

LangChain

Chat with Your Data

Harrison Chase



LangChain

Overview



Open-source developer framework for building LLM applications

Python and TypeScript packages

Focused on composition and modularity

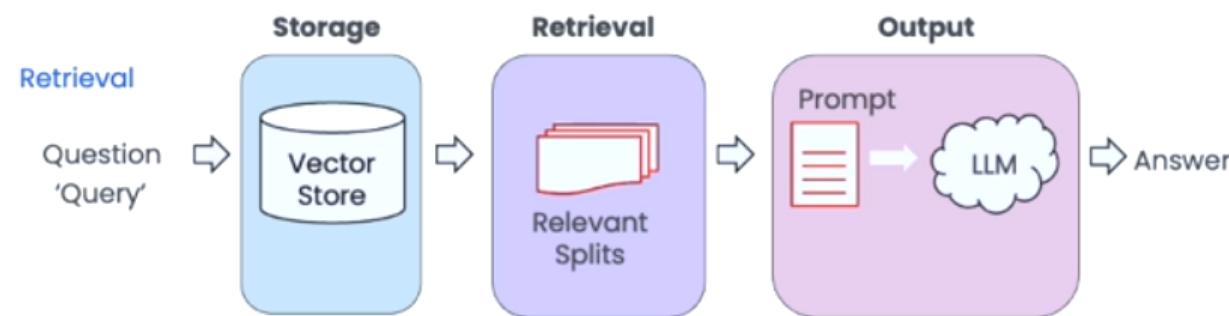
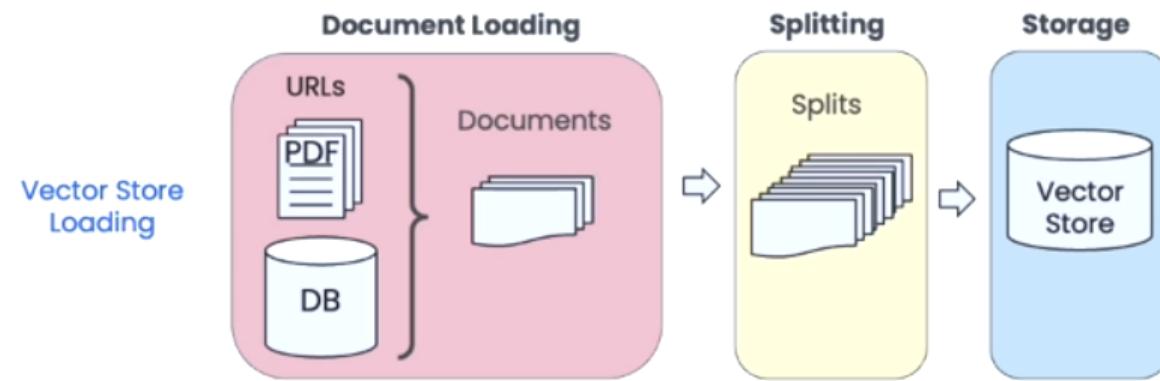
Key value adds:

1. Modular components (and implementations of those components)
2. Use cases - common ways to combine those components together

Components

- **Prompts**
 - Prompt Templates
 - Output Parsers: 5+ implementations
 - Retry/fixing logic
 - Example Selectors: 5+ implementations
- **Models**
 - LLM's: 20+ integrations
 - Chat Models
 - Text Embedding Models: 10+ integrations
- **Indexes**
 - Document Loaders: 50+ implementations
 - Text Splitters: 10+ implementations
 - Vector stores: 10+ integrations
 - Retrievers: 5+ integrations/implementations
- **Chains**
 - Can be used as building blocks for other chains
 - More application specific chains: 20+ different types
- **Agents**
 - Agent Types: 5+ types
 - Algorithms for getting LLMs to use tools
 - Agent Toolkits: 10+ implementations
 - Agents armed with specific tools for a specific application

Retrieval Augmented Generation



Acknowledgments



Ankush Gola
LangChain



Lance Martin
LangChain



Geoff Ladwig
DeepLearning.AI



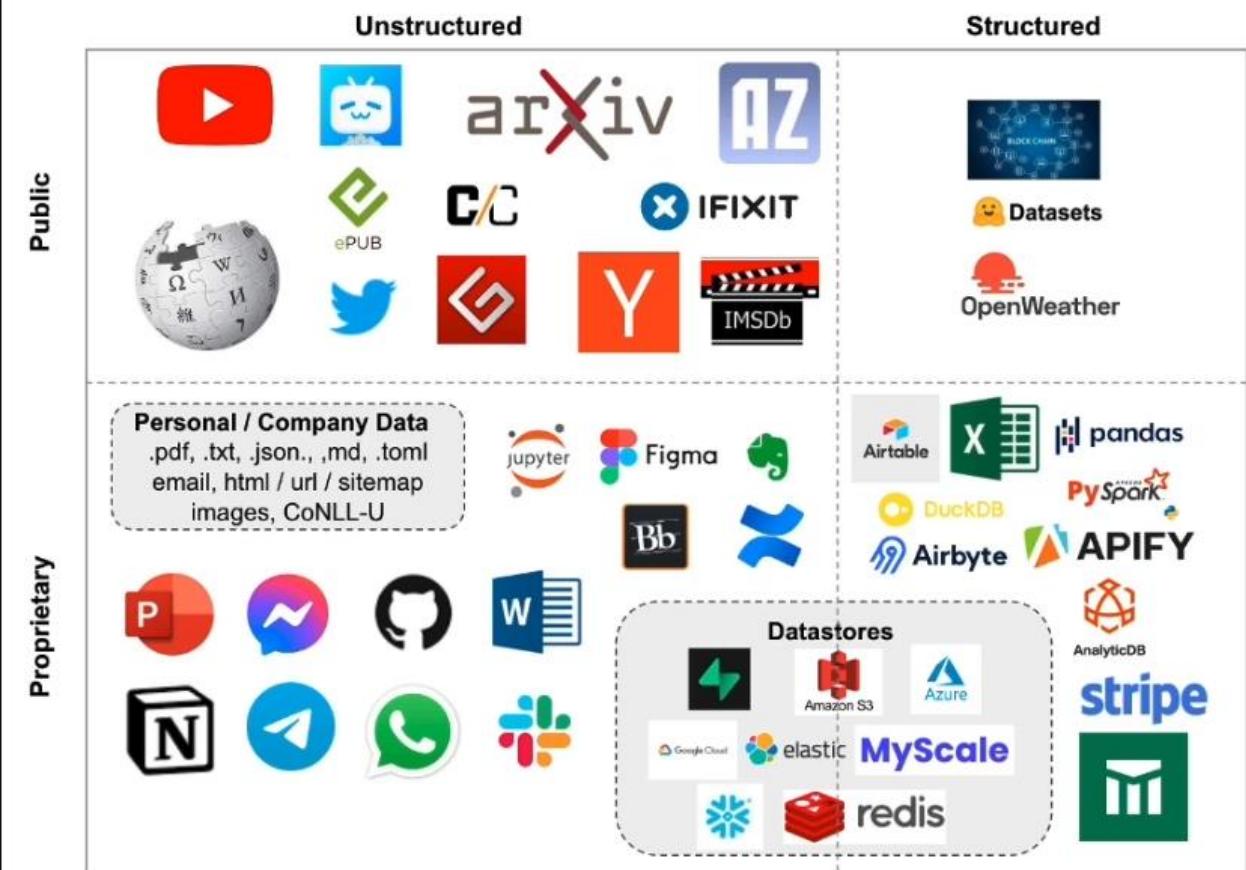
Diala Ezzeddine
DeepLearning.AI

Loaders

- Loaders deal with the specifics of accessing and converting data
 - Accessing
 - Web Sites
 - Data Bases
 - YouTube
 - arXiv
 - ...
 - Data Types
 - PDF
 - HTML
 - JSON
 - Word, PowerPoint...
- Returns a list of `Document` objects:

```
[  
Document(page_content='MachineLearning-Lecture01 \nInstructor (Andrew Ng): Okay.  
Good morning. Welcome to CS229....',  
metadata={'source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 0})  
...  
Document(page_content='[End of Audio] \nDuration: 69 minutes ',  
metadata={'source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 21})  
]
```

Document Loaders



```
from langchain.document_loaders import PyPDFLoader
```

```
loader = PyPDFLoader("docs/cs229_lectures/MachineLearning-Lecture01.pdf")
```

```
pages = loader.load()
```

```
len(pages)
```

```
22
```

```
page = pages[0]
```

```
print(page.page_content[:500])
```

MachineLearning-Lecture01

Instructor (Andrew Ng): Okay. Good morning. Welcome to CS229, the machine learning class. So what I wanna do today is just spend a little time going over the logistics

of the class, and then we'll start to talk a bit about machine learning. By way of introduction, my name's Andrew Ng and I'll be instructor for this class. And so

I personally work in machine learning, and I've worked on it for about 15 years now, and I actually think that machine learning is

```
page.metadata
```

```
{'source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 0}
```

I personally work in machine learning, and I've worked on it for about 15 years now, and I actually think that machine learning is

```
page.metadata
```

```
{'source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 0}
```

```
from langchain.document_loaders.generic import GenericLoader
from langchain.document_loaders.parsers import OpenAIWhisperParser
from langchain.document_loaders.blob_loaders.youtube_audio import YoutubeAudio
```

```
url="https://www.youtube.com/watch?v=jGw0_UgTS7I"
```

```
save_dir="docs/youtube/"
```

```
loader = GenericLoader(
```

```
YoutubeAudioLoader([url], save_dir),  
OpenAIWhisperParser()
```

```
)
```

```
docs = loader.load()
```

```
[youtube] Extracting URL: https://www.youtube.com/watch?v=jGw0_UgTS7I
```

```
[youtube] jGw0_UgTS7I: Downloading webpage
```

```
[youtube] jGw0_UgTS7I: Downloading android player API JSON
```

```
[info] jGw0_UgTS7I: Downloading 1 format(s): 140
```

```
[download] docs/youtube//Stanford CS229: Machine Learning Course, Lecture 1 – Andrew Ng (Autumn 2018).m4a has already been downloaded
```

```
[download] 100% of 69.71MiB
```

```
[ExtractAudio] Not converting audio docs/youtube//Stanford CS229: Machine Learning Course, Lecture 1 – Andrew Ng (Autumn 2018).m4a; file is already in target format m4a
```

Transcribing part 1! 



```
docs[0].page_content[0:500]
```

"Welcome to CS229 Machine Learning. Some of you know that this is a class that's taught at Stanford for a long time. And this is often the class that, um, I most look forward to teaching each year because this is where we've helped, I think several generations of Stanford students become experts in machine learning, got- built many of their  products and services and startups that I'm sure, many of you or probably all of you are using, uh, uh, today. Um, so what I want to do today was spend some t"



Transcript is done via OpenAI Whisper
<https://openai.com/research/whisper>

Whisper is an automatic speech recognition (ASR) system trained on 680,000 hours of multilingual and multitask supervised data collected from the web. We show that the use of such a large and diverse dataset leads to improved robustness to accents, background noise and technical language. Moreover, it enables transcription in multiple languages, as well as translation from those languages into English. We are open-sourcing models and inference code to serve as a foundation for building useful applications and for further research on robust speech processing.

Run it locally:

<https://www.youtube.com/watch?v=U2TqRKjw1-k>

Multitask training data (680k hours)

English transcription

"Ask not what your country can do for ..."

Ask not what your country can do for ...

Any-to-English speech translation

"El rápido zorro marrón salta sobre ..."

The quick brown fox jumps over ...

Non-English transcription

"언덕 위에 올라 내려다보면 너무나 넓고 넓은 ..."

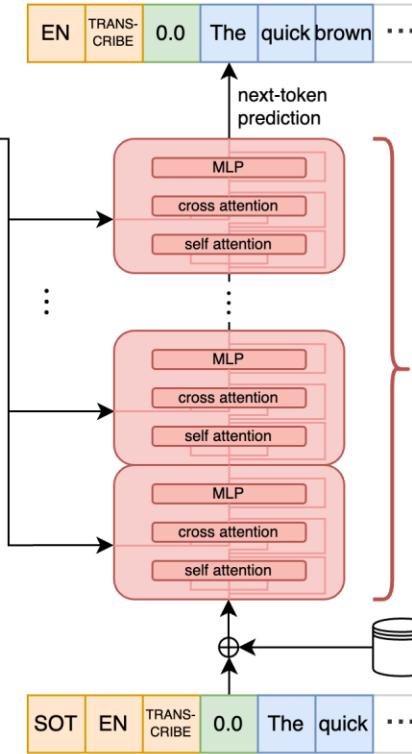
언덕 위에 올라 내려다보면 너무나 넓고 넓은 ...

No speech

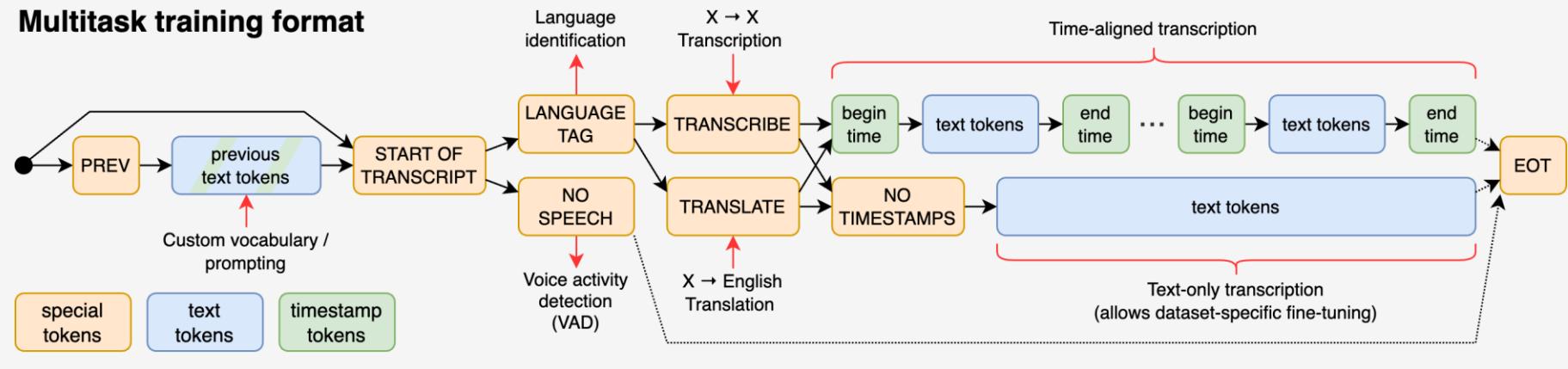
(background music playing)

Ø

Sequence-to-sequence learning



Multitask training format

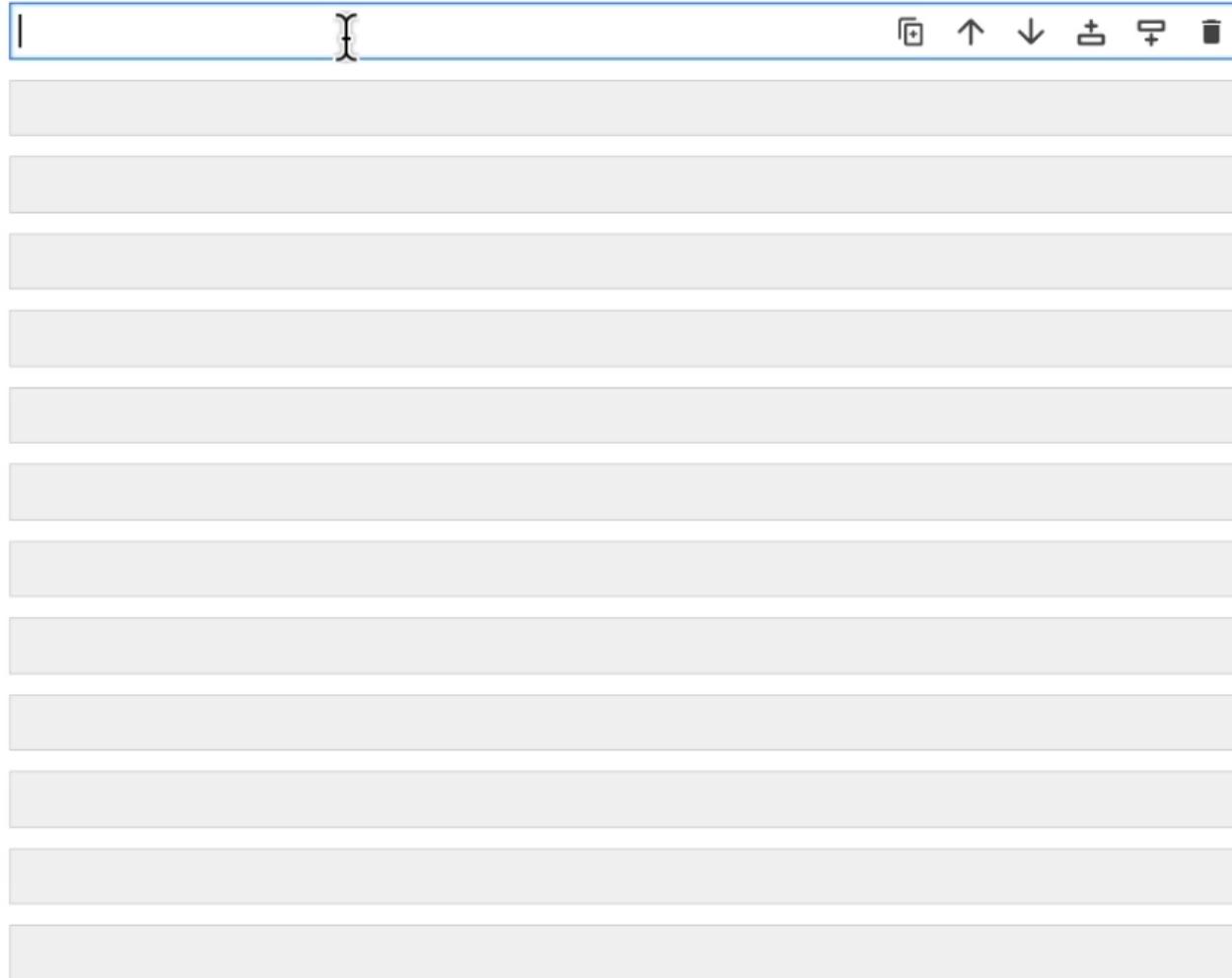


```
from langchain.document_loaders import NotionDirectoryLoader
loader = NotionDirectoryLoader("docs/Notion_DB")
docs = loader.load()
```

```
print(docs[0].page_content[0:200])
```

```
# Blendle's Employee Handbook
```

This is a living document with everything we've learned working with people while running a startup. And, of course, we continue to learn. Therefore it's a document that



LangChain

Chat with Your Data

Document Splitting



LangChain



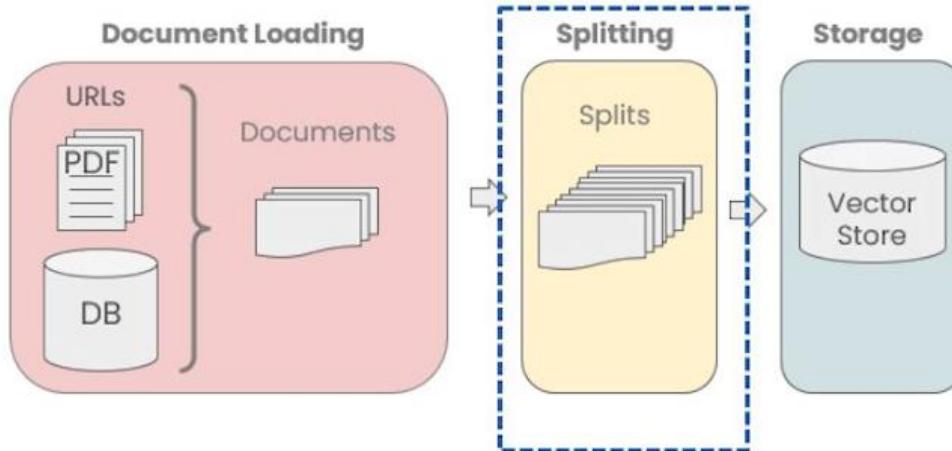
DeepLearning.AI

Document Splitting

It's tricky to split and not to end with sentence splitting to get full information.

Sometime its splitting with chunk size and overlapping

- Splitting Documents into smaller chunks
 - Retaining meaningful relationships!



...
on this model. The Toyota Camry has a head-snapping
80 HP and an eight-speed automatic transmission that will

...
Chunk 1: on this model. The Toyota Camry has a head-snapping

Chunk 2: 80 HP and an eight-speed automatic transmission that will

Question: What are the specifications on the Camry?

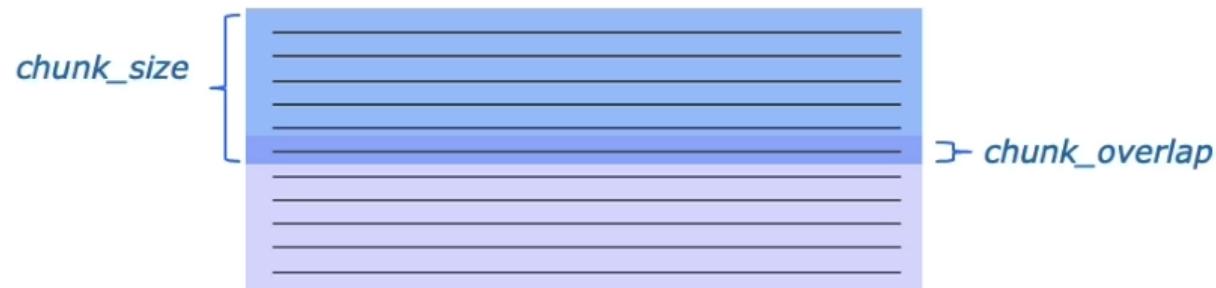
Example Splitter

```
langchain.text_splitter.CharacterTextSplitter(  
    separator: str = "\n\n"  
    chunk_size=4000,  
    chunk_overlap=200,  
    length_function=<builtin function len>,  
)
```

Methods:

`create_documents()` - Create documents from a list of texts.

`split_documents()` - Split documents.



Types of splitters

`langchain.text_splitter.`

- **CharacterTextSplitter()**- Implementation of splitting text that looks at characters.
- **MarkdownHeaderTextSplitter()** - Implementation of splitting markdown files based on specified headers.
- **TokenTextSplitter()** - Implementation of splitting text that looks at tokens.
- **SentenceTransformersTokenTextSplitter()** - Implementation of splitting text that looks at tokens.
- **RecursiveCharacterTextSplitter()** - Implementation of splitting text that looks at characters. Recursively tries to split by different characters to find one that works.
- **Language()** – for CPP, Python, Ruby, Markdown etc
- **NLTKTextSplitter()** - Implementation of splitting text that looks at sentences using NLTK (Natural Language Tool Kit)
- **SpacyTextSplitter()** - Implementation of splitting text that looks at sentences using Spacy

```
_ = load_dotenv(find_dotenv()) # read local .env file
openai.api_key = os.environ['OPENAI_API_KEY']

from langchain.text_splitter import RecursiveCharacterTextSplitter, CharacterT
[REDACTED]
chunk_size = 26
chunk_overlap = 4
[REDACTED]
r_splitter = RecursiveCharacterTextSplitter(chunk_size=chunk_size,chunk_overl
c_splitter = CharacterTextSplitter(chunk_size=chunk_size,chunk_overlap=chunk_o
[REDACTED]
text1 = 'abcdefghijklmnopqrstuvwxyz'
[REDACTED]
r_splitter.split_text(text1)
[REDACTED]
['abcdefghijklmnopqrstuvwxyz']
[REDACTED]
text2 = 'abcdefghijklmnopqrstuvwxyzabcdefghijklmnopqrstuvwxyz'
[REDACTED]
r_splitter.split_text(text2)
[REDACTED]
['abcdefghijklmnopqrstuvwxyz', 'wxyzabcdefghijklmnopqrstuvwxyz']
[REDACTED]
text3 = "a b c d e f g h i j k l m n o p q r s t u v w x y z"
[REDACTED]
r_splitter.split_text(text3)
[REDACTED]
['a b c d e f g h i j k l m', 'l m n o p q r s t u v w x', 'w x y z']
```

→ Splitting characters by position, can be empty - space

→ There are 3 chunks with overlapping characters like „l m” , „w x”

```
['abcdefghijklmnopqrstuvwxyz', 'wxyzabcdefg']
```

```
text3 = "a b c d e f g h i j k l m n o p q r s t u v w x y z"
```

```
r_splitter.split_text(text3)
```

```
['a b c d e f g h i j k l m', 'l m n o p q r s t u v w x', 'w x y z']
```

```
c_splitter.split_text(text3)
```

```
['a b c d e f g h i j k l m n o p q r s t u v w x y z']
```

```
c_splitter = CharacterTextSplitter(  
    chunk_size=chunk_size,  
    chunk_overlap=chunk_overlap,  
    separator = ''  
)
```

```
c_splitter.split_text(text3)
```

```
['a b c d e f g h i j k l m', 'l m n o p q r s t u v w x', 'w x y z']
```

Splitting characters, can be empty - space

Splitting characters, separator can be empty space, so there is no empty space in from or end of a chunk

```
some_text = """When writing documents, writers will use document structure to group content. This can convey to the reader, which idea's are related. For example, closely related ideas are in sentences. Similar ideas are in paragraphs. Paragraphs form a document. Paragraphs are often delimited with a carriage return or two carriage returns. Carriage returns are the "backslash n" you see embedded in this string. \n Sentences have a period at the end, but also, have a space.\n and words are separated by space."""
```

```
len(some_text)
```

496

```
c_splitter = CharacterTextSplitter(  
    chunk_size=450,  
    chunk_overlap=0,  
    separator = ''  
)  
r_splitter = RecursiveCharacterTextSplitter(  
    chunk_size=450,  
    chunk_overlap=0,  
    separators=["\n\n", "\n", " ", ""]  
)
```

```
c_splitter.split_text(some_text)
```

```
['When writing documents, writers will use document structure to group content. This can convey to the reader, which idea's are related. For example, closely related ideas are in sentences. Similar ideas are in paragraphs. Paragraphs form a document. \n\n Paragraphs are often delimited with a carriage return or two carriage returns. Carriage returns are the "backslash n" you see embedded in this string. Sentences have a period at the end, but also, have a space. and words are separated by space.']
```

→ Recursive splitter first splits by “\n\n” and if it needs further to split then it does it with “\n” and further with “ ”.

It splits with the “ ” and creates

```
'have a space.and words are separated by space.]
```

```
r_splitter.split_text(some_text)
```

["When writing documents, writers will use document structure to group content. This can convey to the reader, which idea's are related. For example, closely related ideas are in sentences. Similar ideas are in paragraphs. Paragraphs form a document."],

'Paragraphs are often delimited with a carriage return or two carriage returns. Carriage returns are the "backslash n" you see embedded in this string. Sentences have a period at the end, but also, have a space.and words are separated by space.]'



Recursive splitter first splits by "\n\n" and if it needs further to split then it does it with "\n" and further with "..". It's properly splitted in two chunks just before 450 characters.

```
chunk_overlap=0,  
separators=["\n\n", "\n", "\.", " ", ""]  
)  
  
r_splitter.split_text(some_text)  
  
["When writing documents, writers will use document structure to group content. This can convey to the reader, which idea's are related.",  
'. For example, closely related ideas are in sentences. Similar ideas are in paragraphs. Paragraphs form a document.',  
'Paragraphs are often delimited with a carriage return or two carriage returns.',  
. Carriage returns are the "backslash n" you see embedded in this string.',  
. Sentences have a period at the end, but also, have a space. and words are separated by space.]
```

Not properly split with the “.”

```
r_splitter = RecursiveCharacterTextSplitter(  
    chunk_size=150,  
    chunk_overlap=0,  
    separators=["\n\n", "\n", "(?<=\. )", " ", ""]  
)  
  
r_splitter.split_text(some_text)  
  
["When writing documents, writers will use document structure to group content. This can convey to the reader, which idea's are related.",  
'. For example, closely related ideas are in sentences. Similar ideas are in paragraphs. Paragraphs form a document.',  
'Paragraphs are often delimited with a carriage return or two carriage returns.',  
. Carriage returns are the "backslash n" you see embedded in this string.',  
. Sentences have a period at the end, but also, have a space. and words are separated by space.]
```

With regex separator sentences are properly splitted.



```
from langchain.document_loaders import PyPDFLoader  
loader = PyPDFLoader("docs/cs229_lectures/MachineLearning-Lecture01.pdf")  
pages = loader.load()
```

Splitting PDF documents

```
from langchain.text_splitter import CharacterTextSplitter  
text_splitter = CharacterTextSplitter(  
    separator="\n",  
    chunk_size=1000,  
    chunk_overlap=150,  
    length_function=len  
)
```

```
docs = text_splitter.split_documents(pages)
```

```
len(docs)
```

```
77
```

```
len(pages)
```

```
22
```

Using method `split_documents` for PDFs

```
loader = PyPDFLoader("docs/cs229_lectures/MachineLearning-Lecture01.pdf")
pages = loader.load()
```

```
from langchain.text_splitter import CharacterTextSplitter
text_splitter = CharacterTextSplitter(
    separator="\n",
    chunk_size=1000,
    chunk_overlap=150,
    length_function=len
)
```

```
docs = text_splitter.split_documents(pages)
```

```
len(docs)
```

```
77
```

```
len(pages)
```

```
22
```

```
from langchain.document_loaders import NotionDirectoryLoader
loader = NotionDirectoryLoader("docs/Notion_DB")
notion_db = loader.load()
```

```
docs = text_splitter.split_documents(notion_db)
```

```
len(notion_db)
```

```
52
```

```
len(docs)
```

```
353
```



Using method for database loading and splitting

```
from langchain.text_splitter import TokenTextSplitter

text_splitter = TokenTextSplitter(chunk_size=1, chunk_overlap=0)

text1 = "foo bar bazzyfoo"

text_splitter.split_text(text1)

['foo', ' bar', ' b', 'az', 'zy', 'foo']

text_splitter = TokenTextSplitter(chunk_size=10, chunk_overlap=0)

docs = text_splitter.split_documents(pages)

docs[0]

Document(page_content='MachineLearning-Lecture01 \n', metadata={'source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 0})

pages[0].metadata

{'source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 0}
```

Token splitting

We can also split on token count explicitly, if we want. This can be useful because LLMs often have context windows designated in tokens. Tokens are often ~4 characters.

Metadata of the source is preserved, so for QA we can do a reference to a document from a chunk.

```
from langchain.text_splitter import MarkdownHeaderTextSplitter

markdown_document = """# Title\n\n \
## Chapter 1\n\n \
Hi this is Jim\n\n Hi this is Joe\n\n \
### Section \n\n \
Hi this is Lance \n\n \
## Chapter 2\n\n \
Hi this is Molly"""

headers_to_split_on = [
    ("#", "Header 1"),
    ("##", "Header 2"),
    ("###", "Header 3"),
]

markdown_splitter = MarkdownHeaderTextSplitter(
    headers_to_split_on=headers_to_split_on
)
md_header_splits = markdown_splitter.split_text(markdown_document)

md_header_splits[0]
```

Document(page_content='Hi this is Jim \nHi this is Joe', metadata={'Header 1': 'Title', 'Header 2': 'Chapter 1'})

```
md_header_splits[1]
```

Document(page_content='Hi this is Lance', metadata={'Header 1': 'Title', 'Header 2': 'Chapter 1', 'Header 3': 'Section'})



LangChain

Chat with Your Data

Vector Stores and Embeddings

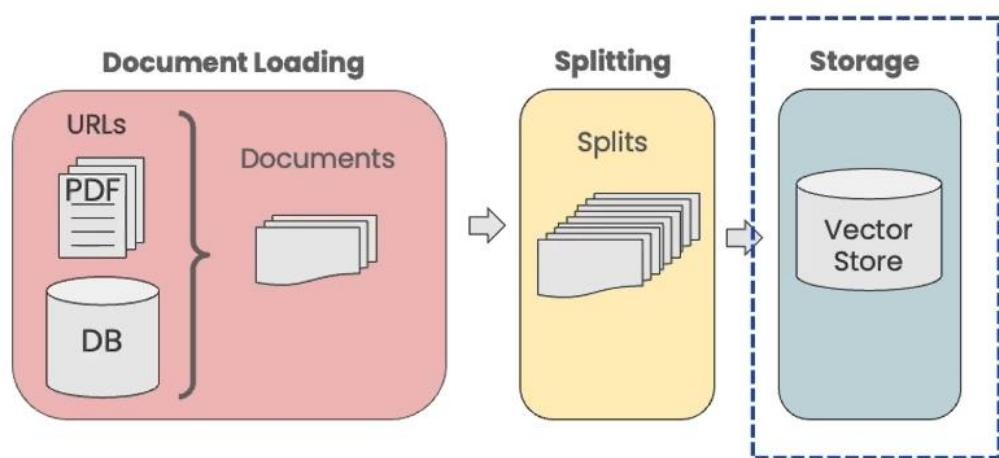


LangChain

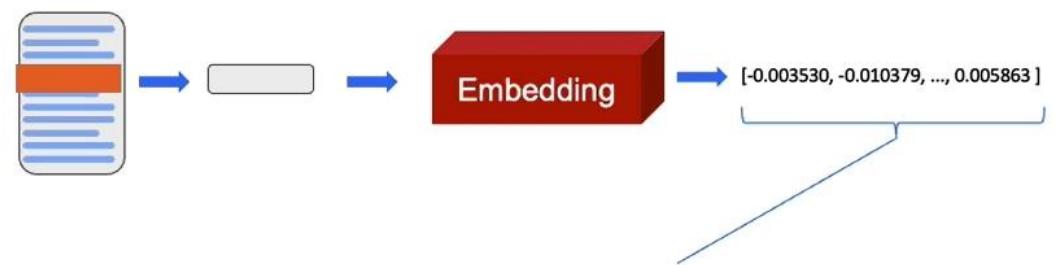


DeepLearning.AI

Vector Stores

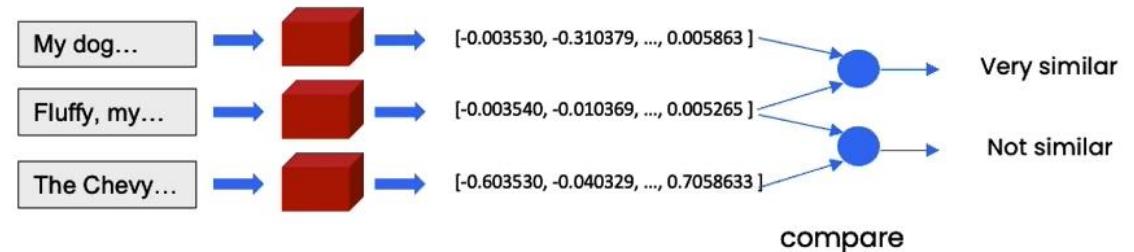


Embeddings

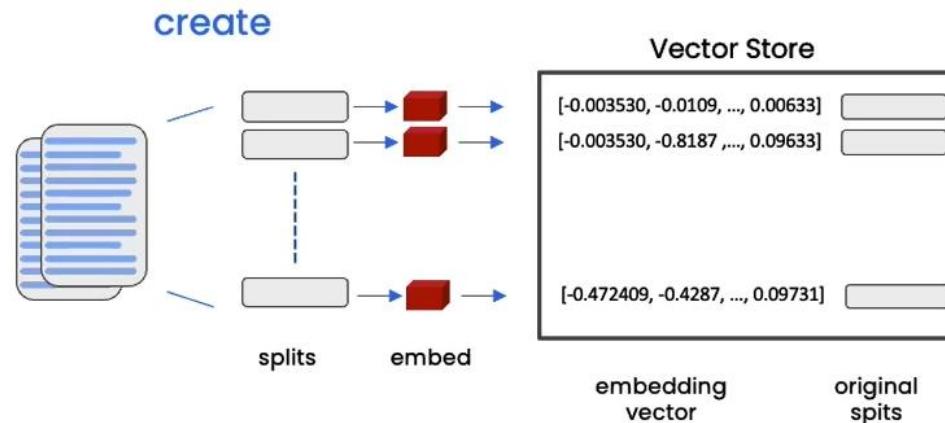


- Embedding vector captures content/meaning
- Text with similar content will have similar vectors

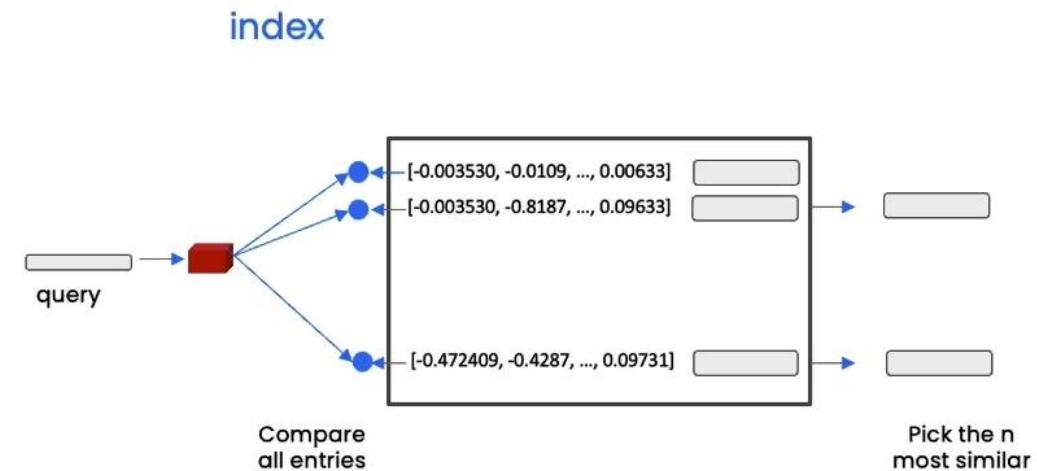
- 1) My dog Rover likes to chase squirrels.
2) Fluffy, my cat, refuses to eat from a can.
3) The Chevy Bolt accelerates to 60 mph in 6.7 seconds.



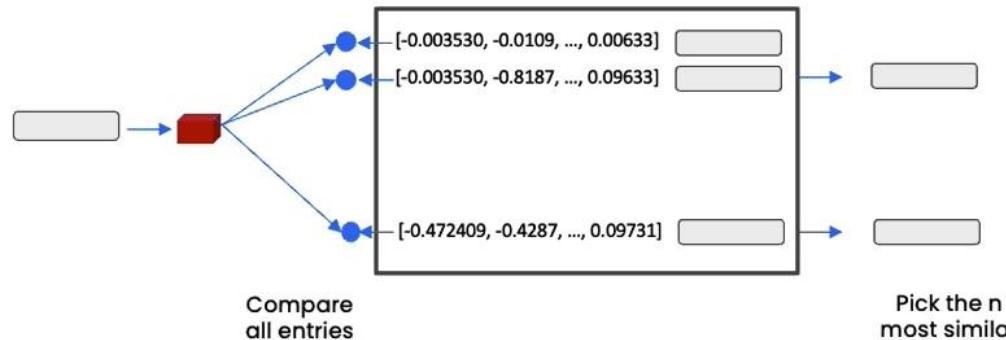
Vector Store



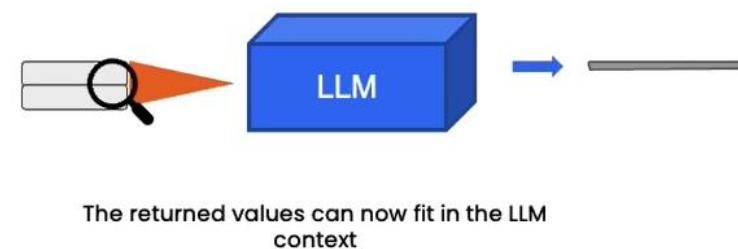
Vector Store/Database



index



Process with llm



```
from langchain.embeddings.openai import OpenAIEmbeddings
embedding = OpenAIEmbeddings()
```

```
sentence1 = "i like dogs"
sentence2 = "i like canines"
sentence3 = "the weather is ugly outside"
```

```
embedding1 = embedding.embed_query(sentence1)
embedding2 = embedding.embed_query(sentence2)
embedding3 = embedding.embed_query(sentence3)
```

```
import numpy as np
```

```
np.dot(embedding1, embedding2)
```

```
0.9631853877103519
```

```
np.dot(embedding1, embedding3)
```

```
0.770999765129468
```

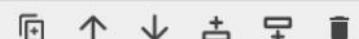
```
np.dot(embedding2, embedding3)
```

```
0.7596334120325541
```

```
from langchain.vectorstores import Chroma
persist_directory = 'docs/chroma/'
! rm -rf ./docs/chroma
vectordb = Chroma.from_documents(
    documents=splits,
    embedding=embedding,
    persist_directory=persist_directory
)
```

```
print(vectordb._collection.count())
```

```
209
```



```
question = "is there an email i can ask for help"
```

```
vectordb.persist()
```

```
docs = vectordb.similarity_search(question, k=3)
```

↑ ↓ ← → ⌂ ⌂

```
len(docs)
```

```
3
```

```
print(docs[0].page_content)
```

cs229-qa@cs.stanford.edu. This goes to an account that's read by all the TAs and me. So

rather than sending us email individually, if you send email to this account, it will

actually let us get back to you maximally quickly with answers to your questions.

If you're asking questions about homework problems, please say in the subject line which assignment and which question the email refers to, since that will also help us to route your question to the appropriate TA or to me appropriately and get the response back to you quickly.

Let's see. Skipping ahead – let's see – for homework, one midterm, one open and term project. Notice on the honor code. So one thing that I think will help you to succeed and do well in this class and even help you to enjoy this class more is if you form a study group.

So start looking around where you're sitting now or at the end of class today, mingle a little bit and get to know your classmates. I strongly encourage you to form study groups and sort of have a group of people to study with and have a group of your fellow students to talk over these concepts with. You can also post on the class news group if you want to use that to try to form a study group.

But some of the problem sets in this class are reasonably difficult. People

```
question = "what did they say about matlab?"
```

```
docs = vectordb.similarity_search(question,k=5)
```

```
docs[0]
```

```
Document(page_content='those homeworks will be done in either MATLAB or in Octave, which is sort of - I \nknow some people call it a free version of MATLAB, which it sort of is, sort of isn\'t. \nSo I guess for those of you that haven\'t seen MATLAB before, and I know most of you \nhave, MATLAB is I guess part of the programming language that makes it very easy to write codes using matrices, to write code for numerical routines, to move data around, to \nplot data. And it\'s sort of an extremely easy to learn tool to use for implementing a lot of \nlearning algorithms. \nAnd in case some of you want to work on your own home computer or something if you \ndon\'t have a MATLAB license, for the purposes of this class, there\'s also - [inaudible] \nwrite that down [inaudible] MATLAB - there\'s also a software package called Octave \nthat you can download for free off the Internet. And it has somewhat fewer features than MATLAB, but it\'s free, and for the purposes of this class, it will work for just about \neverything. \nSo actually I, well, so yeah, just a side comment for those of you that haven\'t seen \nMATLAB before I guess, once a colleague of mine at a different university, not at \nStanford, actually teaches another machine learning course. He\'s taught it for many years. \nSo one day, he was in his office, and an old student of his from, like, ten years ago came \ninto his office and he said, "Oh, professor, professor, thank you so much for your", metadata={'source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 8})
```

```
docs[1]
```

```
Document(page_content='those homeworks will be done in either MATLAB or in Octave, which is sort of - I \nknow some people call it a free version of MATLAB, which it sort of is, sort of isn\'t. \nSo I guess for those of you that haven\'t seen MATLAB before, and I know most of you \nhave, MATLAB is I guess part of the programming language that makes it very easy to write codes using matrices, to write code for numerical routines, to move data around, to \nplot data. And it\'s sort of an extremely easy to learn tool to use for implementing a lot of \nlearning algorithms. \nAnd in case some of you want to work on your own home computer or something if you \ndon\'t have a MATLAB license, for the purposes of this class, there\'s also - [inaudible] \nwrite that down [inaudible] MATLAB - there\'s also a software package called Octave \nthat you can download for free off the Internet. And it has somewhat fewer features than MATLAB, but it\'s free, and for the purposes of this class, it will work for just about \neverything. \nSo actually I, well, so yeah, just a side comment for those of you that haven\'t seen \nMATLAB before I guess, once a colleague of mine at a different university, not at \nStanford, actually teaches another machine learning course. He\'s taught it for many years. \nSo one day, he was in his office, and an old student of his from, like, ten years ago came \ninto his office and he said, "Oh, professor, professor, thank you so much for your", metadata={'source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 8})
```

Two same answers because 2 same documents were loaded, and 2 chunks were returned.

```
question = "what did they say about regression in the third lecture?"
```

```
docs = vectordb.similarity_search(question,k=5)
```

```
for doc in docs:  
    print(doc.metadata)
```

```
['source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 0}  
['source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 14}  
['source': 'docs/cs229_lectures/MachineLearning-Lecture02.pdf', 'page': 0}  
['source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 6}  
['source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 8}
```



```
docs = vectordb.similarity_search(question,k=5)
```

```
for doc in docs:  
    print(doc.metadata)
```

```
{'source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 0}  
{'source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 14}  
{'source': 'docs/cs229_lectures/MachineLearning-Lecture02.pdf', 'page': 0}  
{'source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 6}  
{'source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 8}
```

```
print(docs[4].page_content)
```

into his office and he said, "Oh, professor, professor, thank you so much for your machine learning class. I learned so much from it. There's this stuff that I learned in your class, and I now use every day. And it's helped me make lots of money, and here's a picture of my big house."

So my friend was very excited. He said, "Wow. That's great. I'm glad to hear this

machine learning stuff was actually useful. So what was it that you learned? Was it

logistic regression? Was it the PCA? Was it the data networks? What was it that you

learned that was so helpful?" And the student said, "Oh, it was the MATLAB." So for those of you that don't know MATLAB yet, I hope you do learn it. It's not hard,

and we'll actually have a short MATLAB tutorial in one of the discussion sections for

those of you that don't know it.

Okay. The very last piece of logistical thing is the discussion sections. So discussion

sections will be taught by the TAs, and attendance at discussion sections is optional, although they'll also be recorded and televised. And we'll use the discussion sections

mainly for two things. For the next two or three weeks, we'll use the discussion sections

to go over the prerequisites to this class or if some of you haven't seen probability or

LangChain

Chat with your data

Retrieval



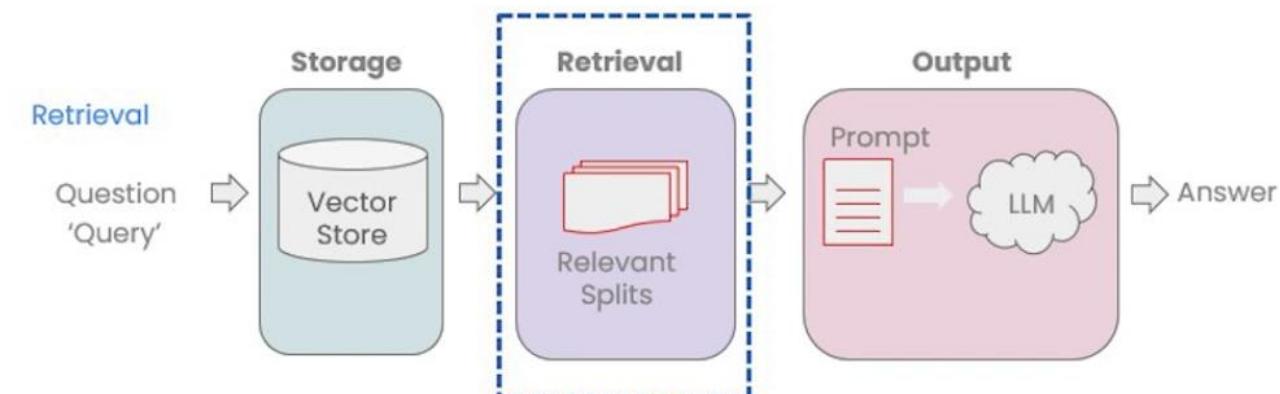
LangChain



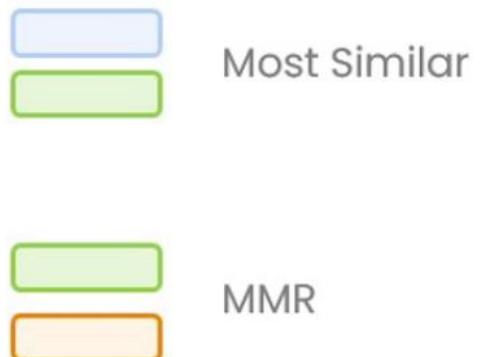
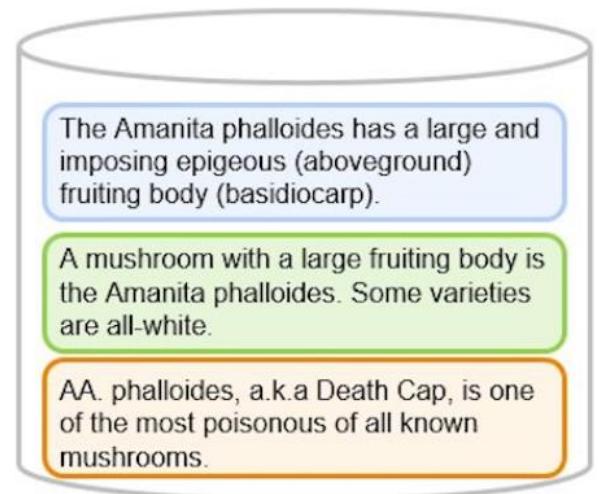
DeepLearning.AI

Retrieval

Maximum marginal relevance(MMR)

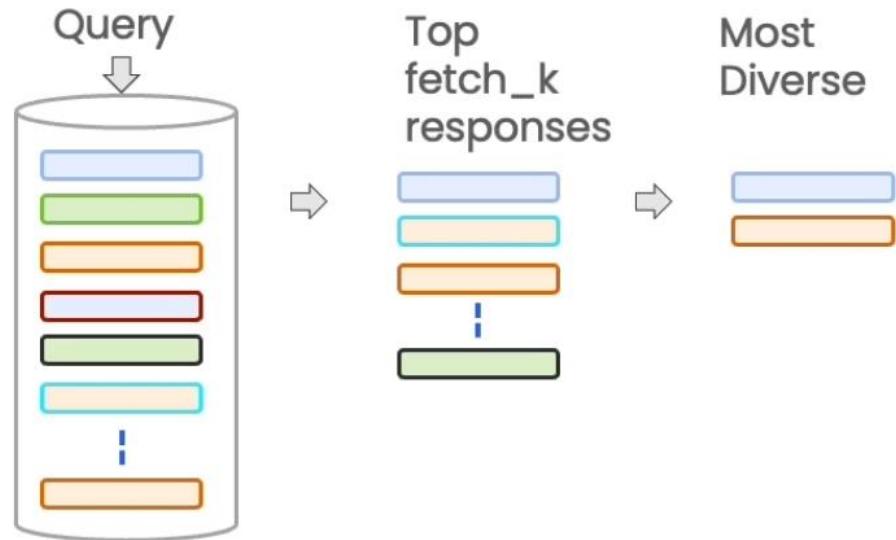


- You may not always want to choose the most similar responses



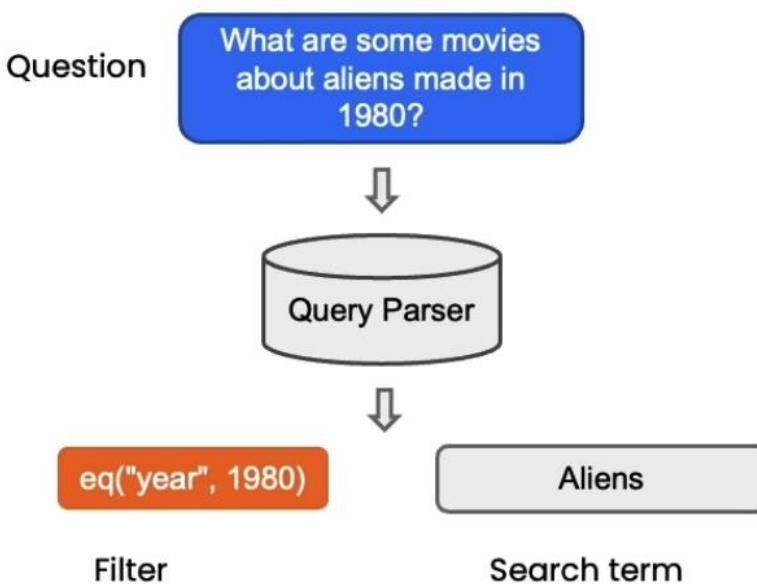
MMR algorithm

- Query the Vector Store
- Choose the `fetch_k` most similar responses
- Within those responses choose the `k` most diverse



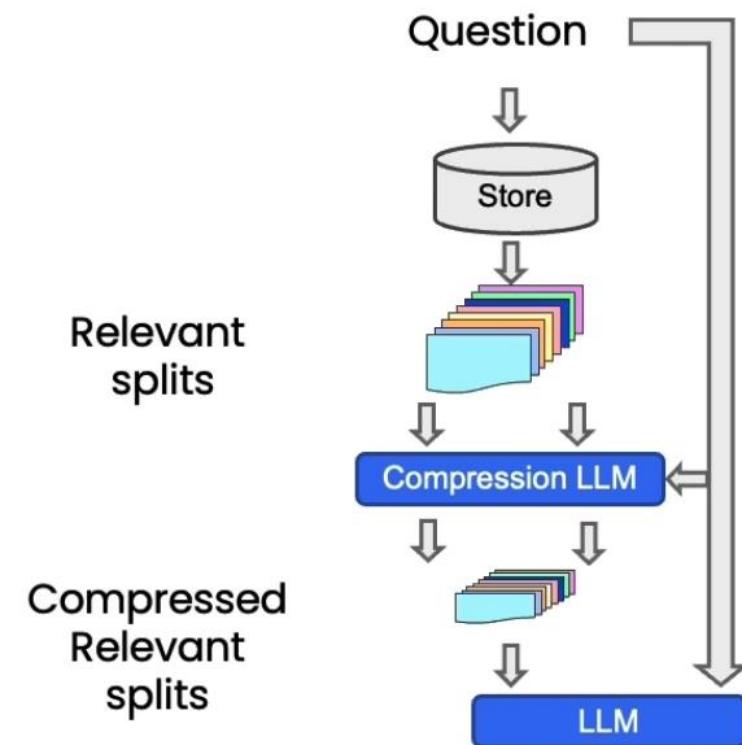
LLM Aided Retrieval

- There are several situations where the **Query** applied to the DB is more than just the **Question** asked.
- One is SelfQuery, where we use an LLM to convert the user question into a query



Compression

- Increase the number of results you can put in the context by shrinking the responses to only the relevant information.



```
import os
import openai
import sys
sys.path.append('..../..')

from dotenv import load_dotenv, find_dotenv
_ = load_dotenv(find_dotenv()) # read local .env file

openai.api_key = os.environ['OPENAI_API_KEY']

from langchain.vectorstores import Chroma
from langchain.embeddings.openai import OpenAIEMBEDDINGS
persist_directory = 'docs/chroma/'

embedding = OpenAIEMBEDDINGS()
vectordb = Chroma(
    persist_directory=persist_directory,
    embedding_function=embedding
)

print(vectordb._collection.count())
209
```

```
texts = [
    """The Amanita phalloides has a large and imposing epigeous (aboveground)
    """A mushroom with a large fruiting body is the Amanita phalloides. Some varieties are all-white.
    """A. phalloides, a.k.a Death Cap, is one of the most poisonous of all known mushrooms."""
]

smalldb = Chroma.from_texts(texts, embedding=embedding)

question = "Tell me about all-white mushrooms with large fruiting bodies"

smalldb.similarity_search(question, k=2)
[Document(page_content='A mushroom with a large fruiting body is the Amanita phalloides. Some varieties are all-white.', metadata={}),
 Document(page_content='The Amanita phalloides has a large and imposing epigeous (aboveground) fruiting body (basidiocarp).', metadata={})]

smalldb.max_marginal_relevance_search(question,k=2, fetch_k=3)
[Document(page_content='A mushroom with a large fruiting body is the Amanita phalloides. Some varieties are all-white.', metadata={}),
 Document(page_content='A. phalloides, a.k.a Death Cap, is one of the most poisonous of all known mushrooms.', metadata={})]

question = "what did they say about matlab?"
docs_ss = vectordb.similarity_search(question,k=3)
```

```
sonous of all known mushrooms.', metadata={})]
```

```
question = "what did they say about matlab?"  
docs_ss = vectordb.similarity_search(question,k=3)
```

```
docs_ss[0].page_content[:100]
```

```
'those homeworks will be done in either MATLAB or in Octave, which is sort of  
- I \nknow some people '
```

```
docs_ss[1].page_content[:100]
```

```
'those homeworks will be done in either MATLAB or in Octave, which is sort of  
- I \nknow some people '
```

```
docs_mmr = vectordb.max_marginal_relevance_search(question,k=3)
```

```
docs_mmr[0].page_content[:100]
```

```
'those homeworks will be done in either MATLAB or in Octave, which is sort of  
- I \nknow some people '
```

```
docs_mmr[1].page_content[:100]
```

```
'algorithm then? So what's different? How come I was making all that noise ea  
rlier about \nleast squa'
```

```
question = "what did they say about regression in the third lecture?"
```

```
docs = vectordb.similarity_search(  
    question,  
    k=3,  
    filter={"source":"docs/cs229_lectures/MachineLearning-Lecture03.pdf"})
```

```
for d in docs:  
    print(d.metadata)
```

```
{'source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 0}  
'source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 14}  
'source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 4}
```

```
from langchain.llms import OpenAI  
from langchain.retrievers.self_query.base import SelfQueryRetriever  
from langchain.chains.query_constructor.base import AttributeInfo
```

```
metadata_field_info = [  
    AttributeInfo(  
        name="source",  
        description="The lecture the chunk is from, should be one of `docs/cs229_lectures/MachineLearning-Lecture03.pdf`",  
        type="string",  
    ),  
    AttributeInfo(  
        name="page",  
        description="The page from the lecture",  
        type="integer",  
    ),  
]
```

Manually defined source

LLM auto defined source

Metadata info for source and page is defined

```
metadata_field_info = [
    AttributeInfo(
        name="source",
        description="The lecture the chunk is from, should be one of `docs/cs229_lectures/MachineLearning-Lecture03.pdf`",
        type="string",
    ),
    AttributeInfo(
        name="page",
        description="The page from the lecture",
        type="integer",
    ),
]

```

```
document_content_description = "Lecture notes"
llm = OpenAI(temperature=0)
retriever = SelfQueryRetriever.from_llm(
    llm,
    vectordb,
    document_content_description,
    metadata_field_info,
    verbose=True
)
```

Self query retrieval defines the source and page that we want to search and return, it will be done by LLM

```
docs = retriever.get_relevant_documents(question)
query='regression' filter=Comparison(comparator=<Comparator.EQ>, attribute='source', value='docs/cs229_lectures/MachineLearning-Lecture03.pdf') limit=None

for d in docs:
    print(d.metadata)

{'source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 14}
{'source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 0}
{'source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 10}
{'source': 'docs/cs229_lectures/MachineLearning-Lecture03.pdf', 'page': 10}
```

Self query retrieval shows from which documents the answer is produced

```
from langchain.retrievers import ContextualCompressionRetriever
from langchain.retrievers.document_compressors import LLMChainExtractor
```



```
def pretty_print_docs(docs):
    print(f"\n{'-' * 100}\n".join([f"Document {i+1}:\n\n" + d.page_content for
```



```
llm = OpenAI(temperature=0)
compressor = LLMChainExtractor.from_llm(llm)
```

```
compression_retriever = ContextualCompressionRetriever(
    base_compressor=compressor,
    base_retriever=vectordb.as_retriever()
)
```

```
question = "what did they say about matlab?"
compressed_docs = compression_retriever.get_relevant_documents(question)
pretty_print_docs(compressed_docs)
```

Document 1:

"MATLAB is I guess part of the programming language that makes it very easy to write codes using matrices, to write code for numerical routines, to move data around, to plot data. And it's sort of an extremely easy to learn tool to use for implementing a lot of learning algorithms."

Document 2:

"MATLAB is I guess part of the programming language that makes it very easy to write codes using matrices, to write code for numerical routines, to move data around, to plot data. And it's sort of an extremely easy to learn tool to use for implementing a lot of learning algorithms."

Document 3:

"And the student said, "Oh, it was the MATLAB." So for those of you that don't

ContextualCompressionRetriever
compress answers from all documents.
LLMChainExtractor will return only
relevant parts of document

We got shorter documents, but it can
be repeated

```
from langchain.retrievers import ContextualCompressionRetriever
from langchain.retrievers.document_compressors import LLMChainExtractor

def pretty_print_docs(docs):
    print(f"\n{'-' * 100}\n".join([f"Document {i+1}:\n\n" + d.page_content for

llm = OpenAI(temperature=0)
compressor = LLMChainExtractor.from_llm(llm)

compression_retriever = ContextualCompressionRetriever(
    base_compressor=compressor,
    base_retriever=vectordb.as_retriever(search_type="mmr")
)

question = "what did they say about matlab?"
compressed_docs = compression_retriever.get_relevant_documents(question)
pretty_print_docs(compressed_docs)
```

Document 1:

"MATLAB is I guess part of the programming language that makes it very easy to write codes using matrices, to write code for numerical routines, to move data around, to plot data. And it's sort of an extremely easy to learn tool to use for implementing a lot of learning algorithms."

Document 2:

"And the student said, "Oh, it was the MATLAB." So for those of you that don't know MATLAB yet, I hope you do learn it. It's not hard, and we'll actually have a short MATLAB tutorial in one of the discussion sections for those of you that don't know it."

With MMR we will not get duplicate documents



Other types of retrieval

Not using a vector database, such as:

- SVM Support Vector Machine retrieval,
- TF-IDF term frequency-inverse document frequency, is a measure of importance of a word to a document in a collection or corpus, adjusted for the fact that some words appear more frequently in general.
- ...

```
from langchain.retrievers import SVMRetriever
from langchain.retrievers import TFIDFRetriever
from langchain.document_loaders import PyPDFLoader
from langchain.text_splitter import RecursiveCharacterTextSplitter
```

```
# Load PDF
loader = PyPDFLoader("docs/cs229_lectures/MachineLearning-Lecture01.pdf")
pages = loader.load()
all_page_text=[p.page_content for p in pages]
joined_page_text=".join(all_page_text)

# Split
text_splitter = RecursiveCharacterTextSplitter(chunk_size = 1500,chunk_overlap
splits = text_splitter.split_text(joined_page_text)

svm_retriever = SVMRetriever.from_texts(splits,embedding)
tfidf_retriever = TFIDFRetriever.from_texts(splits)
```



```
question = "what did they say about matlab?"
```

```
docs_svm=svm_retriever.get_relevant_documents(question)
```

```
docs_svm[0]
```

Document(page_content='don\'t have a MATLAB license, for the purposes of this class, there's also - [inaudible] \nwrite that down [inaudible] MATLAB - there's also a software package called Octave\nthat you can download for free off the Internet. And it has somewhat fewer features than MATLAB, but it's free, and for the purposes of this class, it will work for just about \nanything. \nSo actually I, well, so yeah, just a side comment for those of you that haven\'t seen \nMATLAB before I guess, once a colleague of mine at a different university, not at \nStanford, actually teaches another machine learning course. He's taught it for many years. \nSo one day, he was in his office, and an old student of his from, like, ten years ago came \ninto his office and he said, "Oh, professor, professor, thank you so much for your \nmachine learning class. I learned so much from it. There's this stuff that I learned in your \nclass, and I now use every day. And it's helped me make lots of money, and here's a \npicture of my big house." \nSo my friend was very excited. He said, "Wow. That's great. I'm glad to hear this \nmachine learning stuff was actually useful. So what was it that you learned? Was it \nlogistic regression? Was it the PCA? Was it the data networks? What was it that you \nlearned that was so helpful?" And the student said, "Oh, it was the MATLAB." \nSo for those of you that don't know MATLAB yet, I hope you do learn it. It's not hard.', metadata={})



se. He's taught it for many years. \nSo one day, he was in his office, and an old student of his from, like, ten years ago came \ninto his office and he said, "Oh, professor, professor, thank you so much for your \nmachine learning class. I learned so much from it. There's this stuff that I learned in your \nclass, and I now use every day. And it's helped me make lots of money, and here's a \npicture of my big house." \nSo my friend was very excited. He said, "Wow. That's great. I'm glad to hear this \nmachine learning stuff was actually useful. So what was it that you learned? Was it \nlogistic regression? Was it the PCA? Was it the data networks? What was it that you \nlearned that was so helpful?" And the student said, "Oh, it was the MATLAB." \nSo for those of you that don't know MATLAB yet, I hope you do learn it. It's not hard.', metadata={})

```
question = "what did they say about matlab?"
```

```
docs_tfidf=tfidf_retriever.get_relevant_documents(question)
```

```
docs_tfidf[0]
```

Document(page_content="Saxena and Min Sun here did, which is given an image like this, right? This is actually a \npicture taken of the Stanford campus. You can apply that sort of clustering algorithm and \ngroup the picture into regions. Let me actually blow that up so that you can see it more \nclarly. Okay. So in the middle, you see the lines sort of grouping the image together, \ngrouping the image into [inaudible] regions. \nAnd what Ashutosh and Min did was they then applied the learning algorithm to say can \nwe take this clustering and use it to build a 3D model of the world? And so using the \nclustering, they then had a learning algorithm try to learn what the 3D structure of the \nworld looks like so that they could come up with a 3D model that you can sort of fly \nthrough, okay? Although many people used to think it's not possible to take a single \nimage and build a 3D model, but using a learning algorithm and that sort of clustering\nalgorithm is the first step. They were able to. \nI'll just show you one more example. I like this because it's a picture of Stanford with our \nbeautiful Stanford campus. So again, taking the same sort of clustering algorithms, taking \nthe same sort of unsupervised learning algorithm, you can group the pixels into different \nregions. And using that as a pre-processing step, they eventually built this sort of 3D model of Stanford campus in a single picture. You can sort of walk into the ceiling, look", metadata={})



LangChain

Chat with Your Data

Question Answering ChatBot

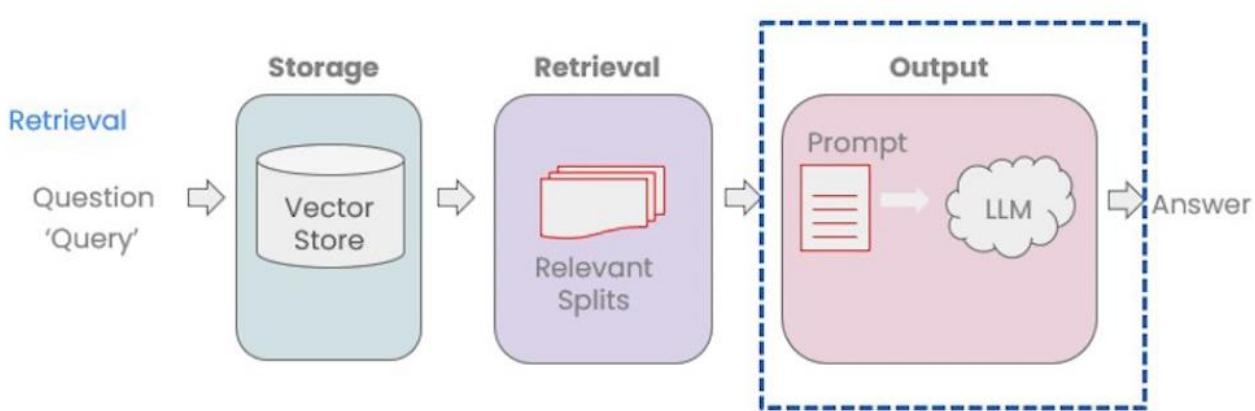


LangChain



DeepLearning.AI

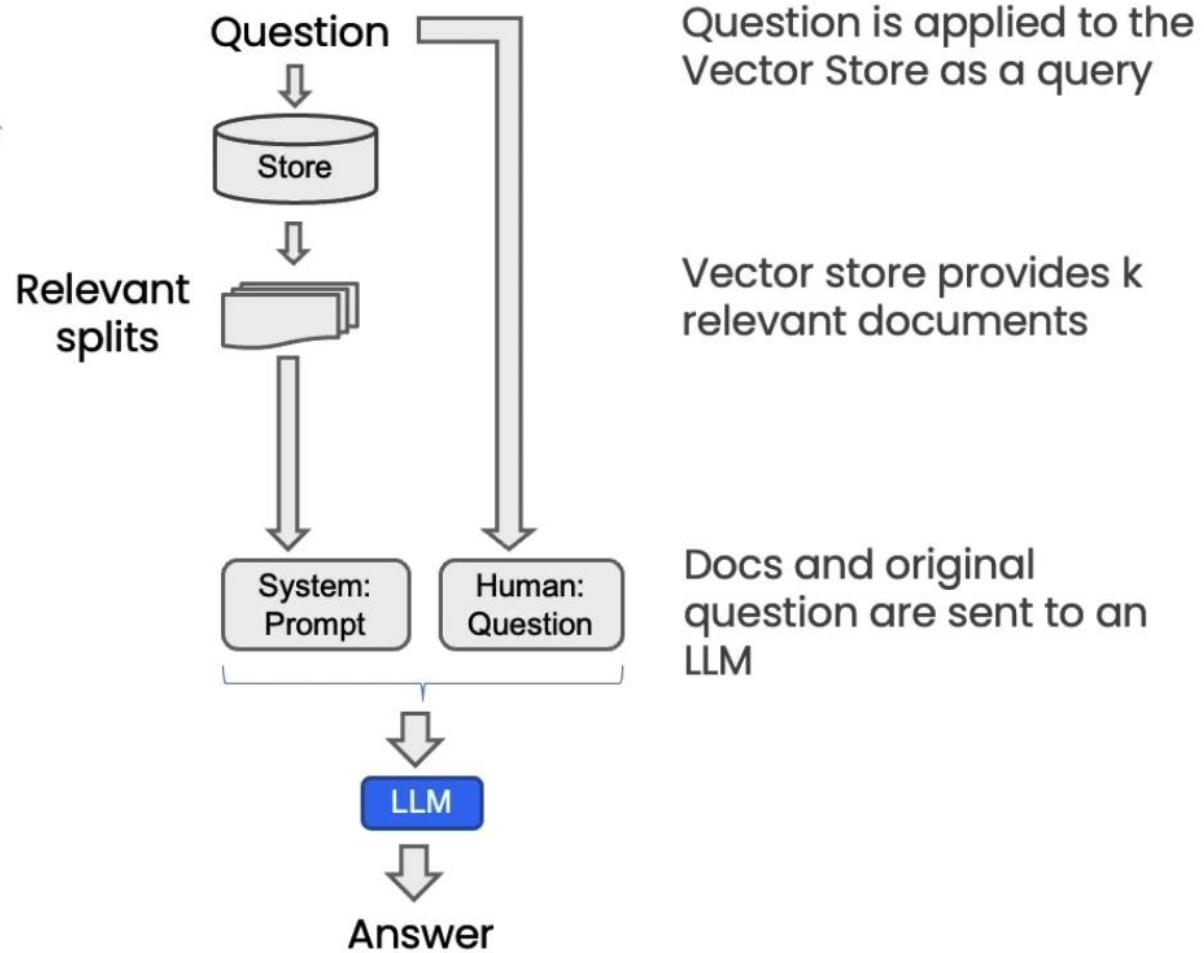
Question Answering



- Multiple relevant documents have been retrieved from the vector store
- Potentially compress the relevant splits to fit into the LLM context
- Send the information along with our question to an LLM to select and format an answer

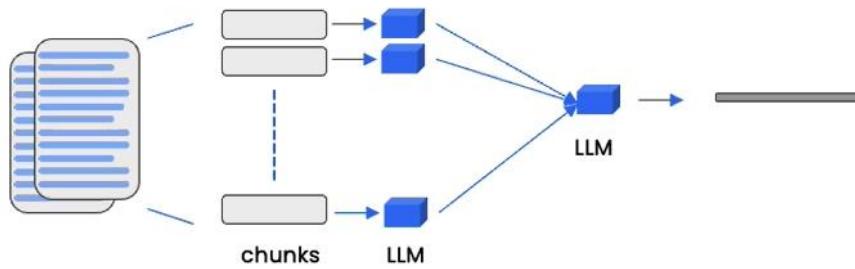
RetrievalQA chain

`RetrievalQA.from_chain_type(, chain_type="stuff",...)`

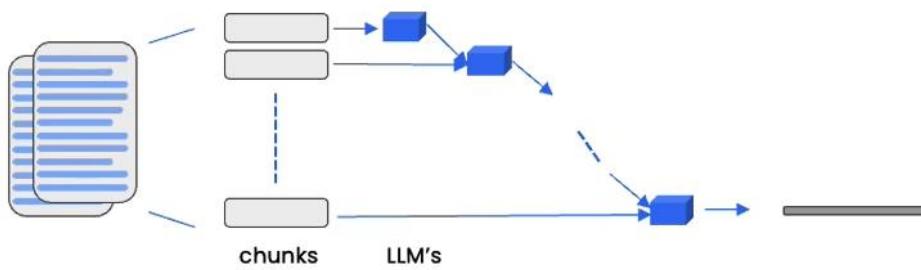


3 additional methods

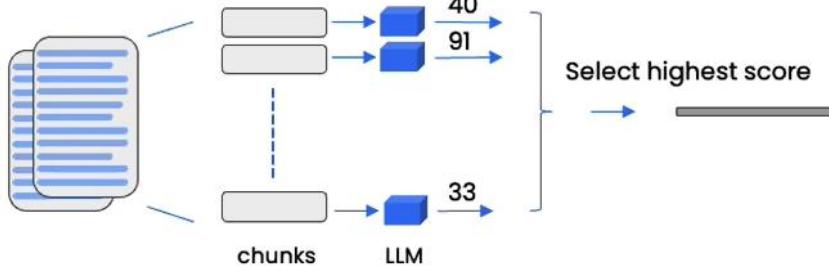
1. Map_reduce



2. Refine



3. Map_rerank



Multiple methods to use to tackle issue of short context window

```
import openai
import sys
sys.path.append('..../..')

from dotenv import load_dotenv, find_dotenv
_ = load_dotenv(find_dotenv()) # read local .env file

openai.api_key = os.environ['OPENAI_API_KEY']

from langchain.vectorstores import Chroma
from langchain.embeddings.openai import OpenAIEmbeddings
persist_directory = 'docs/chroma/'
embedding = OpenAIEmbeddings()
vectordb = Chroma(persist_directory=persist_directory, embedding_function=embedding)

print(vectordb._collection.count())
209

question = "What are major topics for this class?"
docs = vectordb.similarity_search(question,k=3)
len(docs)
3

from langchain.chat_models import ChatOpenAI
llm = ChatOpenAI(model_name="gpt-3.5-turbo", temperature=0)

from langchain.chains import RetrievalQA

qa_chain = RetrievalQA.from_chain_type(
    llm,
    retriever=vectordb.as_retriever()
)

qa_chain = RetrievalQA.from_chain_type(
    llm,
    retriever=vectordb.as_retriever()
)

result = qa_chain({"query": question})
result["result"]

'The major topic for this class is machine learning. Additionally, the class may cover statistics and algebra as refreshers in the discussion sections. Later in the quarter, the discussion sections will also cover extensions for the material taught in the main lectures.'
```

Temperature 0 for lowest variability.

Import the RetrievalQA chain to send result chunks to LLM for answer



aterial taught in the main lectures.

```
from langchain.prompts import PromptTemplate
```

```
# Build prompt
```

```
template = """Use the following pieces of context to answer the question at the end.  
{context}  
Question: {question}  
Helpful Answer: """  
QA_CHAIN_PROMPT = PromptTemplate.from_template(template)
```

```
qa_chain = RetrievalQA.from_chain_type(  
    llm,  
    retriever=vector_db.as_retriever(),  
    return_source_documents=True,  
    chain_type_kwargs={"prompt": QA_CHAIN_PROMPT})
```

```
question = "Is probability a class topic?"
```

```
result = qa_chain({"query": question})
```

```
result["result"]
```

'Yes, probability is assumed to be a prerequisite for this class. The instructor assumes familiarity with basic probability and statistics, and will go over some of the prerequisites in the discussion sections as a refresher course. Thanks for asking!'

```
result["source_documents"][0]
```



```
chain_type_kwargs={"prompt": QA_CHAIN_PROMPT})
```

```
question = "Is probability a class topic?"
```

```
result = qa_chain({"query": question})
```

```
result["result"]
```

'Yes, probability is assumed to be a prerequisite for this class. The instructor assumes familiarity with basic probability and statistics, and will go over some of the prerequisites in the discussion sections as a refresher course. Thanks for asking!'

```
result["source_documents"][0]
```

Document(page_content="of this class will not be very programming intensive, although we will do some \nprogramming, mostly in either MATLAB or Octave. I'll say a bit more about that later. \nI also assume familiarity with basic probability and statistics. So most undergraduate \nstatistics class, like Stat 116 taught here at Stanford, will be more than enough. I'm gonna \nassume all of you know what random variables are, that all of you know what expectation \nis, what a variance or a random variable is. And in case of some of you, it's been a while \nsince you've seen some of this material. At some of the discussion sections, we'll actually \ngo over some of the prerequisites, sort of as a refresher course under prerequisite class. \nI'll say a bit more about that later as well. \nLastly, I also assume familiarity with basic linear algebra. And again, most undergraduate \nlinear algebra courses are more than enough. So if you've taken courses like Math 51, \n103, Math 113 or CS205 at Stanford, that would be more than enough. Basically, I'm \ngonna assume that all of you know what matrices and vectors are, that you know how to \nmultiply matrices and vectors and multiply matrix and matrices, that you know what a matrix inverse is. If you know what an eigenvector of a matrix is, that'd be even better. \nBut if you don't quite know or if you're not quite sure, that's fine, too. We'll go over it in \nthe review sections.", metadata={'source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 4})



```
qa_chain_mr = RetrievalQA.from_chain_type(  
    llm,  
    retriever=vectordb.as_retriever(),  
    chain_type="map_reduce"  
)
```

```
result = qa_chain_mr({"query": question})
```

```
result["result"]
```

'There is no clear answer to this question based on the given portion of the document. The document mentions familiarity with basic probability and statistics as a prerequisite for the class, and there is a brief mention of probability in the text, but it is not clear if it is a main topic of the class. The instructor mentions using a probabilistic interpretation to derive a learning algorithm, but does not go into further detail about probability as a topic.'



- It can take multiple documents
- It takes a lot slower
- Response is even worse, it's responding based on different documents, context is distributed over multiple documents

Run: RetrievalQA

[★ Rate Run](#) [+ Add to Dataset](#) [Share](#)
✓ CHAIN a minute ago

Chain Input

[Jump to bottom](#) [Copy](#)

query: Is probability a class topic?

Chain Output

[Jump to bottom](#) [Copy](#)

result: There is no clear answer to this question based on the given portion of the document. The document mentions familiarity with basic probability and statistics as a prerequisite for the class, and there is a brief mention of probability in the text, but it is not clear if it is a main topic of the class. The instructor mentions using a probabilistic interpretation to derive a learning algorithm, but does not go into further detail about probability as a topic.

Run Metadata

[Copy](#)

completion_tokens: 87

prompt_tokens: 1503

total_tokens: 1590

RUNTIME

library: "langchain"

library_version: "0.0.213"

platform: "macOS-13.4-arm64-arm-64bit"

runtime: "python"

Run Stats

RBC 1,503 → 87 1,590 tokens

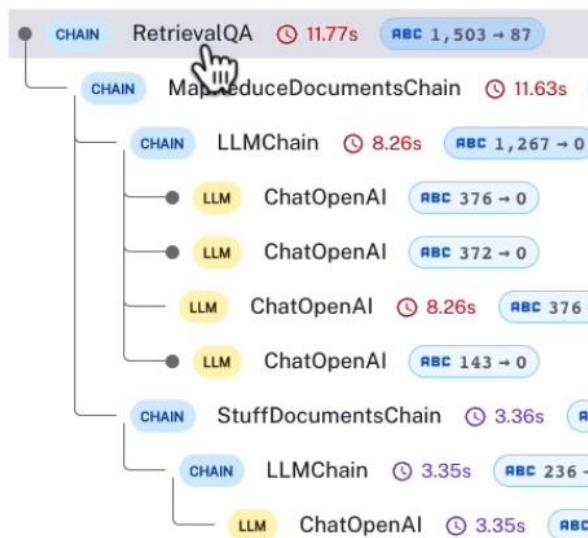
⌚ Start: 06/25/2023, 10:23:49 PM

⌚ End: 06/25/2023, 10:24:01 PM

⌚ Response time: 11.77s

- Check with LangSmith what is happening with map reduce

Child Runs



```
import os
os.environ["LANGCHAIN_TRACING_V2"] = "true"
os.environ["LANGCHAIN_ENDPOINT"] = "https://api.langchain.plus"
os.environ["LANGCHAIN_API_KEY"] = LCP_API_KEY
```

```
qa_chain_mr = RetrievalQA.from_chain_type(
    llm,
    retriever=vectordb.as_retriever(),
    chain_type="map_reduce"
)
result = qa_chain_mr({"query": question})
result["result"]
```

'There is no clear answer to this question based on the given portion of the document. The document mentions familiarity with basic probability and statistics as a prerequisite for the class, and there is a brief mention of probability in the text, but it is not clear if it is a main topic of the class. The instructor mentions using a probabilistic interpretation to derive a learning algorithm, but does not go into further detail about probability as a topic.'

```
qa_chain_mr = RetrievalQA.from_chain_type(
    llm,
    retriever=vectordb.as_retriever(),
    chain_type="refine"
)
result = qa_chain_mr({"query": question})
result["result"]
```

"The class assumes familiarity with basic probability and statistics, but will have review sections to refresh the prerequisites. The instructor also mentions that they will use the discussion sections to go over extensions for the material that they are teaching in the main lectures. These extensions will cover topics that they didn't have time to cover in the main lectures, as machine learning is a vast field."

Refine retrieval method gives better results than map reduce because it combines information sequentially and takes over more information than map reduce



```
a_chain = RetrievalQA.from_chain_type(  
    llm,  
    retriever=vectordb.as_retriever()
```

```
question = "Is probability a class topic?"  
result = qa_chain({"query": question})  
result["result"]
```

Yes, probability is a topic that will be assumed to be familiar to students in this class. The instructor assumes that students have a basic understanding of probability and statistics, and will go over some of the prerequisites as a refresher course in the discussion sections.'

```
question = "why are those prerequisites needed?"  
result = qa_chain({"query": question})  
result["result"]
```

The prerequisites for the class are assumed to be basic knowledge of computer science and basic computer skills and principles. The instructor assumes that all students know about big-O notation and have a basic understanding of computer science. These prerequisites are needed to ensure that all students can follow the technical content of the class.'



Not a good answer as the CHAIN does not have a memory what was a previous question or answer

LangChain

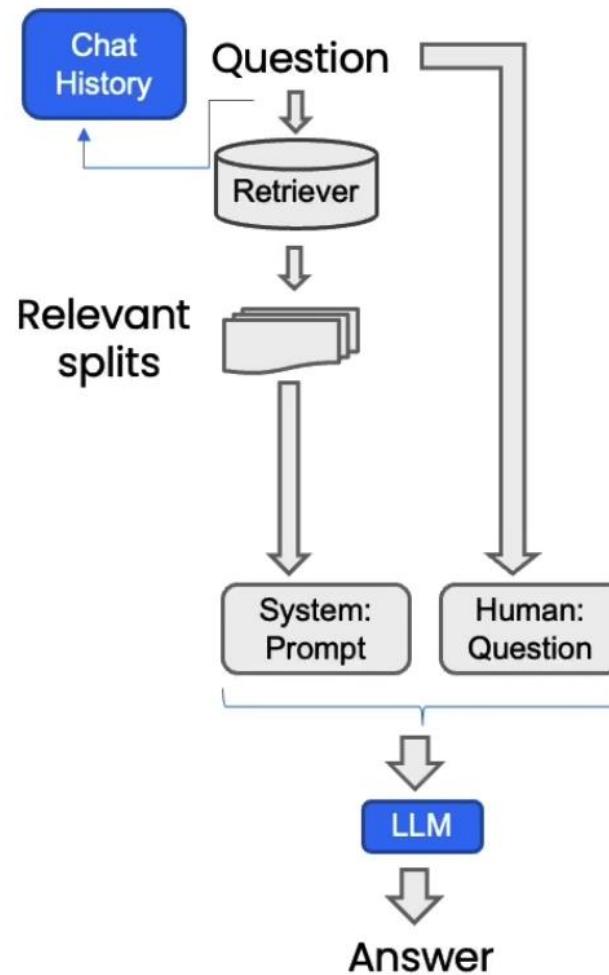
Chat with Your Data

Question Answering ChatBot



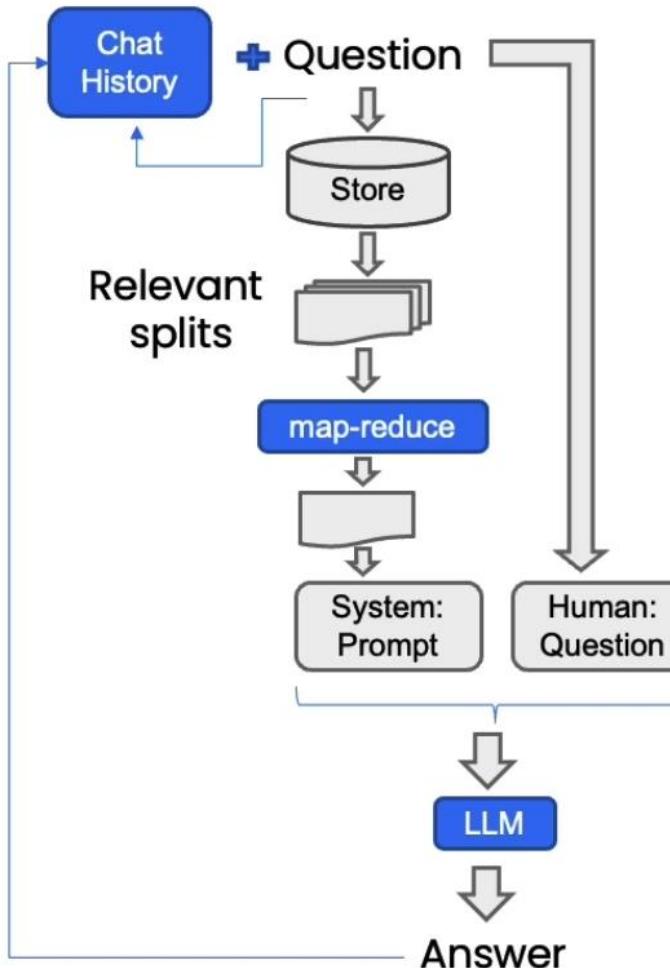
ConversationalRetrievalChain

```
qa = ConversationalRetrievalChain.from_llm(ChatOpenAI(temperature=0),  
vectorstore.as_retriever(), memory=memory)
```



Chat history will give context from previous answers/question and continue the conversation with followup.

ConversationalRetrievalChain



All types of retrieval can be used.

You can add additional Retriever and compression features as needed

```
from langchain.vectorstores import Chroma
from langchain.embeddings.openai import OpenAIEMBEDDINGS
persist_directory = 'docs/chroma/'
embedding = OpenAIEMBEDDINGS()
vectordb = Chroma(persist_directory=persist_directory, embedding_function=embedding)

question = "What are major topics for this class?"
docs = vectordb.similarity_search(question,k=3)
len(docs)

3

from langchain.chat_models import ChatOpenAI
llm = ChatOpenAI(model_name="gpt-3.5-turbo", temperature=0)
llm.predict("Hello world!")

'Hello there! How can I assist you today?'
```



```
# Build prompt
from langchain.prompts import PromptTemplate
template = """Use the following pieces of context to answer the question at the
{context}
Question: {question}
Helpful Answer:"""
QA_CHAIN_PROMPT = PromptTemplate(input_variables=["context", "question"], template=template)

# Run chain
from langchain.chains import RetrievalQA
question = "Is probability a class topic?"
qa_chain = RetrievalQA.from_chain_type(llm,
                                       retriever=vectordb.as_retriever(),
                                       return_source_documents=True,
                                       chain_type_kwargs={"prompt": QA_CHAIN_PROMPT})

result = qa_chain({"query": question})
result["result"]
```

'Yes, probability is assumed to be a prerequisite for this class. The instructor assumes familiarity with basic probability and statistics, and will go over some of the prerequisites in the discussion sections as a refresher course. Thanks for asking!'



```
from langchain.memory import ConversationBufferMemory
memory = ConversationBufferMemory(
    memory_key="chat_history",
    return_messages=True
)
```

```
from langchain.chains import ConversationalRetrievalChain
retriever=vectoradb.as_retriever()
qa = ConversationalRetrievalChain.from_llm(
    llm,
    retriever=retriever,
    memory=memory
)
```

```
question = "Is probability a class topic?"
result = qa({"question": question})
```

```
result["answer"]
```

'Yes, probability is a topic that will be assumed to be familiar to students in this class. The instructor mentions that basic probability and statistics are prerequisites for the class and assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough.'

```
question = "why are those prerequisites needed?"
result = qa({"question": question})
```

```
result["answer"]
```

'The reason for requiring basic probability and statistics as prerequisites for the class is that the class assumes familiarity with these topics. The instructor assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough. The class will cover machine learning, which is a huge field, and assumes that students already know what random variables are, what expectation is, what a variance or a random variable is.'

New chain type -> conversation retrieval chain, adds memory, takes history and the new question and condenses into stand-alone question to pass to the vector store to look up for relevant documents.

Follow up question on previous answer.

[Debug](#) > Project: default

Project: default

P50 Latency: 3.91s P99 Latency: 37.00s Total Tokens: 35,269

eg. eq(run_type, "chain")

Quick Searches



Columns Root Runs All Runs

<input type="checkbox"/> Run Type	Name	Input	Output	Start Time	Tags
<input type="checkbox"/> > CHAIN	ConversationalRe...	{"question": "Is probability and statistics as prerequisites for the class?"}	The reason for requiring basic probability and statistics as prerequisites for the class is that the class assumes familiarity with these topics. The instructor assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough.	6/25/2023, 3:55:41 ...	
<input type="checkbox"/> > CHAIN	ConversationalRe...	{"question": "Is probability and statistics as prerequisites for the class?"}	Yes, probability and statistics are prerequisites for the class and assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough.	6/25/2023, 3:55:39 ...	
<input type="checkbox"/> > CHAIN	ConversationalRe...	{"question": "Is probability and statistics as prerequisites for the class?"}	The reason for requiring basic probability and statistics as prerequisites for the class is that the class assumes familiarity with these topics. The instructor assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough.	6/25/2023, 3:55:00 ...	
<input type="checkbox"/> > CHAIN	ConversationalRe...	{"question": "Is probability and statistics as prerequisites for the class?"}	Yes, probability and statistics are prerequisites for the class and assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough.	6/25/2023, 3:54:41 ...	
<input type="checkbox"/> > CHAIN	RetrievalQA	Is probability and statistics as prerequisites for the class?	"result": "Yes, probability and statistics are prerequisites for the class and assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough."	6/25/2023, 3:52:03 ...	
<input type="checkbox"/> LLM	ChatOpenAI	human: Hello, I'm a student taking a statistics course. Can you tell me if probability and statistics are prerequisites for this class?	human: Hello, I'm a student taking a statistics course. Can you tell me if probability and statistics are prerequisites for this class?	6/25/2023, 3:51:26 ...	
<input type="checkbox"/> > CHAIN	RetrievalQA	why are those prerequisites needed?	The reason for requiring basic probability and statistics as prerequisites for the class is that the class assumes familiarity with these topics. The instructor assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough.	6/25/2023, 3:31:27 ...	
<input type="checkbox"/> > CHAIN	RetrievalQA	Is probability and statistics as prerequisites for the class?	Yes, probability and statistics are prerequisites for the class and assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough.	6/25/2023, 3:31:09 ...	
<input type="checkbox"/> > CHAIN	RetrievalQA	Is probability and statistics as prerequisites for the class?	The class assumes that students already know what random variables are, what expectation is, what a variance or a random variable is.	6/25/2023, 3:25:56 ...	
<input type="checkbox"/> > CHAIN	RetrievalQA	Is probability and statistics as prerequisites for the class?	There is no clear answer provided in the output.	6/25/2023, 3:23:49 ...	

Rows per page: 10 ▾ 1-10 of 43 < >

✓ CHAIN a few seconds ago

Runs

Setup

Chain Input

Jump to bottom

Copy

question: why are those prerequisites needed?
chat history:

- content: Is probability a class topic?
 example: false
 additional_kwargs: {}
- content: Yes, probability is a topic that will be assumed to be familiar to students in this class. The instructor mentions that basic probability and statistics are prerequisites for the class and assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough.
 example: false
 additional_kwargs: {}

Jump to bottom

Copy

answer: The reason for requiring basic probability and statistics as prerequisites for the class is that the class assumes familiarity with these topics. The instructor assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough. The class will cover machine learning, which is a huge field, and assumes that students already know what random variables are, what expectation is, what a variance or a random variable is.

Run Metadata

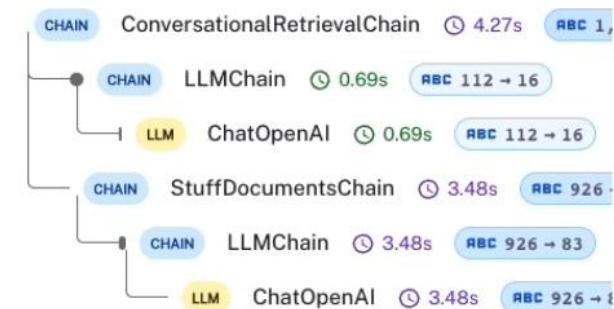
Copy

Run Stats

ABC 1,038 → 99 1,137 tokens

- ⌚ Start: 06/25/2023, 10:55:41 PM
- ⌚ End: 06/25/2023, 10:55:45 PM
- ⌚ Response time: 4.27s

Child Runs



1. Call to llm

[Home](#) [Debug](#) [Testing](#) [Monitoring](#) [Datasets](#)
[Documentation](#) [Home](#) [Debug](#) [Testing](#) [Monitoring](#) [Datasets](#)
[Documentation](#)

Run: ChatOpenAI

✓ LLM a minute ago

Chat Model Inputs

HUMAN Copy

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: Is probability a class topic?
 Assistant: Yes, probability is a topic that will be assumed to be familiar to students in this class. The instructor mentions that basic probability and statistics are prerequisites for the class and assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough.
 Follow Up Input: why are those prerequisites needed?
 Standalone question:

Chat Model Outputs

ASSISTANT Copy

What is the reason for requiring basic probability and statistics as prerequisites for the class?

Run Stats

RBC 112 → 16 128 tokens

- ⌚ Start: 06/25/2023, 10:55:41 PM
- ⌚ End: 06/25/2023, 10:55:42 PM
- ⌚ Response time: 0.69s

Run Metadata

completion_tokens: 16

prompt_tokens: 112

total_tokens: 128

INVOCATION_PARAMS

_type: "openai-chat"

max_tokens: null

model: "gpt-3.5-turbo"

model_name: "gpt-3.5-turbo"

n: 1

request_timeout: null

stop: null

stream: false

temperature: 0

OPTIONS

stop: null

RUNTIME

library: "langchain"

library_version: "0.0.213"

platform: "macOS-13.4-arm64-arm-64bit"

⌚ > Debug > Project: default > Run: ConversationalRetriever...

Run: ConversationalRetrievalChain

✓ CHAIN 2 minutes ago

Chain Input

Jump to bottom
Copy

Run Stats

RBC 1,038 → 99 1,137 tokens

- ⌚ Start: 06/25/2023, 10:55:41 PM
- ⌚ End: 06/25/2023, 10:55:45 PM
- ⌚ Response time: 4.27s

Child Runs



Chain Output

Jump to bottom
Copy

answer: The reason for requiring basic probability and statistics as prerequisites for the class is that the class assumes familiarity with these topics. The instructor assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough. The class will cover machine learning, which is a huge field, and assumes that students already know what random variables are.

Debug > Project: default > Run: ConversationalRet...

Run: ConversationalRetrievalChain

CHAIN 2 minutes ago

Chain Input

[Jump to bottom](#) [Copy](#)

```
question: why are those prerequisites needed?
chat_history:
- content: Is probability a class topic?
  example: false
  additional_kwargs: {}
- content: Yes, probability is a topic that will be assumed to be familiar to students in this class. The instructor mentions that basic probability and statistics are prerequisites for the class and assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough.
  example: false
  additional_kwargs: {}
```

Chain Output

[Jump to bottom](#) [Copy](#)

```
answer: The reason for requiring basic probability and statistics as prerequisites for the class is that the class assumes familiarity with these topics. The instructor assumes that most undergraduate statistics classes, like Stat 116 taught at Stanford, will be more than enough. The class will cover machine learning, which is a huge field, and assumes that students already know what random variables are.
```

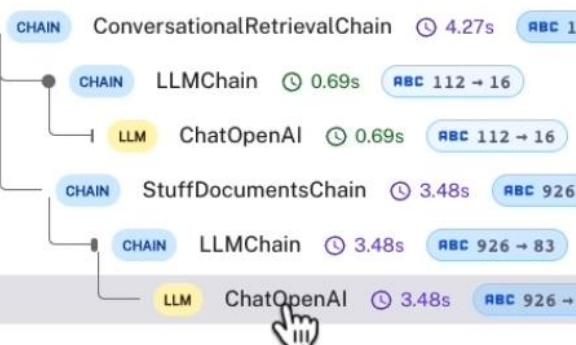
[Rate Run](#) [+ Add to Dataset](#) [Share](#)

Run Stats

RBC 1,038 → 99 1,137 tokens

⌚ Start: 06/25/2023, 10:55:41 PM
⌚ End: 06/25/2023, 10:55:45 PM
⌚ Response time: 4.27s

Child Runs



that you don't know, don't try to make up an answer.

of this class will not be very programming intensive, although we will do some programming, mostly in either MATLAB or Octave. I'll say a bit more about that later.

I also assume familiarity with basic probability and statistics. So most undergraduate statistics class, like Stat 116 taught here at Stanford, will be more than enough. I'm gonna assume all of you know what random variables are, that all of you know what expectation

is, what a variance or a standard deviation is. And in case of some of you, it's been a while since you've seen some of this material. At some of the discussion sections, we'll actually

go over some of the prerequisites, sort of as a refresher course under prerequisite class.

I'll say a bit more about that later as well.

Lastly, I also assume familiarity with basic linear algebra. And again, most undergraduate linear algebra courses are more than enough. So if you've taken courses like Math 51,

103, Math 113 or CS205 at Stanford, that would be more than enough. Basically, I'm gonna assume that all of you know what matrices and vectors are, that you know how to multiply matrices and vectors and multiply matrix and matrices, that you know what a matrix inverse is. If you know what an eigenvector of a matrix is, that'd be

Run Metadata

completion_tokens: 83

prompt_tokens: 926

total_tokens: 1009

INVOCATION_PARAMS

```
_type: "openai-chat"
max_tokens: null
model: "gpt-3.5-turbo"
model_name: "gpt-3.5-turbo"
n: 1
request_timeout: null
stop: null
stream: false
temperature: 0
```

OPTIONS

stop: null

RUNTIME

```
library: "langchain"
library_version: "0.0.213"
platform: "macOS-13.4-arm64-arm-64bit"
runtime: "python"
runtime_version: "3.10.9"
sdk_version: "0.0.17"
```

LLM ChatOpenAI

id:

- langchain
- chat_models
- openai
- ChatOpenAI

lc: 1

[Copy](#)

**Create a chatbot that
works on your
documents**

```

from langchain.embeddings.openai import OpenAIEmbeddings
from langchain.text_splitter import CharacterTextSplitter
from langchain.vectorstores import DocArrayInMemorySearch
from langchain.document_loaders import TextLoader
from langchain.chains import RetrievalQA, ConversationalRetrievalChain
from langchain.memory import ConversationBufferMemory
from langchain.chat_models import ChatOpenAI
from langchain.document_loaders import TextLoader
from langchain.document_loaders import PyPDFLoader

# This will initialize your database and receiver chain
def load_db(file, chain_type, k):
    # load documents
    loader = PyPDFLoader(file)
    documents = loader.load()
    # split documents
    text_splitter = CharacterTextSplitter(chunk_size=650, chunk_overlap=0, sep
    docs = text_splitter.split_documents(documents)
    # define embedding
    embeddings = OpenAIEmbeddings()
    # create vector database from data
    db = DocArrayInMemorySearch.from_documents(docs, embeddings)
    # define retriever
    retriever = db.as_retriever(search_type="similarity", search_kwargs={"k": 4})
    # create a chatbot chain. Memory is managed separately
    qa = ConversationalRetrievalChain.from_llm(
        llm=ChatOpenAI(model_name="gpt-3.5-turbo", temperature=0),
        chain_type=chain_type,
        retriever=retriever,
        return_source_documents=True,
        return_generated_question=True,
    )
    return qa

```



```

def call_load_db(count):
    global qa
    if count == 0 or file_input.value is None: # init or no file
        file = "docs/cs229_lectures/MachineLearning-Lecture01.pdf"
        button_load.button_style="outline"
        qa = load_db(file, "stuff", 4)
        button_load.button_style="solid"
    else:
        file_input.save("temp.pdf")
        file = file_input.filename
        button_load.button_style="outline"
        qa = load_db("temp.pdf", "stuff", 4)
        button_load.button_style="solid"
    clr_history()
    return pn.pane.Markdown(f"Loaded File: {file}")

def convchain(count):
    global qa, chat_history, last_question, last_source
    if count == 0:
        return
    query = inp.value_input
    inp.value = ''
    button_send.button_style="outline"
    result = qa({"question": query, "chat_history": chat_history})
    button_send.button_style="solid"
    chat_history.extend([(query, result["answer"])])
    last_question[0] = result["generated_question"]
    last_source[0] = result["source_documents"]
    print(result['answer'])
    panels.extend([
        pn.Row('User:', pn.pane.Markdown(query, width=600)),
        pn.Row('ChatBot:', pn.pane.Markdown(result['answer'], width=600, style:
    ])

    return pn.WidgetBox(*panels, scroll=True)

def get_lquest(count):
    global last_question
    if count == 0:

```

```
)  
tab3= pn.Column(  
    pn.panel(lhist),  
    pn.layout.Divider(),  
)  
tab4=pn.Column(  
    pn.Row( file_input, button_load, bound_button_load),  
    pn.Row( button_clearhistory, pn.pane.Markdown("Clears chat history. Can use  
    pn.layout.Divider(),  
)  
dashboard = pn.Column(  
    pn.Row(pn.pane.Markdown('# QnA_Bot')),  
    pn.Tabs(('Conversation', tab1), ('Database', tab2), ('Chat History', tab3))  
)  
dashboard
```

QnA_Bot

Conversation Database Chat History Configure

Enter text here... Send Question