

Cheat Sheet: Project: Generative AI Applications with RAG and LangChain

1. 1

1. </tr>

Copied!

Package/Method	Description	Code example
Load method	Loads data from a server and puts the returned data into the selected element.	<pre>1. 1 1. data = loader.load()</pre> <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">Copied!</div>
Document object	Contains information about data in LangChain. It has two attributes: <ul style="list-style-type: none"> • page_content: str: This attribute holds the content of the document. • metadata: dict: This attribute contains arbitrary metadata associated with the document. It can be used to track various details such as the document id, file name, and so on. 	<pre>1. 1 2. 2 3. 3 1. from langchain_core.documents import Document 2. Document(page_content="Python is an interpreted high-level general-purpose programming language. Python's design philosophy emphasizes code readability with its not metadata={ 'my_document_id' : 234234, 'my_document_source' : "About Python", 'my_document_create_time' : 1680013019 } 3.</pre> <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">Copied!</div>
pprint function	A function in Python used to “pretty-print” data structures, making them more readable and easier to understand.	<pre>1. 1 1. pprint(data[0].page_content[:1000])</pre> <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">Copied!</div>
PyPDFLoader	Simplifies the process of loading PDF documents into a format that can be easily manipulated and analyzed within your applications.	<pre>1. 1 2. 2 3. 3 4. 4 1. pdf_url = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/Q81D33CdRLK6L 2. 3. loader = PyPDFLoader(pdf_url) 4. pages = loader.load_and_split()</pre> <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">Copied!</div>
PyMuPDFLoader	The fastest of the PDF parsing options. It provides detailed metadata about the PDF and its pages and returns one document per page.	<pre>1. 1 2. 2 3. 3 4. 4 5. 5 1. loader = PyMuPDFLoader(pdf_url) 2. loader 3. 4. data = loader.load() 5. print(data[0])</pre> <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">Copied!</div>
UnstructuredMarkdownLoader	A powerful tool within the LangChain framework that facilitates the loading of Markdown documents into a structured format suitable for downstream processing.	<pre>1. 1 2. 2 3. 3 4. 4 5. 5 1. !wget 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/eMSP5vJjj9y0fAaC1: 2. 3. markdown_path = "markdown-sample.md" loader = UnstructuredMarkdownLoader(markdown_path) loader 4. data = loader.load() 5. data</pre> <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">Copied!</div>
JSONLoader	A module that builds a straightforward Python object from loaded JSON or similar dict-based data loading. It also checks if the input-loaded JSON has all the necessary attributes for the pipeline and that it has the right types.	<pre>1. 1 1. !wget 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/hAmzVJeOUAMHzmhUH1</pre> <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">Copied!</div>

Package/Method	Description	Code example
CSVLoader	CSV files are a common format for storing tabular data. The CSVLoader provides a convenient way to read and process this data.	<pre data-bbox="654 105 1578 228"> 1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 1. loader = UnstructuredCSVLoader(2. file_path="mlb-teams-2012.csv", mode="elements" 3.) 4. data = loader.load() 5. 6. data[0].page_content 7. print(data[0].metadata["text_as_html"]) </pre>
UnstructuredCSVLoader	The UnstructuredCSVLoader considers the entire CSV file as a single unstructured table element. This approach is beneficial when you want to analyze the data as a complete table rather than as separate entries.	<pre data-bbox="654 228 1578 508"> 1. 1 2. 2 3. 3 4. 4 5. 5 1. import requests 2. from bs4 import BeautifulSoup 3. 4. url = 'https://www.ibm.com/topics/langchain' response = requests.get(url) 5. soup = BeautifulSoup(response.content, 'html.parser') print(soup.prettify()) </pre>
BeautifulSoup	A Python library used for web scraping purposes to pull the data out of HTML and XML files. It creates a parse tree for parsed pages that can be used to extract data easily.	<pre data-bbox="654 508 1578 846"> 1. 1 2. 2 3. 3 4. 4 5. 5 1. For single page: 2. loader = WebBaseLoader("https://www.ibm.com/topics/langchain") 3. 4. data = loader.load() 5. data 6. For multiple pages: loader = WebBaseLoader(["https://www.ibm.com/topics/langchain", "https://www.redhat.com/en/"]) data </pre>
WebBaseLoader	LangChain's tool designed to extract all text from HTML webpages and convert it into a document format suitable for further processing.	<pre data-bbox="654 846 1578 1220"> 1. 1 2. 2 3. 3 4. 4 5. 5 1. !wget https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/94hiHUNLZdb0bLMkrC 2. 3. loader = Docx2txtLoader("file-sample.docx") 4. data = loader.load() 5. data </pre>
Docx2txtLoader	Utilized to convert Word documents into a document format suitable for further processing.	<pre data-bbox="654 1220 1578 1564"> 1. 1 2. 2 3. 3 4. 4 5. 5 1. loader = UnstructuredFileLoader("companypolicies.txt") 2. data = loader.load() 3. data </pre>
Load .txt file	Supports the loading of .txt files when you need to load content from various text sources and formats without writing a separate loader for each one.	<pre data-bbox="654 1564 1578 1854"> 1. 1 2. 2 3. 3 1. loader = UnstructuredFileLoader("markdown-sample.md") 2. data = loader.load() 3. data </pre>
Load .md file	Supports the loading of .md files when you need to load content from various text sources and formats without writing a separate loader for each one.	

Package/Method	Description	Code example
		Copied!
Load multiple files with different formats	Supports the loading of multiple file types when you need to load content from various text sources and formats without writing a separate loader for each one.	<pre> 1. 1 2. 2 3. 3 4. 4 5. 5 6. 7. 8. 9. 10. 11. 12. 13. 14. 15. 16. 17. 18. 19. 20. 21. 22. </pre>
Model ID	In LangChain, the model ID is used to specify which language model you want to use. This ID can vary depending on the model provider and the specific model you are accessing.	<pre> 1. files = ["markdown-sample.md", "companypolicies.txt"] 2. 3. loader = UnstructuredFileLoader(files) 4. data = loader.load() 5. data </pre> <p data-bbox="633 432 714 464">Copied!</p> <pre> 1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8 9. 9 10. 10 11. 11 12. 12 13. 13 14. 14 15. 15 16. 16 17. 17 18. 18 19. 19 20. 20 21. 21 22. 22 </pre> <pre> 1. def llm_model(model_id): 2. parameters = { 3. GenParams.MAX_NEW_TOKENS: 256, # this controls the maximum number of tokens in the 4. GenParams.TEMPERATURE: 0.5, # this randomness or creativity of the model's response 5. } 6. 7. credentials = { 8. "url": "https://us-south.ml.cloud.ibm.com" 9. } 10. 11. project_id = "skills-network" 12. 13. model = ModelInference(14. model_id=model_id, 15. params=parameters, 16. credentials=credentials, 17. project_id=project_id 18.) 19. llm = WatsonxLLM(watsonx_model = model) 20. return llm 21. 22. </pre>
Load source document	Loading a source document into a large language model (LLM) involves providing the model with specific data or text that it can be used to generate responses or perform tasks.	<p data-bbox="633 1087 714 1119">Copied!</p> <pre> 1. 1 2. !wget "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/d_ahNwb1L2duIxBR6 </pre> <p data-bbox="633 1763 714 1795">Copied!</p>
LangChain prompt template	A prompt template is set up using LangChain to make it reusable.	<pre> 1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8 </pre>

Package/Method	Description	Code example
Use mixtral model	A sparse mixture-of-experts (SMoE) network developed by Mistral AI. It is a decoder-only transformer model with a unique architecture that includes 8 experts per feedforward block, totaling 45 billion parameters.	<pre> 9. 9 10. 10 11. 11 1. template = """According to the document content here 2. {content}, 3. answer this question 4. {question}. 5. Do not try to make up the answer. 6. 7. YOUR RESPONSE: 8. 9. """ 10. prompt_template = PromptTemplate(template=template, input_variables=['content', 'question'] 11. prompt_template) Copied!</pre> <p>1. 1 2. 2 3. 3 4. 4</p> <p>1. mixtral_llm = llm_model('mistralai/mixtral-8x7b-instruct-v01') 2. 3. query_chain = LLMChain(llm=mixtral_llm, prompt=prompt_template) 4. query = "It is in which year of our nation?" response = query_chain.invoke(input={'content': content, 'question': query}) print(response['text'])</p>
Use Llama 3 model	The Llama model (Large	<p>Copied!</p> <p>1. 1 2. 2</p> <p>1. query_chain = LLMChain(llm=llama_llm, prompt=prompt_template) 2. query_chain</p>
Use one piece of information	Language Model Meta AI) is a family of autoregressive large language models developed by Meta AI.	<p>Copied!</p> <p>1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8</p> <p>1. content = """ 2. The only nation that can be defined by a single word: possibilities. 3. 4. So on this night, in our 245th year as a nation, I have come to report on the State of the ! 5. 6. And my report is this: the State of the Union is strong—because you, the American people, a 7. 8. """ Copied!</p> <p>1. 1 2. 2 3. 3</p>
Split by Character	This is the simplest method of splitting text, which splits the text based on characters (by default "\n\n") and measures chunk length by the number of characters.	<p>1. from langchain.text_splitter import CharacterTextSplitter 2. 3. text_splitter = CharacterTextSplitter(separator="", chunk_size=200, chunk_overlap=20, length_function=len,)</p>
Recursively Split by Character	A text splitter recommended for generic text. It is parameterized by a list of characters, and it tries to split them in order until the chunks are small enough.	<p>Copied!</p> <p>1. 1 2. 2 3. 3</p> <p>1. from langchain.text_splitter import RecursiveCharacterTextSplitter 2. 3. text_splitter = RecursiveCharacterTextSplitter(chunk_size=100, chunk_overlap=20, length_function=len,)</p>

Package/Method	Description	Code example
Split Code	This method allows you to split your code, supporting multiple programming languages. It is based on the Recursively Split by Character strategy.	Copied!
		<pre> 1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8 1. PYTHON_CODE = """ 2. def hello_world(): 3. print("Hello, World!") 4. 5. # Call the function 6. hello_world() 7. 8. """ python_splitter = RecursiveCharacterTextSplitter.from_language(language=Language.PYTHON, chunk_size=50, chunk_overlap=0) python_docs = python_splitter.create_documents([PYTHON_CODE]) python_docs </pre>
Markdown Header Text Splitter	A Markdown file is organized by headers. Creating chunks within specific header groups is an intuitive approach. This splitter will divide a	Copied!
	Markdown file based on a specified set of headers.	<pre> 1. 1 2. 2 3. 3 1. markdown_splitter = MarkdownHeaderTextSplitter(headers_to_split_on=headers_to_split_on) 2. md_header_splits = markdown_splitter.split_text(md) 3. md_header_splits </pre>
Split by HTML	This splitting method is a "structure-aware" chunker that splits text at the element level and adds metadata for each header "relevant" to any given chunk.	Copied!
		<pre> 1. 1 2. 2 3. 3 1. html_splitter = HTMLHeaderTextSplitter(headers_to_split_on=headers_to_split_on) 2. html_header_splits = html_splitter.split_text(html_string) 3. html_header_splits </pre>
embed_query using watsonx	A method used to embed a single piece of text (e.g., for the purpose of comparing it to other embedded pieces of text).	Copied!
		<pre> 1. query = "How are you?" 2. 3. query_result = watsonx_embedding.embed_query(query) </pre>
embed_documents using watsonx	A method commonly used in various contexts for embedding documents within other documents, or in machine learning for embedding text data.	Copied!
		<pre> 1. 1 2. doc_result = watsonx_embedding.embed_documents(chunks) </pre>
TextLoader	LangChain's TextLoader is a useful tool for loading and processing text data, making it ready for use with large language models (LLMs).	Copied!
		<pre> 1. 1 1. !wget "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/BYlUHaillwM8EUIta </pre>
Embedding model	Embedding models are specifically designed to interface with text embeddings.	Copied!
	Embeddings generate a vector representation for a given piece of text. This is advantageous as it allows you to conceptualize text within a vector space. Consequently, you can perform operations such as semantic search, where you identify pieces of text that are most similar within the vector space.	<pre> 1. from ibm_watsonx_ai.metanames import EmbedTextParamsMetaNames 2. from langchain_ibm import WatsonxEmbeddings 3. 4. embed_params = { EmbedTextParamsMetaNames.TRUNCATE_INPUT_TOKENS: 3, EmbedTextParamsMetaNames.RETURN_OPTIONS: {"input_text": True}, } 5. watsonx_embedding = WatsonxEmbeddings(model_id="ibm/slate-125m-english-rtrvr", url="https://us-south.ml.cloud.ibm.com", project_id="skills-network", params=embed_params,) </pre>
		Copied!

Package/Method	Description	Code example
Using Chroma DB to store embeddings	<p>Refers to using the embedding model to create embeddings for each chunk and then storing them in the Chroma database.</p>	<pre>1. 1 2. vectordb = Chroma.from_documents(chunks, watsonx_embedding, ids=ids)</pre> <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">Copied!</div>
	<p>A vector database that involves finding items that are most similar to a given query item based on their vector representations.</p>	
Similarity search	<p>In this process, data objects are converted into vectors (which you've already done), and the search algorithm identifies and retrieves those with the closest vector distances to the query, enabling efficient and accurate identification of similar items in large datasets.</p>	<pre>1. 1 2. 2 3. 3 4. query = "Email policy" 5. docs = vectordb.similarity_search(query) 6. docs</pre> <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">Copied!</div>
Using FAISS DB to store embeddings	<p>Here is an example of how to perform a similarity search based on the query "Email policy."</p>	
	<p>FAISS is another vector database that is supported by LangChain.</p>	
Defining helper functions	<p>The process of building and using FAISS is similar to Chroma DB.</p>	<pre>1. 1 2. faissdb = FAISS.from_documents(chunks, watsonx_embedding, ids=ids)</pre> <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">Copied!</div>
	<p>However, there may be differences in the retrieval results between FAISS and Chroma DB.</p>	
mixtral-8x7b-instruct-v01	<p>Helper functions are smaller, reusable functions that perform specific tasks and can be called within other functions to simplify code and avoid repetition. They help make code more modular, readable, and maintainable.</p>	<pre>1. 1 2. 2 3. 3 4. 4 5. 5 6. def warn(*args, **kwargs): 7. pass 8. import warnings 9. warnings.warn = warn 10. warnings.filterwarnings('ignore')</pre> <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">Copied!</div>
	<p>An LLM model developed by Mistral AI. It's a Sparse Mixture of Experts (SMoE) model, which means it uses a combination of different expert models to generate high-quality text outputs.</p>	

Package/Method	Description	Code example
	<pre> 9. credentials = { 10. "url": "https://us-south.ml.cloud.ibm.com" 11. } 12. project_id = "skills-network" 13. 14. model = ModelInference(15. model_id=model_id, 16. params=parameters, 17. credentials=credentials, 18. project_id=project_id 19.) 20. 21. mixtral_llm = WatsonxLLM(model = model) 22. return mixtral_llm 23. 24. </pre>	Copied!
MMR retrieval	<p>MMR in vector stores is a technique used to balance the relevance and diversity of retrieved results. It selects documents that are both highly relevant to the query and minimally similar to previously selected documents.</p>	<pre> 1. 1 2. 2 3. 3 </pre> <pre> 1. retriever = vectordb.as_retriever(search_type="mmr") 2. docs = retriever.invoke(query) 3. docs </pre>
Similarity score threshold retrieval	<p>You can set a retrieval method that defines a similarity score threshold, returning only documents with a score above that threshold.</p>	<pre> 1. 1 2. 2 3. 3 4. 4 5. 5 </pre> <pre> 1. retriever = vectordb.as_retriever(2. search_type="similarity_score_threshold", search_kwargs={"score_threshold": 0.4} 3.) 4. docs = retriever.invoke(query) 5. docs </pre>
Self-Querying Retriever	<p>A Self-Querying Retriever has the ability to query itself. Specifically, given a natural language query, the retriever uses a query-constructing LLM chain to generate a structured query. It then applies this structured query to its underlying vector store. This enables the retriever to not only use the user-input query for semantic similarity comparison with the contents of stored documents but also to extract and apply filters based on the metadata of those documents.</p>	<pre> 1. 1 2. 2 3. 3 4. 4 </pre> <pre> 1. from langchain_core.documents import Document 2. from langchain.chains.query_constructor.base import AttributeInfo 3. from langchain.retrievers.self_query.base import SelfQueryRetriever 4. from lark import lark </pre>
Parent Document Retriever	<p>When splitting documents for retrieval, there are often conflicting desires:</p> <ul style="list-style-type: none"> • You may want to have small documents so that their embeddings can most accurately reflect their meaning. If the documents are too long, the embeddings can lose meaning. • You want to have long enough documents so that the context of each chunk is retained. 	<pre> 1. 1 2. 2 3. 3 </pre> <pre> 1. from langchain.retrievers import ParentDocumentRetriever 2. from langchain_text_splitters import CharacterTextSplitter 3. from langchain.storage import InMemoryStore </pre>

Package/Method	Description	Code example
Multi-Query Retriever	<p>The Parent Document Retriever strikes that balance by splitting and storing small chunks of data.</p> <p>The Multi Query Retriever uses an LLM to generate multiple queries from different perspectives for a given user input query. For each query, it retrieves a set of relevant documents and then takes the unique union of these results to form a larger set of potentially relevant documents.</p>	<pre> 1. 1 2. 2 3. 3 4. 4 5. 6. 7. 8. 9. 10. 11. 12. 13. 1. import gradio as gr 2. 3. def add_numbers(Num1, Num2): return Num1 + Num2 4. # Define the interface 5. demo = gr.Interface(6. fn=add_numbers, 7. inputs=[gr.Number(), gr.Number()], # Create two numerical input fields where users can 8. outputs=gr.Number() # Create numerical output fields 9.) 10. 11. 12. # Launch the interface 13. demo.launch(server_name="127.0.0.1", server_port= 7860) </pre>
sum calculator	<p>An application that can calculate the sum of your input numbers in Gradio.</p>	<pre> Copied! 1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 8. 9. 10. 11. 12. 13. 14. 1. # Import necessary packages 2. from ibm_watsonx_ai.foundation_models import ModelInference 3. from ibm_watsonx_ai.metanames import GenTextParamsMetaNames as GenParams 4. from ibm_watsonx_ai import Credentials 5. from langchain_ibm import WatsonxLLM 6. import gradio as gr </pre>
Integrate application into Gradio	<p>You can integrate an application with Gradio to leverage a web interface for inputting questions and receiving responses.</p>	<pre> Copied! 1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 8. 9. 10. 11. 12. 13. 14. 1. # Model and project settings 2. model_id = 'mistralai/mistral-8x7b-instruct-v01' # Directly specifying the model 3. 1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8 9. 9 10. 10 11. 11 12. 12 13. 13 14. 14 </pre>

Package/Method**Description****Code example**

```

15. 15
16. 16
17. 17
18. 18
19. 19
20. 20
21. 21
22. 22
23. 23
24. 24
25. 25
26. 26
27. 27
28. 28
29. 29
30. 30
31. 31
32. 32
33. 33
34. 34
35. 35

1. # Set necessary parameters
2. parameters = {
3.     GenParams.MAX_NEW_TOKENS: 256, # Specifying the max tokens you want to generate
4.     GenParams.TEMPERATURE: 0.5, # This randomness or creativity of the model's responses
5. }
6.

7. project_id = "skills-network"
8. # Wrap up the model into WatsonxLLM inference
9. watsonx_llm = WatsonxLLM(
10.     model_id=model_id,
11.     url="https://us-south.ml.cloud.ibm.com",
12.     project_id=project_id,
13.     params=parameters,
14. )
15.

16. # Function to generate a response from the model
17. def generate_response(prompt_txt):
18.     generated_response = watsonx_llm.invoke(prompt_txt)
19.     return generated_response
20.

21. # Create Gradio interface
22. chat_application = gr.Interface(
23.     fn=generate_response,
24.     allow_flagging="never",
25.     inputs=gr.Textbox(label="Input", lines=2, placeholder="Type your question here..."),
26.     outputs=gr.Textbox(label="Output"),
27.     title="Watsonx.ai Chatbot",
28.     description="Ask any question and the chatbot will try to answer."
29. )
30.

31. # Launch the app
32. chat_application.launch(server_name="127.0.0.1", server_port= 7860)
33.

34. </td>
35.

```

Copied!**Initialize the LLM**

You can initialize the LLM by creating an instance of WatsonxLLM, a class in langchain_ibm. WatsonxLLM can use several underlying foundational models. In this snippet, you use Mixtral 8x7B.

To initialize the LLM, paste the following code into

1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9
10. 10
11. 11
12. 12

Package/Method**Description****Code example**

qabot.py. Note that you are initializing the model with a temperature of 0.5, and allowing for the generation of a maximum of 256 tokens.

```

13. 13
14. 14
15. 15

1. ## LLM
2. def get_llm():
3.     model_id = 'mistralai/mixtral-8x7b-instruct-v01'
4.     parameters = {
5.         GenParams.MAX_NEW_TOKENS: 256,
6.         GenParams.TEMPERATURE: 0.5,
7.     }
8.     project_id = "skills-network"
9.     watsonx_llm = WatsonxLLM(
10.        model_id=model_id,
11.        url="https://us-south.ml.cloud.ibm.com",
12.        project_id=project_id,
13.        params=parameters,
14.    )
15.    return watsonx_llm

```

Copied!

You use the PyPDFLoader class from the langchain_community library to load PDF documents.

Define the PDF document loader

You create the PDF loader as an instance of PyPDFLoader. Then, you load the document and return the loaded document. To incorporate the PDF loader in your bot, add the following code to qabot.py.

```

1. 1
2. 2
3. 3
4. 4
5. 5

```

```

1. ## Document loader
2. def document_loader(file):
3.     loader = PyPDFLoader(file.name)
4.     loaded_document = loader.load()
5.     return loaded_document

```

Copied!

You define a document splitter that will split the text into chunks. Add the following code to qabot.py to define such a text splitter. Note that, in this example, you are defining a RecursiveCharacterTextSplitter with a chunk size of 1000, although other splitters or parameter values are possible.

Define the text splitter

```

1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9

```

```

1. ## Text splitter
2. def text_splitter(data):
3.     text_splitter = RecursiveCharacterTextSplitter(
4.         chunk_size=1000,
5.         chunk_overlap=50,
6.         length_function=len,
7.     )
8.     chunks = text_splitter.split_documents(data)
9.     return chunks

```

Copied!

Add this code to qabot.py to define a function that embeds the chunks using a yet-to-be-defined embedding model and stores the embeddings in a ChromaDB vector store.

Define the vector store

```

1. 1
2. 2
3. 3
4. 4
5. 5

```

```

1. ## Vector db
2. def vector_database(chunks):
3.     embedding_model = watsonx_embedding()
4.     vectordb = Chroma.from_documents(chunks, embedding_model)
5.     return vectordb

```

Copied!

Define the embedding model

Defines a watsonx_embedding() function that returns an instance of WatsonxEmbeddings, a class from langchain_ibm that generates embeddings. In this case, the embeddings are generated using IBM's Slate 125M English embeddings model. Paste this code into the qabot.py file.

```

1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9
10. 10
11. 11
12. 12
13. 13

```

```

1. ## Embedding model
2. def watsonx_embedding():
3.     embed_params = {
4.         EmbedTextParamsMetaNames.TRUNCATE_INPUT_TOKENS: 3,
5.         EmbedTextParamsMetaNames.RETURN_OPTIONS: {"input_text": True},
6.     }
7.     watsonx_embedding = WatsonxEmbeddings(
8.         model_id="ibm/slate-125m-english-rtrv",
9.         url="https://us-south.ml.cloud.ibm.com",
10.        project_id="skills-network",
11.        params=embed_params,
12.    )

```

Package/Method	Description	Code example
Define a question-answering chain	<p>Use RetrievalQA from LangChain, a chain that performs natural-language question-answering over a data source using retrieval-augmented generation (RAG). Add the following code to qabot.py to define a question-answering chain.</p>	<pre data-bbox="649 107 975 137">13. return watsonx_embedding</pre> <div data-bbox="649 149 714 175" style="border: 1px solid black; padding: 2px;">Copied!</div> <pre data-bbox="657 185 1323 593"> 1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8 9. 9 10. 10 1. ## QA Chain 2. def retriever_qa(file, query): 3. llm = get_llm() 4. retriever_obj = retriever(file) 5. qa = RetrievalQA.from_chain_type(llm=llm, 6. chain_type="stuff", 7. retriever=retriever_obj, 8. return_source_documents=False) 9. response = qa.invoke(query) 10. return response['result'] </pre>
Setup the Gradio interface	<p>A Gradio interface should include:</p> <ul style="list-style-type: none"> • A file upload functionality (provided by the File class in Gradio) • An input textbox where the question can be asked (provided by the Textbox class in Gradio) • An output textbox where the question can be answered (provided by the Textbox class in Gradio) 	<pre data-bbox="649 656 1323 686">1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8 9. 9 10. 10 11. 11 12. 12</pre> <div data-bbox="649 901 997 927" style="border: 1px solid black; padding: 2px;">Copied!</div> <pre data-bbox="649 916 1569 1142"> 1. # Create Gradio interface 2. rag_application = gr.Interface(3. fn=retriever_qa, 4. allow_flagging="never", 5. inputs=[6. gr.File(label="Upload PDF File", file_count="single", file_types=['.pdf'], type="file"), 7. gr.Textbox(label="Input Query", lines=2, placeholder="Type your question here...") 8.], 9. outputs=gr.Textbox(label="Output"), 10. title="RAG Chatbot", 11. description="Upload a PDF document and ask any question. The chatbot will try to answer", 12.) </pre>
Add code to launch the application	<p>Add this line to qabot.py to launch the application using port 7860.</p>	<pre data-bbox="649 1248 1307 1277">1. 1 2. 2</pre> <div data-bbox="649 1281 714 1307" style="border: 1px solid black; padding: 2px;">Copied!</div> <pre data-bbox="649 1271 1307 1300"> 1. # Launch the app 2. rag_application.launch(server_name="0.0.0.0", server_port= 7860) </pre>
Verify	<p>The qabot.py should look like this.</p>	<pre data-bbox="649 1343 714 2025"> 1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8 9. 9 10. 10 11. 11 12. 12 13. 13 14. 14 15. 15 16. 16 17. 17 18. 18 19. 19 20. 20 21. 21 22. 22 23. 23 24. 24 25. 25 26. 26 27. 27 28. 28 29. 29 30. 30 31. 31 32. 32 33. 33 34. 34 </pre>

Package/Method	Description	Code example
	35. 35 36. 36 37. 37 38. 38 39. 39 40. 40 41. 41 42. 42 43. 43 44. 44 45. 45 46. 46 47. 47 48. 48 49. 49 50. 50 51. 51 52. 52 53. 53 54. 54 55. 55 56. 56 57. 57 58. 58 59. 59 60. 60 61. 61 62. 62 63. 63 64. 64 65. 65 66. 66 67. 67 68. 68 69. 69 70. 70 71. 71 72. 72 73. 73 74. 74 75. 75 76. 76 77. 77 78. 78 79. 79 80. 80 81. 81 82. 82 83. 83 84. 84 85. 85 86. 86 87. 87 88. 88 89. 89 90. 90 91. 91 92. 92 93. 93 94. 94 95. 95 96. 96 97. 97 98. 98 99. 99 100. 100 101. 101 102. 102 103. 103 104. 104	

```
1. from ibm_watsonx_ai.foundation_models import ModelInference
2. from ibm_watsonx_ai.metanames import GenTextParamsMetaNames as GenParams
3. from ibm_watsonx_ai.metanames import EmbedTextParamsMetaNames
4. from ibm_watsonx_ai import Credentials
5. from langchain_ibm import WatsonxLLM, WatsonxEmbeddings
6. from langchain.text_splitter import RecursiveCharacterTextSplitter
7. from langchain_community.vectorstores import Chroma
8. from langchain_community.document_loaders import PyPDFLoader
9. from langchain.chains import RetrievalQA
10.
11. import gradio as gr
12. # You can use this section to suppress warnings generated by your code:
13. def warn(*args, **kwargs):
14.     pass
15. import warnings
16. warnings.warn = warn
17. warnings.filterwarnings('ignore')
18.
```

Package/Method	Description	Code example
	<pre> 19. ## LLM 20. def get_llm(): 21. model_id = 'mistralai/mistral-8x7b-instruct-v01' 22. parameters = { 23. GenParams.MAX_NEW_TOKENS: 256, 24. GenParams.TEMPERATURE: 0.5, 25. } 26. project_id = "skills-network" 27. watsonx_llm = WatsonxLLM(28. model_id=model_id, 29. url="https://us-south.ml.cloud.ibm.com", 30. project_id=project_id, 31. params=parameters, 32.) 33. return watsonx_llm 34. 35. ## Document loader 36. def document_loader(file): 37. loader = PyPDFLoader(file.name) 38. loaded_document = loader.load() 39. return loaded_document 40. 41. ## Text splitter 42. def text_splitter(data): 43. text_splitter = RecursiveCharacterTextSplitter(44. chunk_size=1000, 45. chunk_overlap=50, 46. length_function=len, 47.) 48. chunks = text_splitter.split_documents(data) 49. return chunks 50. 51. ## Vector db 52. def vector_database(chunks): 53. embedding_model = watsonx_embedding() 54. vectordb = Chroma.from_documents(chunks, embedding_model) 55. return vectordb 56. 57. ## Embedding model 58. def watsonx_embedding(): 59. embed_params = { 60. EmbedTextParamsMetaNames.TRUNCATE_INPUT_TOKENS: 3, 61. EmbedTextParamsMetaNames.RETURN_OPTIONS: {"input_text": True}, 62. } 63. watsonx_embedding = WatsonxEmbeddings(64. model_id="ibm/slate-125m-english-rtrvr", 65. url="https://us-south.ml.cloud.ibm.com", 66. project_id="skills-network", 67. params=embed_params, 68.) 69. return watsonx_embedding 70. </pre>	

Package/Method**Description****Code example**

```

71. ## Retriever
72. def retriever(file):
73.     splits = document_loader(file)
74.     chunks = text_splitter(splits)
75.     vectordb = vector_database(chunks)
76.     retriever = vectordb.as_retriever()
77.     return retriever
78.

79. ## QA Chain
80. def retriever_qa(file, query):
81.     llm = get_llm()
82.     retriever_obj = retriever(file)
83.     qa = RetrievalQA.from_chain_type(llm=llm,
84.                                         chain_type="stuff",
85.                                         retriever=retriever_obj,
86.                                         return_source_documents=False)
87.     response = qa.invoke(query)
88.     return response['result']
89.

90. # Create Gradio interface
91. rag_application = gr.Interface(
92.     fn=retriever_qa,
93.     allow_flagging="never",
94.     inputs=[
95.         gr.File(label="Upload PDF File", file_count="single", file_types=['.pdf'], type="fi
96.         gr.Textbox(label="Input Query", lines=2, placeholder="Type your question here...")
97.     ],
98.     outputs=gr.Textbox(label="Output"),
99.     title="RAG Chatbot",
100.    description="Upload a PDF document and ask any question. The chatbot will try to answer
101. )
102.

103. # Launch the app
104. rag_application.launch(server_name="0.0.0.0", server_port= 7860)

```

Copied!

1. 1

1. python3.11 qabot.py

Copied!

Serve the application

To serve the application, paste
this code into your Python
terminal:



Skills Network