistralai.client import MistralClient
6  from mistralai.models.chat_completion import ChatMessage
7  from langchain.document_loaders import PyPDFLoader
8  from langchain.text_splitter import RecursiveCharacterTextSplitter

```

## 2. Data preparation

The `data_prep()` function loads data from a PDF, a document, or a specified URL. It extracts text content from a webpage/documentation, removes unwanted elements, and then splits the data into manageable chunks.



Once the data is chunked, we use the Mistral AI embedding endpoint to compute embeddings for every chunk and save them in the document. Afterward, each document is added to a MongoDB collection.

```
1 def data_prep(file):
2     # Set up Mistral client
3     api_key = os.environ["MISTRAL_API_KEY"]
4     client = MistralClient(api_key=api_key)
5
6     # Process the uploaded file
7     loader = PyPDFLoader(file.name)
8     pages = loader.load_and_split()
9
10    # Split data
11    text_splitter = RecursiveCharacterTextSplitter(
12        chunk_size=100,
13        chunk_overlap=20,
14        separators=["\n\n", "\n", "(?<=\. )", " ", ""],
15        length_function=len,
16    )
17    docs = text_splitter.split_documents(pages)
18
19    # Calculate embeddings and store into MongoDB
20    text_chunks = [text.page_content for text in docs]
21    df = pd.DataFrame({'text_chunks': text_chunks})
22    df['embedding'] = df.text_chunks.apply(lambda x: get_embeddin
23
24    collection = connect_mongodb()
25    df_dict = df.to_dict(orient='records')
26    collection.insert_many(df_dict)
27
28    return "PDF processed and data stored in MongoDB."
```



# Connecting to MongoDB server

The `connect_mongodb()` function establishes a connection to a MongoDB server. It returns a collection object that can be used to interact with the database. This function will be called in the `data_prep()` function.

In order to get your MongoDB connection string, you can go to your MongoDB Atlas console, click the “Connect” button on your cluster, and choose the Python driver.

The screenshot shows the MongoDB Atlas console interface. The top navigation bar includes the Atlas logo, a dropdown menu for 'SC-Growth...', and links for 'Access Manager', 'Billing', 'All Clusters', 'Get Help', and a user profile. The main sidebar on the left lists various services and security options. The central panel displays the 'Database Deployments' page for a specific cluster.

**Database Deployments**

Find a database deployment...

**Cluster0** Connect View Monitoring Browse Collections ... DEDICATED

**Enhance Your Experience**

For high performance production applications, upgrade to a M30 cluster now!

Dismiss Upgrade

**Connections**

100.0% 100.0 100.0 B/s

**Disk Usage**

43.5 GB / 128.0 GB (34%)  
Last 30 days  
128.0 GB

VERSION	REGION	CLUSTER TIER	TYPE	BACKUPS	LINKED APP SERVICES	ATLAS SQL	ONLINE ARCHIVE	ATLAS SEARCH
7.0.5	GCP / Paris (europe-west9)	M20 (General)	Replica Set - 3 nodes	Active	<a href="#">Multiple applications linked</a>	<a href="#">Connect</a>	Active	<a href="#">24 search indexes</a>

+ Add Tag

## Connect to Cluster0



### Connecting with MongoDB Driver

#### 1. Select your driver and version

We recommend installing and using the latest driver version.

Driver	Version
Python	3.12 or later

#### 2. Install your driver

Run the following on the command line

Note: Use appropriate Python 3 executable

```
python -m pip install "pymongo[srv]"
```

[View MongoDB Python Driver installation instructions.](#)

#### 3. Add your connection string into your application code

Connect to Cluster and Online Archive

Connect to Cluster

Connect to Online Archive

☐ View full code sample

```
mongodb+srv://<username>:<password>@cluster0.bofm7.mongodb.net/?retryWrites=true&w=majority
```

Replace **<password>** with the password for the **<username>** user. Ensure any option params are [URL encoded](#).

#### RESOURCES

[Get started with the Python Driver](#)

[Python Starter Sample App](#)

```
1 def connect_mongodb():
2     # Your MongoDB connection string
3     mongo_url = os.environ["MONGO_URI"]
4     client = pymongo.MongoClient(mongo_url)
5     db = client["mistralpdf"]
6     collection = db["pdfRAG"]
7     return collection
```

You can import your mongo\_url by doing the following command in shell.

```
1 export MONGO_URI="Your_cluster_connection_string"
```



## Getting the embedding

The `get_embedding(text)` function generates an embedding for a given text. It replaces newline characters and then uses Mistral AI “La plateforme” embedding endpoints to get the embedding. This function will be called in both data preparation and question and answering processes.

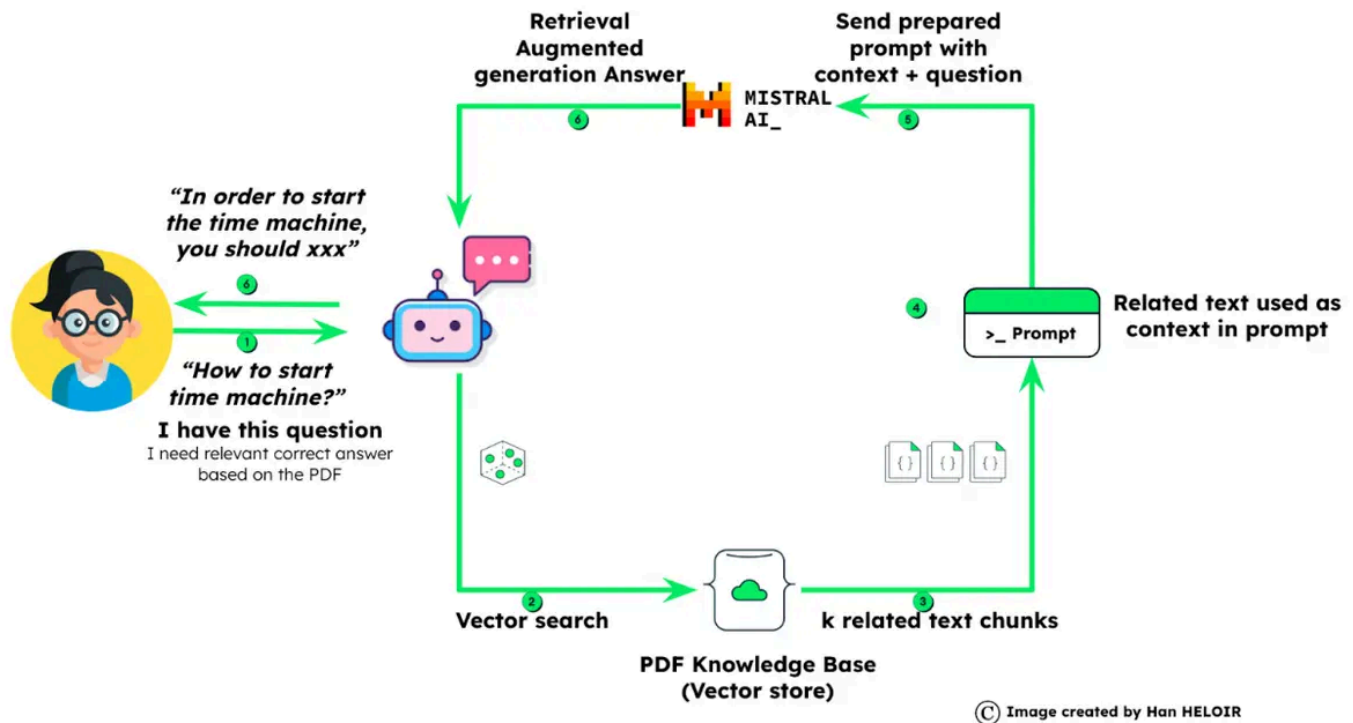
```
1 def get_embedding(text, client):
2     text = text.replace("\n", " ")
3     embeddings_batch_response = client.embeddings(
4         model="mistral-embed",
5         input=text,
6     )
7     return embeddings_batch_response.data[0].embedding
```



## 3. Question and answer function

This function is the core of the program. It processes a user's question and creates a response using the context supplied by Mistral AI.





This process involves several key steps. Here's how it works:

- Firstly, we generate a numerical representation, called an embedding, through a Mistral AI embedding endpoint, for the user's question.
- Next, we run a vector search in the MongoDB collection to identify the documents similar to the user's question.
- It then constructs a contextual background by combining chunks of text from these similar documents. We prepare an assistant instruction by combining all this information.
- The user's question and the assistant's instruction are prepared into a prompt for the Mistral AI model.
- Finally, Mistral AI will generate responses to the user thanks to the retrieval-augmented generation process.

```
1 def qna(users_question):
2     # Set up Mistral client
```

```
3     api_key = os.environ["MISTRAL_API_KEY"]
4     client = MistralClient(api_key=api_key)
5
6     question_embedding = get_embedding(users_question, client)
7     print("-----Here is user question-----")
8     print(users_question)
9     documents = find_similar_documents(question_embedding)
10
11     print("-----Retrieved documents-----")
12     print(documents)
13     for doc in documents:
14         doc['text_chunks'] = doc['text_chunks'].replace('\n', ' ')
15
16     for document in documents:
17         print(str(document) + "\n")
18
19     context = " ".join([doc["text_chunks"] for doc in documents])
20     template = f"""
21     You are an expert who loves to help people! Given the followi
22     question using only the given context. If you are unsure and
23     explicitly written in the documentation, say "Sorry, I don't
24
25     Context sections:
26     {context}
27
28     Question:
29     {users_question}
30
31     Answer:
32     """
```

```
33     messages = [ChatMessage(role="user", content=template)]
34     chat_response = client.chat(
35         model="mistral-tiny",
36         messages=messages,
37     )
38     formatted_documents = '\n'.join([doc['text_chunks'] for doc in documents])
39
40     return chat_response.choices[0].message, formatted_documents
```



## The last configuration on the MongoDB vector search index

In order to run a vector search query, you only need to create a vector search index in MongoDB Atlas as follows. (You can also learn more about [how to create a vector search index.](#))

```
1  {
2    "type": "vectorSearch",
3    "fields": [
4      {
5        "numDimensions": 1536,
6        "path": "'embedding'",
7        "similarity": "cosine",
8        "type": "vector"
9      }
10   ]
11 }
```



## Finding similar documents

The `find_similar_documents(embedding)` function runs the vector search query in a MongoDB collection. This function will be called when the user asks a question. We will use this function to find similar documents to the questions in the question and answering process.

```
1 def find_similar_documents(embedding):
2     collection = connect_mongodb()
3     documents = list(
4         collection.aggregate([
5             {
6                 "$vectorSearch": {
7                     "index": "vector_index",
8                     "path": "embedding",
9                     "queryVector": embedding,
10                    "numCandidates": 20,
11                    "limit": 10
12                }
13            },
14            {"$project": {"_id": 0, "text_chunks": 1}}
15        ]))
16     return documents
```



## 4. Gradio user interface

In order to have a better user experience, we wrap the PDF upload and chatbot into two tabs using Gradio. Gradio is a Python library that enables the fast creation of customizable web applications for machine learning models and data processing workflows. You can put this code at the end of your Python file. Inside of this function, depending on which tab you are using, either data preparation or

question and answering, we will call the explained dataprep() function or qna() function.

```
1  # Gradio Interface for PDF Upload
2  pdf_upload_interface = gr.Interface(
3      fn=data_prep,
4      inputs=gr.File(label="Upload PDF"),
5      outputs="text",
6      allow_flagging="never"
7  )
8
9  # Gradio Interface for Chatbot
10 chatbot_interface = gr.Interface(
11     fn=qna,
12     inputs=gr.Textbox(label="Enter Your Question"),
13     outputs=[
14         gr.Textbox(label="Mistral Answer"),
15         gr.Textbox(label="Retrieved Documents from MongoDB Atlas")
16     ],
17     allow_flagging="never"
18 )
19
20 # Combine interfaces into tabs
21 iface = gr.TabbedInterface(
22     [pdf_upload_interface, chatbot_interface],
23     ["Upload PDF", "Chatbot"]
24 )
25
26 # Launch the Gradio app
```

77

Upload PDF

Chatbot

Upload PDF

↑

Déposer le Fichier Ici

- ou -

Cliquer pour Télécharger

Clear

Submit

output

Upload PDF

Chatbot

Enter Your Question

how Mongoddb support edge computing which feature

Clear

Submit

Mistral Answer

role='assistant' content='MongoDB supports edge computing by providing data management capabilities and handling data synchronization in edge environments. It can be used to combine edge computing solutions with its flexible data model and robust syncing mechanisms, benefiting industries such as energy, telecommunication, and others by enabling real-time data processing and analysis at the edge. (Refer to the context sections titled "MongoDB supports edge computing," "MongoDB\'s role in edge computing," "How does MongoDB support in edge environments," and "Benefits of using MongoDB for edge computing.")'

Retrieved Document from MongoDB Atlas

MongoDB support edge computing architectures? How does MongoDB handle data synchronization in edge use cases for edge computing? What is MongoDB's role in edge computing? How does MongoDB support in edge environments? What are the benefits of using MongoDB for edge computing? Can MongoDB be used solutions. How does MongoDB support edge computing architectures? MongoDB supports edge computing What are the benefits of using MongoDB for edge computing? MongoDB provides benefits such as flexible is MongoDB's role in edge computing? MongoDB's role in edge computing involves providing the data

## Conclusion

This detailed guide has delved into the dynamic combination of Mistral AI and MongoDB, showcasing how to develop a bespoke large language model GenAI application. Integrating the advanced capabilities of Mistral AI with MongoDB's

robust data management features enables the creation of a custom AI assistant that caters to unique requirements.

We have provided a straightforward, step-by-step methodology, covering everything from initial data gathering and segmentation to the generation of embeddings and efficient data querying. This guide serves as a comprehensive blueprint for implementing the system, complemented by practical code examples and instructions for setting up Mistral AI on a GPU-powered machine and linking it with MongoDB.

Leveraging Mistral AI and MongoDB Atlas, users gain access to the expansive possibilities of AI applications, transforming our interaction with technology and unlocking new, secure ways to harness data insights while maintaining privacy.

## Learn more

To learn more about how Atlas helps organizations integrate and operationalize GenAI and LLM data, take a look at our [Embedding Generative AI whitepaper](#) to explore RAG in more detail.

If you want to know more about how to deploy a self-hosted Mistral AI with MongoDB, you can refer to my previous articles:

[Unleashing AI Sovereignty: Getting Mistral.ai 7B Model Up and Running in Less Than 10 Minutes](#)

and

[Starting Today with Mistral AI & MongoDB: A Beginner's Guide to a Self-Hosted LLM Generative AI Application](#)

. [Mixture of Experts Explained](#)

---

Rate this tutorial

## Related

CODE EXAMPLE

### A Spotify Song and Playlist Recommendation Engine

---

Nov 13, 2023 | 6 min read

QUICKSTART

### Getting Started with MongoDB and Starlette

---

Sep 23, 2022 | 5 min read

TUTORIAL

### Symfony and MongoDB Workshop: Building a Rental Listing Application

---

Apr 08, 2024 | 3 min read

ARTICLE

### Real-Time Card Fraud Solution Accelerator with MongoDB and Databricks

---

Jul 11, 2023 | 7 min read

[Request a Tutorial](#)



## About

[Careers](#)[Investor Relations](#)[Legal Notices](#)[Privacy Notices](#)[Security Information](#)[Trust Center](#)

## Support

[Contact Us](#)[Customer Portal](#)[Atlas Status](#)[Customer Support](#)[Manage Cookies](#)

## Social

[GitHub](#)[Stack Overflow](#)[LinkedIn](#)[YouTube](#)[X](#)[Twitch](#)



© 2024 MongoDB, Inc.