

build_chat()

Builds and returns a Conversational Retrieval Chain to answer a user's question

The chain should support streaming text generation

The chain should use random variations of component parts

Input

How much?

Conversation
Buffer Memory

Condense Question Chain

Prompt

Given the following conversation and a follow up question, rephrase the follow up question to be a standalone question, in its original language.

Chat History:
Human: What country produces the most spice?
Assistant: India produces the most spice
Follow Up Input: How much?
Standalone question:

OpenAI
GPT-3.5-Turbo

Exactly how much
spice does India
produce?

Refined Question

Combine Docs Chain

Exactly how much spice
does India produce?

[-.31, .45, .91, ...]

**Pinecone
Retriever**

in 2011 india
produced 1.5m
tons of spice

Use the
following
context to
answer the
user's question.
Context: In 2011,
India produced
1.5m tons of
spice
User's question:
Exactly how

OpenAI
GPT-3.5-Turbo

India
produced
1.5m tons in

RetrievalChain

Retriever

Memory

LLM

Retrievers

Pinecone Retriever

Chroma Retriever

Memory

Conversation Buffer Memory

*Conversation Buffer Window
Memory*

LLM

GPT-3.5-turbo

GPT-4

Input

How much?

Retrievers

Pinecone Retriever

LLM

GPT-3.5-turbo

Memory

Conversation Buffer Memory

Conversation Buffer Window Memory

Conversation Buffer Window Memory

Condense Question Chain

Prompt

Given the following conversation and a follow up question, rephrase the follow up question to be a standalone question, in its original language.

Chat History:
Human: What country produces the most spice?
Assistant: India produces the most spice
Follow Up Input: How much?
Standalone question:

GPT-4

Exactly how much spice does India produce?

Refined Question

Combine Docs Chain

Exactly how much spice does India produce?

[-.31, .45, .91, ...]

Chroma Retriever

in 2011 india produced 1.5m tons of spice

Use the following context to answer the user's question.
Context: In 2011, India produced 1.5m tons of spice
User's question: Exactly how

GPT-4

India produced 1.5m tons in

What is this PDF about?

RetrievalChain

retriever that fetches 2 relevant documents

Conversation Buffer Memory

GPT-3.5-turbo

Bad Answer

This PDF is, like, about stuff you put on food or something.

What is this PDF about?

RetrievalChain

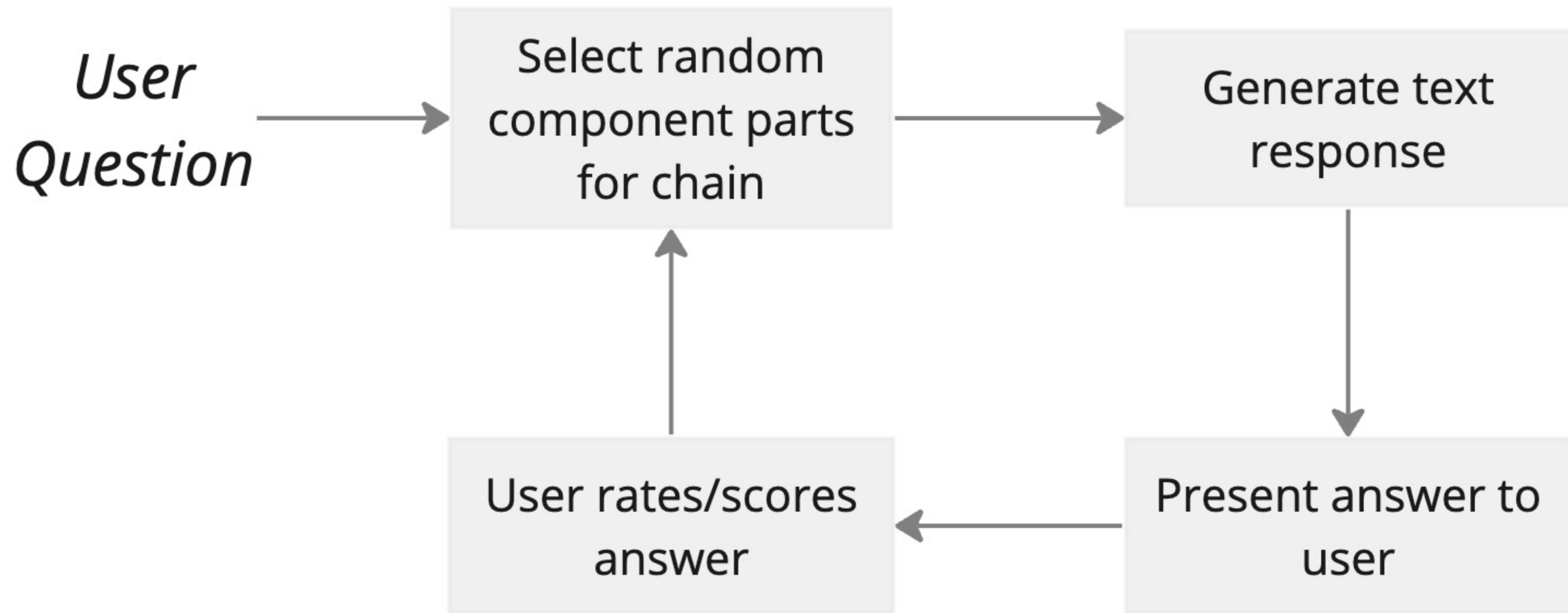
Pinecone Retriever

Conversation Buffer Window Memory

GPT-4

Good Answer

This PDF is about the production and usage of spice. It also focuses on the origins and top producers of spice.



Component Maps

Retrievers

name	component
pinecone_2	<i>Pinecone Retriever that fetches top 2 docs</i>
pinecone_3	<i>Pinecone Retriever that fetches top 3 docs</i>

Memory

name	component
buffer	<i>Conversation Buffer Memory</i>
window	<i>Conversation Buffer Window</i>

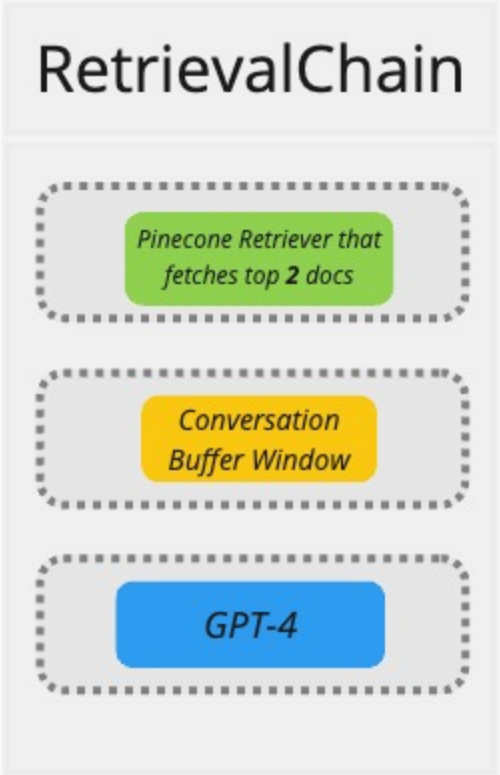
LLM

name	component
gpt-4	<i>GPT-4</i>
gpt-3.5	<i>GPT-3.5-turbo</i>

Dictionaries where the key is an identifier for a component configured in a particular way

First message of conversation

User Message tied to Conversation '123'



Retrievers

name	component
pinecone_2	Pinecone Retriever that fetches top 2 docs
pinecone_3	Pinecone Retriever that fetches top 3 docs

Memory

Conversation Buffer Window

LLM

name	component
gpt-4	GPT-4
gpt-3.5	GPT-3.5-turbo

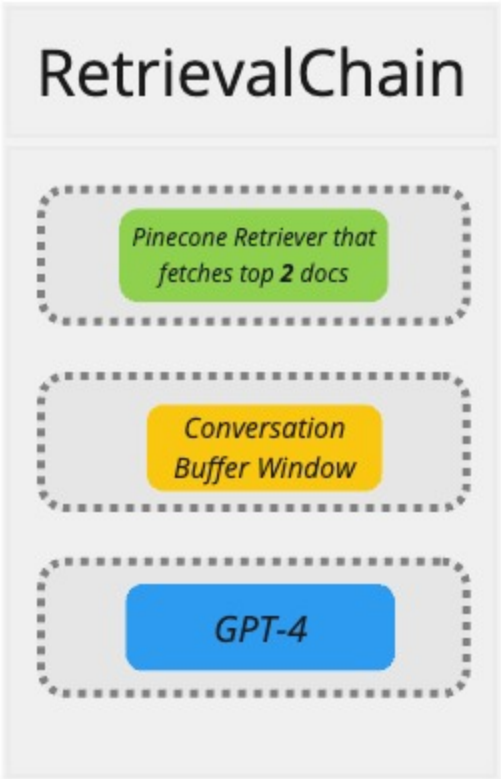
SQLite Database List of Conversations

id	llm	memory	retriever
123	gpt-4	window	pinecone_2

This table already exists in the database

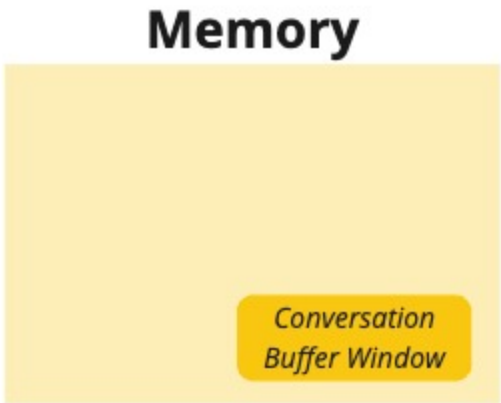
*Second, third, fourth, etc
message of conversation*

User Message
tied to
Conversation
'123'



Retrievers

name	component
pinecone_2	Pinecone Retriever that fetches top 2 docs
pinecone_3	Pinecone Retriever that fetches top 3 docs



LLM

name	component
gpt-4	GPT-4
gpt-3.5	GPT-3.5-turbo



SQLite Database List of Conversations

id	llm	memory	retriever
123	gpt-4	window	pinecone_2

app.web Module

app.chat

get_messages_by_conversation_id()

add_message_to_conversation()

set_conversation_components()

get_conversation_components()

*Updates a conversation
with the names of
components we used*

*Retrieves the names of
components used in a
conversation*

```
1  def build_chat(chat_args):
2      # Assume it is a follow-up message
3      components = get_conversation_components(chat_args.conversation_id)
4
5      print(components) # -> { "llm": "gpt-3.5-turbo" }
6
7      if components["llm"] == "gpt-3.5-turbo":
8          retriever = build_retriever(
9              chat_args,
10             k=2
11         )
12     elif components["llm"] == "gpt-4":
13         retriever = build_retriever(
14             chat_args,
15             k=3
16         )
```



```
1  def build_retriever(chat_args, k):
2      search_kwargs = {
3          "filter": { "pdf_id": chat_args.pdf_id },
4          "k": k
5      }
6      return vector_store.as_retriever(
7          search_kwargs=search_kwargs
8      )
```


Different Retrievers might require different config

name	component
pinecone_2	Pinecone Retriever that fetches top 2 docs
pinecone_3	Pinecone Retriever that fetches top 3 docs
chroma_spice_2	Chroma Retriever that finds first 2 docs that contain the word "spice"
chroma_3	Chroma Retriever that finds first 3 docs

We aren't going to add in Chroma, this is just an example


```
1  {
2    "filter": {"pdf_id": chat_args.pdf_id },
3    "k": 2
4  }
```

```
1  {
2    "filter": {"pdf_id": chat_args.pdf_id },
3    "k": 3
4  }
```

```
1  {
2    "where": {"pdf_id": chat_args.pdf_id}
3    "where_doc": {"$contains": {"text": "spice"} }
4  }
```

```
1  {
2    "where": {"pdf_id": chat_args.pdf_id}
3    "where_doc": {"$contains": {"text": "technical"}}
4  }
```

```
1 def build_chat(chat_args):
2     conversation_id = chat_args.conversation_id
3     # Assume it is a follow-up message
4     components = get_conversation_components(conversation_id
5     )
6
7     print(components) # -> { "retriever": "pinecone_2" }
8
9     if components["retriever"] == "pinecone_2":
10         retriever = build_pinecone_retriever(
11             chat_args,
12             k=2
13         )
14     elif components["retriever"] == "pinecone_3":
15         retriever = build_pinecone_retriever(
16             chat_args,
17             k=3
18         )
19     elif components["retriever"] == "chroma_spice_2":
20         retriever = build_chroma_retriever(
21             chat_args,
22             n_results=2,
23             where_doc={"$contains": {"text": "spice" }}
24         )
25     elif components["retriever"] == "chroma_3":
26         retriever = build_chroma_retriever(
27             chat_args,
28             n_results=3
29         )
```



```
1 def build_pinecone_retriever(chat_args, k):
2     search_kwargs = {
3         "filter": { "pdf_id": chat_args.pdf_id },
4         "k": k
5     }
6     return vector_store.as_retriever(
7         search_kwargs=search_kwargs
8     )
```



```
1 def build_chroma_retriever(chat_args, n_results, where_doc):
2     search_kwargs = {
3         "where": { "pdf_id": chat_args.pdf_id },
4         "n_results": n_results,
5         "where_document": where_doc
6     }
7     return vector_store.as_retriever(
8         search_kwargs=search_kwargs
9     )
```



```

1  def build_chat(chat_args):
2      conversation_id = chat_args.conversation_id
3      # Assume it is a follow-up message
4      components = get_conversation_components(conversation_id
5      )
6
7      print(components) # -> { "retriever": "pinecone_2" }
8
9      if components["retriever"] == "pinecone_2":
10         retriever = build_pinecone_retriever(
11             chat_args,
12             k=2
13         )
14     elif components["retriever"] == "pinecone_3":
15         retriever = build_pinecone_retriever(
16             chat_args,
17             k=3
18         )
19     elif components["retriever"] == "chroma_spice_2":
20         retriever = build_chroma_retriever(
21             chat_args,
22             n_results=2,
23             where_doc={"$contains": {"text": "spice" }}
24         )
25     elif components["retriever"] == "chroma_3":
26         retriever = build_chroma_retriever(
27             chat_args,
28             n_results=3
29         )

```

**Many arguments
are hard coded and
tied to the name of
the component**