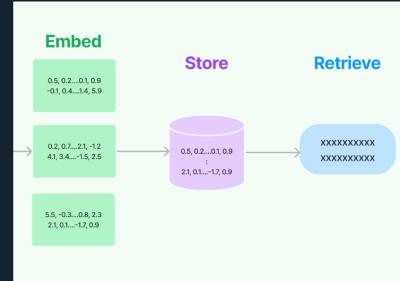


Enterprise Data Science, Machine Learning, and AI / Nov 2, 2023

How to Build a Retrieval-Augmented Generation Chatbot

Andrew Huang and Sophia Yang



Retrieval-augmented generation (RAG) has been empowering conversational AI by allowing models to access and leverage external knowledge bases. In this post, we delve into how to build a RAG chatbot with LangChain and Panel. You will learn:

- What is retrieval-augmented generation (RAG)?
- How to develop a retrieval-augmented generation (RAG) application in LangChain
- How to use Panel's chat interface for our RAG application

By the end of this blog, you will be able to build a RAG chatbot like this:

The screenshot shows a Panel Application window with the following components:

- Left Sidebar:** Labeled "Panel Application". It includes a "Protect your secrets! Make sure you trust the publisher of this app before entering your secrets." message, an "OPENAI API KEY" input field, and a "Number of chunks: 2" slider.
- Central Chat Area:** Shows a conversation between "System" and "User". The System says: "Please first enter an OpenAI Api key and upload a PDF!". The User has uploaded a PDF titled "Artificial Intelligence Index Report 2023" (CHAPTER 1: Research and Development).
- Bottom Input:** A "Pdf" tab and a "TextInput" field with placeholder "Ask questions here!". Buttons for "Send", "Return", "Undo", and "Clear" are at the bottom.

What is Retrieval-Augmented Generation (RAG)?

Are you interested in making a chatbot that can make use of your own collections of data when answering questions? Retrieval-augmented generation (RAG) is an AI framework that combines the strengths of pre-trained language models and information

retrieval systems to generate responses in a conversational AI system or to create content by leveraging external knowledge. It integrates the retrieval of relevant information from a knowledge source and the generation of responses based on that retrieved information.

In a typical RAG setup:

- **Retrieval:** Given a user query or prompt, the system searches through a knowledge source (a vector store with text embeddings) to find relevant documents or text snippets. The retrieval component typically employs some form of similarity or relevance scoring to determine which portions of the knowledge source are most pertinent to the input query.
- **Generation:** The retrieved documents or snippets are then provided to a large language model, which uses them as additional context for generating a more detailed, factual, and relevant response.

RAG can be particularly useful when the pre-trained language model alone may not have the necessary information to generate accurate or sufficiently detailed responses since standard language models like GPT-4 are not capable of accessing real-time or post-training external information directly.

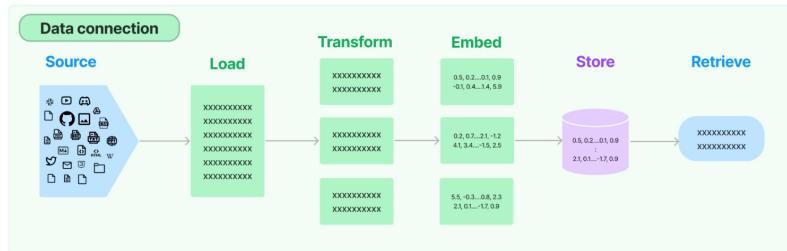
Basic setup

Before we get started in building a RAG application, you will need to install `panel 1.3` and other packages you might need including `jupyterlab`, `pypdf`, `chromadb`, `tiktoken`, `langchain`, and `openai`.

The [complete code for this blog](#) can be found in GitHub.

Developing a Retrieval-Augmented Generation (RAG) application in LangChain

There are actually multiple ways to do RAG in LangChain. Check out our previous blog post, [4 Ways to Do Question Answering in LangChain](#), for details. In this example, we will use the `RetrievalQA` chain. There are several steps in this process:



Source: https://python.langchain.com/docs/modules/data_connection/

- **Load documents:** LangChain provides multiple built-in [document loaders](#), that work with PDF files, JSON files, or a Python file in your file directory. We can use LangChain's PyPDFLoader to import your PDF seamlessly.
 - **Split documents into chunks:** When our document is long, it's necessary to split up our document text into chunks. There are various ways to split your text. Let's just use the simplest method `CharacterTextSplitter` to split based on characters and measure chunk length by the number of characters.
 - **Create text embeddings:** The text chunks are then translated into numerical vectors through embeddings, allowing us to work with text data like semantic search in a computationally efficient manner. We can choose an embedding model provider like OpenAI for this task.
 - **Create a vector store:** We then need to store our embedding vectors in a vector store, which allows us to search and retrieve the relevant vectors at query time.
 - **Create a retriever interface:** We can expose the vector store in a retriever interface. To retrieve text, we can choose a search type like "similarity" to use similarity search in the retriever object where it selects text chunk vectors that are most similar to the question vector. `k=2` lets us find the top 2 most relevant text chunk vectors.
- Create a RetrievalQA chain to answer questions:** A RetrievalQA chain chains a large language model with our retriever interface. You can also define the chain type as one of the four options: "stuff", "map_reduce", "refine", "map_rerank".

- The default `chain_type="stuff"` incorporates ALL text from the documents into the prompt.

The “map_reduce” type breaks texts into groups, poses the question to the LLM for each batch separately, and derives the ultimate answer based on the replies from each batch.

The “refine” type partitions texts into batches, presents the first batch to the LLM, and then submits the answer along with the second batch to the LLM. It progressively refines the answer by processing through all the batches.

The “map-rerank” type divides texts into batches, submits each one to the LLM, returns a score indicating how comprehensively it answers the question, and determines the final answer based on the highest-scoring replies from each batch.

```
import os
from langchain.chains import RetrievalQA
from langchain.document_loaders import PyPDFLoader
from langchain.embeddings import OpenAIEMBEDDINGS
from langchain.llms import OpenAI
from langchain.text_splitter import CharacterTextSplitter
from langchain.vectorstores import Chroma

os.environ["OPENAI_API_KEY"] = "Type your OpenAI API key here"
# load documents
loader = PyPDFLoader("example.pdf")
documents = loader.load()
# split the documents into chunks
text_splitter = CharacterTextSplitter(chunk_size=1000, chunk_overlap=0)
texts = text_splitter.split_documents(documents)
# select which embeddings we want to use
embeddings = OpenAIEMBEDDINGS()
# create the vectorestore to use as the index
db = Chroma.from_documents(texts, embeddings)
# expose this index in a retriever interface
retriever = db.as_retriever(
    search_type="similarity", search_kwargs={"k": 2}
)
# create a chain to answer questions
qa = RetrievalQA.from_chain_type(
    llm=OpenAI(),
    chain_type="map_reduce",
    retriever=retriever,
    return_source_documents=True,
    verbose=True,
)
```

When we ask a question, we can see the result and two source documents:

Using Panel's chat interface for our RAG application

In our [previous blog post](#), we introduced Panop's brand new chat interface and how to build basic AI chatbots in Panop. We

In our [previous blog post](#), we introduced Panel's brand new Chat interface and how to build basic AI chatbots in Panel. We recommend you check out that blog post if you are interested in learning more about the Panel and the chat interface. To make a Panel chatbot for our RAG application, here are four simple steps:

- ➊ Define Panel widgets
- ➋ Wrap LangChain logic into a function
- ➌ Create a chat interface
- ➍ Customize the look with a template

Step 1. Define Panel widgets

Panel widgets are interactive components that allow you to upload files or select values for your application.

For our RAG application chatbot, we define four Panel widgets:

- ➊ pdf_input: To allow users to upload a PDF file.
- ➋ key_input: To input the OpenAI API key.
- ➌ k_slider: To select the number of relevant text chunks.
- ➍ Chain_selection: To select the chain type for retrieval.

```
import panel as pn
pn.extension()

pdf_input = pn.widgets.FileInput(accept=".pdf", value="", height=50)
key_input = pn.widgets.PasswordInput(
    name="OpenAI Key",
    placeholder="sk-...",
)
k_slider = pn.widgets.IntSlider(
    name="Number of Relevant Chunks", start=1, end=5, step=1, value=2
)
chain_select = pn.widgets.RadioButtonGroup(
    name="Chain Type", options=["stuff", "map_reduce", "refine", "map_rerank"]
)
chat_input = pn.widgets.TextInput(placeholder="First, upload a PDF!")
```

Here is what the widget looks like in a Jupyter Notebook:



Step 2: Wrap LangChain Logic into a Function

Next, let's wrap up the LangChain code above into a function. This function should look very familiar to you. It's worth pointing out that we have replaced some values with the widgets we just defined. Specifically:

- We define OpenAI API key with the `key_input` widget
 - We load the file in the `pdf_input` widget
 - We replace `search_kwargs={"k": 2}` with `search_kwargs={"k": k_slider.value}` so that we can control how many relevant docs we'd like to retrieve
 - We replace `chain_type="map_reduce"` with `chain_type=chain_select.value` to allow us to choose one of the four chain types.

```

def initialize_chain():
    if key_input.value:
        os.environ["OPENAI_API_KEY"] = key_input.value

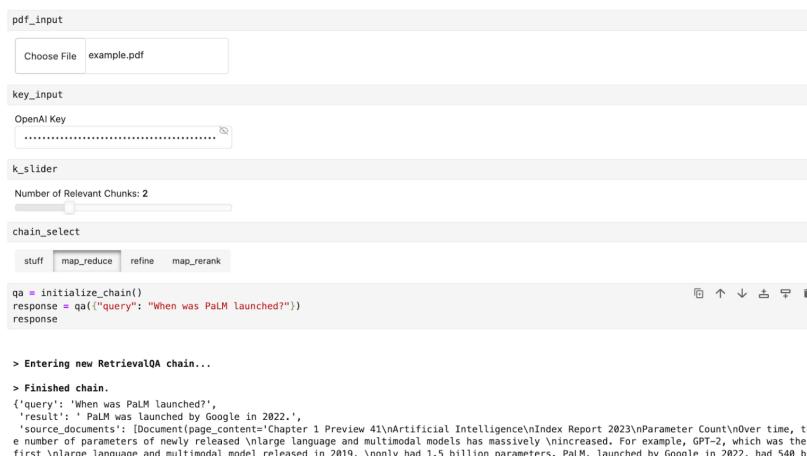
    selections = (pdf_input.value, k_slider.value, chain_select.value)
    if selections in pn.state.cache:
        return pn.state.cache[selections]

    chat_input.placeholder = "Ask questions here!"

    # load document
    with tempfile.NamedTemporaryFile("wb", delete=False) as f:
        f.write(pdf_input.value)
    file_name = f.name
    loader = PyPDFLoader(file_name)
    documents = loader.load()
    # split the documents into chunks
    text_splitter = CharacterTextSplitter(chunk_size=1000, chunk_overlap=0)
    texts = text_splitter.split_documents(documents)
    # select which embeddings we want to use
    embeddings = OpenAIEMBEDDINGS()
    # create the vectorestore to use as the index
    db = Chroma.from_documents(texts, embeddings)
    # expose this index in a retriever interface
    retriever = db.as_retriever(
        search_type="similarity", search_kwargs={"k": k_slider.value}
    )
    # create a chain to answer questions
    qa = RetrievalQA.from_chain_type(
        llm=OpenAI(),
        chain_type=chain_select.value,
        retriever=retriever,
        return_source_documents=True,
        verbose=True,
    )
    return qa

```

After we define the values in the widgets, we can call this function and ask questions about the document we uploaded in the `pdf_input` widget:



lion, nearly 360 times more than GPT-2. The median number of parameters in large language and multimodal models is increasing exponentially over time (Figure 1.2.15). Chapter 1: Research and Development\GPT-2\Grover-Megatron-LM (Original, 8.38B) \n75B \n38T \n11B \nMeenaTuring NLGPT-3 \n175B (davinci) \nGPT-3 \n175B \nNeoX \n20B \nChinchillaPalM (540B) \nDALL-E 2\nStable Diffusion (LM-KL-8-G1OPT-175B) \nJurassic-XM inferno (540B) \nGLM-130BBLOOM \n2019-Feb-2019-Sep2019-Oct-2020-Jan2020-Feb-2020-May2022-Jun-2022-Aug2022-Nov3.2e+81.0e+93.2e+91.0e+183.2e+101.0e+113.2e+111.0e+123.2e+12Number of Parameters (Log Scale)\nNumber of Parameters of Select Large Language and Multimodal Models, 2019-22\nSource: Epoch, 2022 | Chart: 2023 AI Index Report\figure 1.2.15\1.2 Trends in Significant Machine Learning Systems', metadata: {'page': 40, 'source': 'example.pdf'})

Document page content: Chapter 1: Research and Development\GPT-2\Parameter Count\Over time, the number of parameters in newly released large language and multimodal models has massively increased. For example, GPT-3, which was the first large language and multimodal model released in 2019, only had 1.5 billion parameters. PaLM, launched by Google in 2022, had 540 billion, nearly 360 times more than GPT-2. The median number of parameters in large language and multimodal models is increasing exponentially over time (Figure 1.2.15). Chapter 1: Research and Development\GPT-2\Grover-Megatron-LM (Original, 8.38B) \n75B \n38T \n11B \nMeenaTuring NLGPT-3 \n175B (davinci) \nERNIE-GEN (large) \nDALL-E \nWenYinPanGu \nGPT-J \n6BHyperLoRA \nCogyle wu Dao 2.0 \nERNIE 3.0 \nCodexJurassic-1-JumboMegatron-Turing NLG 530B \nGopher \nGPT-NeoX-20B \nChinchillaPalM (540B) \nDALL-E 2\nStable Diffusion (LM-KL-8-G1OPT-175B) \nJurassic-XMInferno (540B) \nGLM-130BBLOOM \n2019-Feb-2019-Sep2019-Oct-2020-Jan2020-Feb-2020-May2022-Jun-2022-Aug2022-Nov3.2e+81.0e+93.2e+91.0e+183.2e+101.0e+113.2e+111.0e+123.2e+12Number of Parameters (Log Scale)\nNumber of Parameters of Select Large Language and Multimodal Models, 2019-22\nSource: Epoch, 2022 | Chart: 2023 AI Index Report\figure 1.2.15\1.2 Trends in Significant Machine Learning Systems', metadata: {'page': 40, 'source': '/var/folders/dg/23jtng_n4h1295wnm3_wkm8000gp/7tmp_kd4t1pr'})

Step 3. Create a chat interface

How do we interact with our documents and ask questions in a chat interface? This is where Panel's ChatInterface widget comes in!

We must create a function `respond` to define how the chatbot responds. This function takes the response from Step 2 and formats it into a Panel object `answers`. We also append the relevant source documents in this Panel object. Then we can simply call the function in the `pn.chat.ChatInterface` `callback`.

```
async def respond(contents, user, chat_interface):
    if not pdf_input.value:
        chat_interface.send(
            {"user": "System", "value": "Please first upload a PDF!"}, respond=False
        )
        return
    elif chat_interface.active == 0:
        chat_interface.active = 1
        chat_interface.active_widget.placeholder = "Ask questions here!"
        yield {"user": "OpenAI", "value": "Let's chat about the PDF!"}
        return

    qa = initialize_chain()
    response = qa({"query": contents})
    answers = pn.Column(response["result"])
    answers.append(pn.layout.Divider())
    for doc in response["source_documents"][:-1]:
        answers.append(f"**Page {doc.metadata['page']}**:")
        answers.append(f"```\n{doc.page_content}\n```")
    yield {"user": "OpenAI", "value": answers}

    chat_interface = pn.chat.ChatInterface(
        callback=respond, sizing_mode="stretch_width", widgets=[pdf_input, chat_input]
    )
    chat_interface.send(
        {"user": "System", "value": "Please first upload a PDF and click send!"},
        respond=False,
    )
```

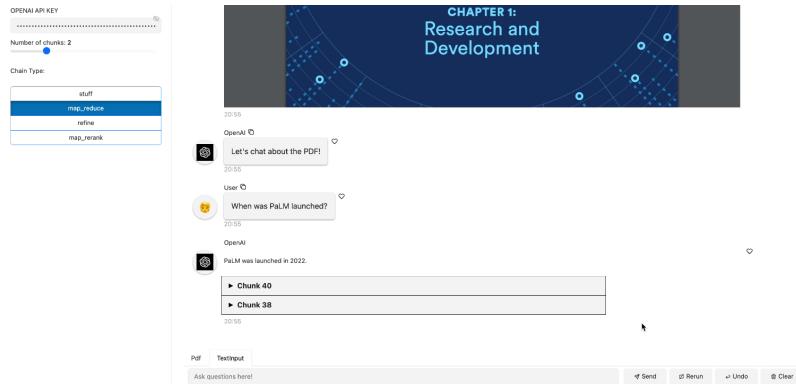
Step 4. Customize the look with a template

The final step is to combine the widget and the chat interface together in an application. Panel comes with multiple templates that allow us to quickly and easily create web apps with better aesthetics. Here we use the `BootstrapTemplate` to organize the widgets in the sidebar and display the chat interface in the center of the app.

```
template = pn.template.BootstrapTemplate(
    sidebar=[key_input, k_slider, chain_select], main=[chat_interface]
)
template.servable()
```

To serve the app, run `panel serve app.py` or `panel serve app.ipynb`. You will get the app shown at the beginning of this blog post:





Conclusion

Retrieval-augmented generation (RAG) is a fascinating blend of information retrieval and generative techniques in the AI landscape. In this blog, we've broken down its essentials, walked through creating a RAG application using LangChain, and capped it off with integrating Panel's user-friendly chat interface. This serves as a practical guide for anyone looking to understand or implement RAG in their projects. As we continue to navigate the evolving tech world, understanding and utilizing tools like RAG can be a beneficial step forward. We hope you found value in this introduction and guide. Happy coding!

Note:

- The [complete code for this blog](#) can be found in GitHub (updated with caching):

- All these tools are open source and free for everyone to use, but if you'd like some help getting started from Anaconda's AI and Python app experts, reach out to sales@anaconda.com.

You may also be interested in:

Enterprise Data Science, Machine Learning, and AI

Oct 25, 2023

How to Build Your Own Panel AI Chatbots

News Product Updates

Oct 10, 2023

Anaconda Learning: Turbocharge your Python Journey in Anaconda Notebooks

Enterprise Data Science, Machine Learning, and AI

Sep 28, 2023

Introducing Pandata: The Scalable Open-Source Analysis Stack

Let's Connect

Get in touch to learn more about Anaconda.

[Contact Us >](#)



Why Anaconda?

Data Science & AI Platform

Professional Services

Pricing

Contact Sales

Solutions

BY NEED

Managing Security & Compliance

BY ROLE

Practitioners

IT

About

About Anaconda

Careers

Education

Partners

Our Open-Source Commitment

Customer Reference Program

Press

Events

Contact Us

Resources

Libraries & Packages

Support Center

Open Source

Blog

Resource Center

© 2024 Anaconda Inc. All rights reserved.

[Service Status](#) / [Legal](#) / [Privacy Policy](#) / [Terms of Use](#)



Hey! 🌟 Welcome to Anaconda. I'm here to help. What are you looking for today?

