# 高阶大语言模型课程

Huajun Zeng

12/13/2024 - 1/31/2025
（12月27日和1月3日放假，共计6次课）
每周五 5pm-7pm PT / 8pm-10pm ET

# 课程安排

| Week | Date | Content | Week | Date | Content |
|------|------|---------|------|------|---------|
| 1 | 2024-12-13 | Retrieval Augmented Generation (RAG) for LLM<br>● Why augmenting LLMs?<br>● Methods for LLM augmentation<br>● Augmenting LLMs with retrieval<br>● Augmenting LLMs with fine tuning | 4 | 2025-01-17 | Pipeline for LLM Applications: From Code to Products<br>● Full stack LLM: tools needed for an LLM application<br>● Case study: build an LLM app from ground |
| 2 | 2024-12-20 | Chatbot Building with LLM APIs<br>● Environment setup<br>● Introduction to LLM APIs<br>● Using chat completion APIs<br>● Using fine tuning APIs | 5 | 2025-01-24 | More LLM Applications and Course Project<br>● Showcase of potential LLM applications for productivity, creativity and more<br>● More advanced LLM applications: AI Agent, Multi-modality, etc.<br>● Introduction to the course project: requirement and discussion |
| 3 | 2025-01-10 | Chatbot Building with LLM Frameworks and Vector Database<br>● Introduction to Langchain and LlamaIndex<br>● Case study: a chatbot from Langchain and vector database | 6 | 2025-01-31 | Project Presentation<br>● Student presentation on the course project |

# 家庭作业回顾

- Using your private domain data (for example, a set of research papers, a set of financial report documents, …)
- Write a python program to answer questions about your data
  - Receive question from console and generate responses in streaming way
  - Using OpenAI SDK (or Anthropic)
- Refer to the examples in https://github.com/hzeng-otterai/chatbot-example/tree/main/backend_api

# Comments 1

- Great to see you are using a notebook file for the homework!
- Good to see a brief description of the program before the code
- A good domain-specific document for testing the chatbot
- Good to use load_dotenv
- Multiple PDF -> text file, simple and effective implementation
- nest_asyncio.apply() seems a necessary step in Notebook environment! Something I didn't know!
- Domain specific queries looks great!

# Comments 2

- Good to see using class to wrap the code, looks very clean!
- Good to implement trim_history function (minor improve suggestion: just return history[-self.max_history:] is sufficient)
- Consider use gpt-4o as the default model which is cheaper and more powerful than gpt-4
- Very good to use asyncio.timeout() to manage long running async tasks
- Extract text from PDF files, nice!
- There are some merge conflicts, see those files contains "<<<<<<< HEAD". These need to be resolved before commit to git.

# 第三课: 建立基于RAG的Chatbot应用

- Using LangChain
- Using Pinecone
- Advanced RAG

# 使用API调用LLM的缺点

```python
from openai import OpenAI

system_prompt_template = """You are a virtual assistant ...
<context>{context}</co
"""

with open("news_result.txt") as in_file:
    context_content = in_file.read()

system_prompt = system_prompt_template.format(context=context_content)

client = OpenAI()
result = client.chat.co
    messages=[
        {"role": "system", "content": system_prompt},
        {"role": "user", "content": "What are those news about?"}
    ],
    model="gpt-4", max_tokens=256, temperature=0.5, stream=True
)

for token in result:
    print(token.choice                      flush=True, end="")
```

Load from a single text file

Load all content into memory

Put all content into LLM context

Use specific LLM API

Manage chat history

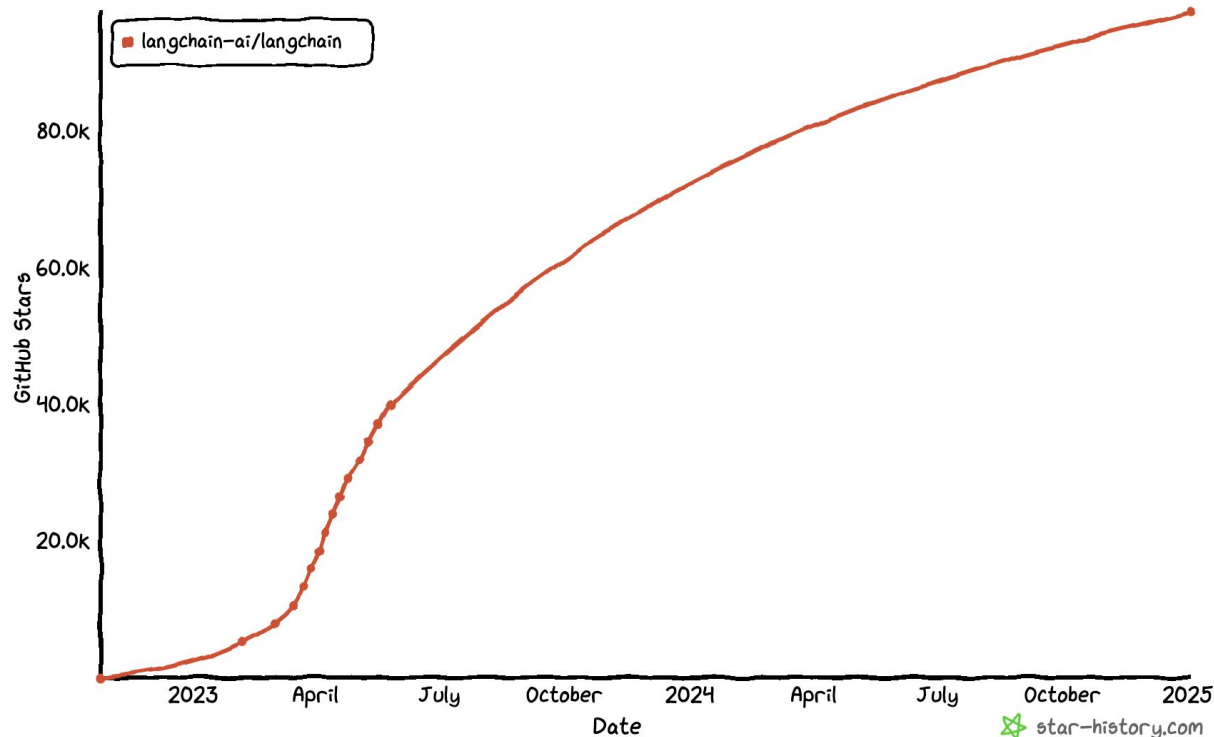Textual result

# 为何要用Langchain框架

- Abstraction - LangChain provides an abstraction layer to interact with LLMs, vector DBs and other modules.
- Fast Prototyping - LangChain provides ready-to-use chains for a lot of applications.
- Built-In Instrumentation for enhanced observability and debuggability
- Langchain can also bring side effect because of its abstraction.

# Langchain的历史

- Release date: October 2022
- Stars: 97.5k

# LangChain代码示例

Refer to `test_00_langchain_example.py`

# LangChain的管道操作符"｜"

A declarative way to easily compose chains together

- Streaming support
- Async support
- Optimized parallel execution


- Good for simple chains
- For complex chains such as branching, cycles, multiple agents, use LangGraph instead.

# 设置环境

```
$ conda create -n langchain_chatbot_env python=3.11
$ cd chatbot-example/backend_langchain
$ pip install -r requirements.txt
$ export OPENAI_API_KEY=...
$ export PINECONE_API_KEY=...
```
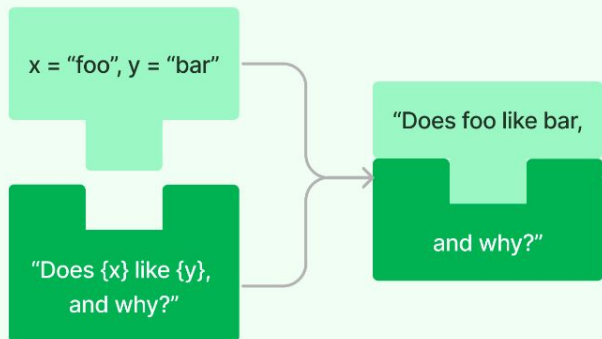
# LangChain的模块

- 模型输入输出: Interface with language models
- 检索: Interface with application-specific data
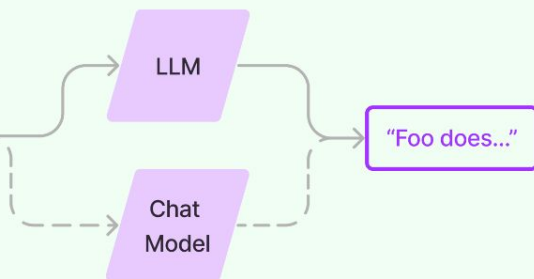- 组合: Construct sequences of calls

# 模型输入输出



**Model I/O**

**Format** · **Predict** · **Parse**

x = "foo", y = "bar"

"Does {x} like {y}, and why?"

"Does foo like bar, and why?"

LLM

Chat Model

"Foo does..."

```
{
    "likes": True,
    "reason": "Because ...."
}
```

# 调用LLM

Refer to test_01_chat.py

# 提示词模版

Refer to test_02_chat_with_prompt.py

# Few-shot提示词模版

Refer to test_03_chat_with_few_shot_prompt.py

# 模版格式

## f-string

```python
fstring_template = "Tell me a {adjective} joke about {content}"
prompt = PromptTemplate.from_template(fstring_template)


prompt.format(adjective="funny", content="basketball")
```

## jinja2

```python
jinja2_template = "Tell me a {{ adjective }} joke about {{ content }}"
prompt = PromptTemplate.from_template(jinja2_template, template_format="jinja2")


prompt.format(adjective="funny", content="basketball")
```

# 输出结构解析

Refer to test_04_output_parser.py.

# 输出结构解析

- Get format instructions: A method which returns a string containing instructions for how the output of a language model should be formatted.
- Parse: A method which takes in a string (assumed to be the response from a language model) and parses it into some structure.

# Pydantic

- Data validation
- Data parsing
- Automatic docs

```python
from datetime import datetime
from pydantic import BaseModel, PositiveInt

class User(BaseModel):
    id: int
    name: str = 'John Doe'
    signup_ts: datetime | None
    tastes: dict[str, int]

external_data = {
    'id': 123,
    'signup_ts': '2019-06-01 12:22',
    'tastes': {
        'wine': 9,
        'cheese': 7
    }
}

user = User(**external_data)
```

# 检索

**Data connection**

**Source**

**Load**

**Transform**

**Embed**

**Store**

**Retrieve**

XXXXXXXXXX
XXXXXXXXXX
XXXXXXXXXX
XXXXXXXXXX
XXXXXXXXXX
XXXXXXXXXX
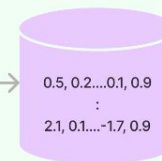
XXXXXXXXXX
XXXXXXXXXX

XXXXXXXXXX
XXXXXXXXXX

XXXXXXXXXX
XXXXXXXXXX

0.5, 0.2....0.1, 0.9
-0.1, 0.4....1.4, 5.9

0.2, 0.7....2.1, -1.2
4.1, 3.4....-1.5, 2.5

5.5, -0.3....0.8, 2.3
2.1, 0.1....-1.7, 0.9

0.5, 0.2....0.1, 0.9
:
2.1, 0.1....-1.7, 0.9

XXXXXXXXXX
XXXXXXXXXX

# 文档加载

- CSV
- File Directory
- HTML
- JSON
- Markdown
- PDF

```python
from langchain.document_loaders import TextLoader

loader = DirectoryLoader(
    '../',
    glob="**/*.md",
    loader_cls=TextLoader
)

docs = loader.load()
```

# 文档转换: Text splitters

- RecursiveCharacterTextSplitter
- It tries to split on them in order until the chunks are small enough. The default list is ["\n\n", "\n", " ", ""]

```
documents = TextLoader("./state_of_the_union.txt").load()

text_splitter = RecursiveCharacterTextSplitter.from_tiktoken_encoder(
    chunk_size=500,
    chunk_overlap=10
)

texts = text_splitter.split_documents(documents)
```

# Embedding模型

- Embedding providers: OpenAI, Voyage, Hugging Face, etc.

```python
from langchain.embeddings import OpenAIEmbeddings

embeddings_model = OpenAIEmbeddings()

embeddings = embeddings_model.embed_documents([
    "Hi there!",
    "Oh, hello!",
    "What's your name?",
])

embedded_query = embeddings_model.embed_query("What is your name?")
```

# 向量数据库

- Local stores: Chroma, FAISS, etc.
- On-premis stores: Qdrant, Milvus, PGVector, etc.
- Cloud stores: Pinecone, etc.

```python
raw_documents = TextLoader('./state_of_the_union.txt').load()
text_splitter = CharacterTextSplitter(chunk_size=1000, chunk_overlap=0)
documents = text_splitter.split_documents(raw_documents)
db = Chroma.from_documents(documents, OpenAIEmbeddings())

query = "What did the president say about Ketanji Brown Jackson"
docs = db.similarity_search(query)
```
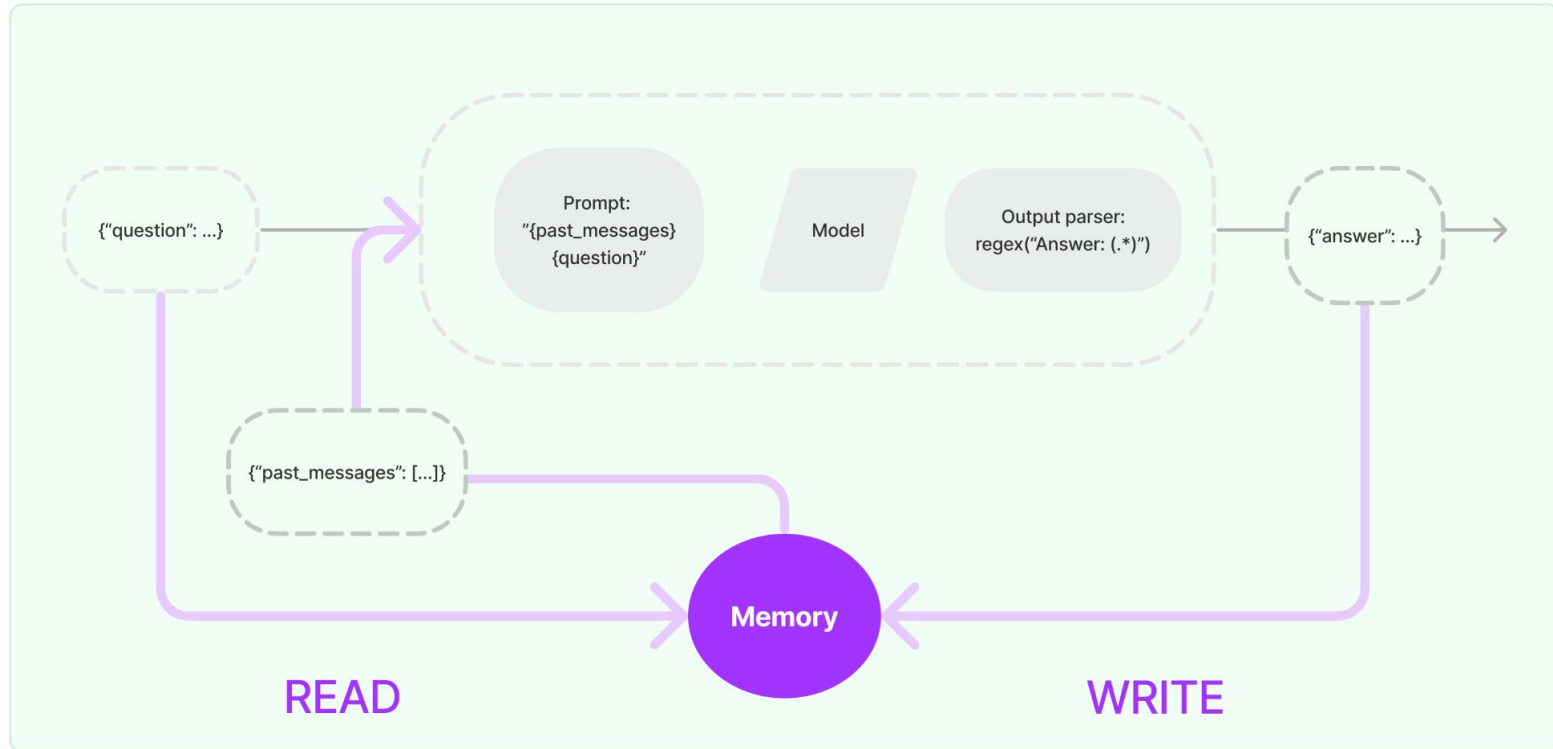
# Retriever

- Retriever class returns Documents given a text query.
- It is more general than a vector store. A retriever does not need to be able to store documents, only to return (or retrieve) it. Vector stores can be used as the backbone of a retriever, but there are other types of retrievers as well.

```python
db = FAISS.from_documents(texts, embeddings)

retriever = db.as_retriever()
docs = retriever.get_relevant_documents("what did he say about jackson")
```

# Memory

# 给Chatbot增加Memory

Refer to test_06_chat_with_memory.py

# 一个最小的Chatbot

Refer to test_07_chat_with_memory_rag_continous.py

# 使用Pinecone作为向量数据库

- For handling large scale of dataset
- For a hosted solution without a lot of DevOps overhead

Refer to test_08_pinecone_add_doc.py and test_09_pinecone.py

# 优化RAG: 基本优化技巧

- Prompt Engineering
- Embeddings
    - MTEB leaderboard
- Chunk Sizes
- Hybrid Search
    - Embedding search plus keyword search
- Metadata Filters
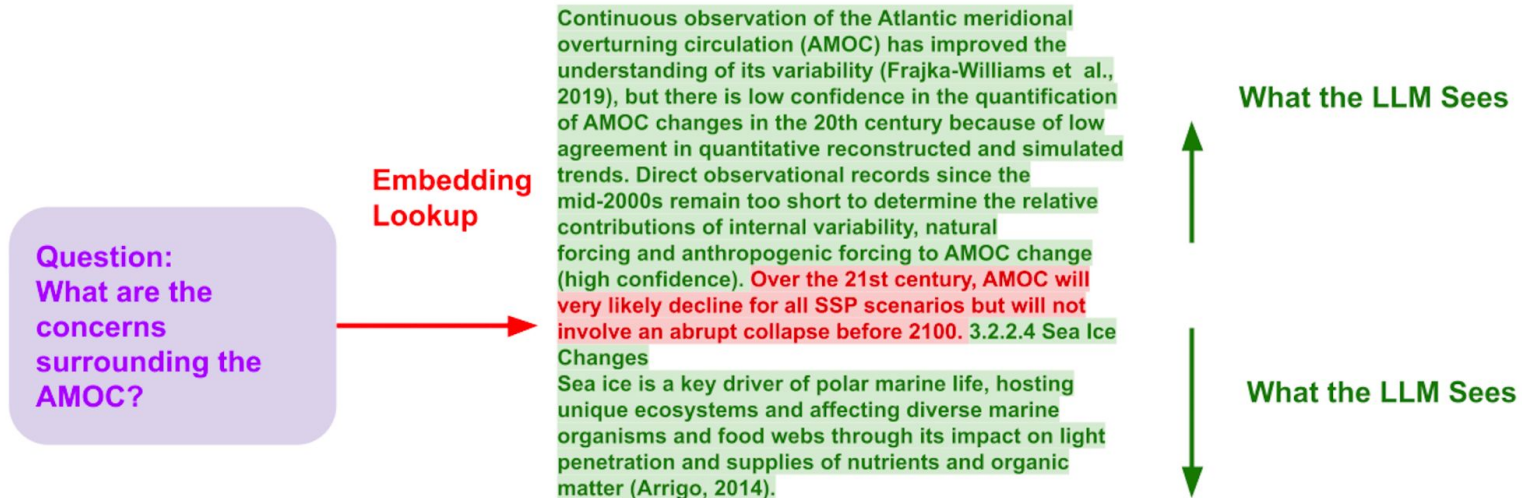    - Filtering based on metadata of documents

# 优化RAG: 高级优化技巧

- Small-to-big retrieval
- Query transformation
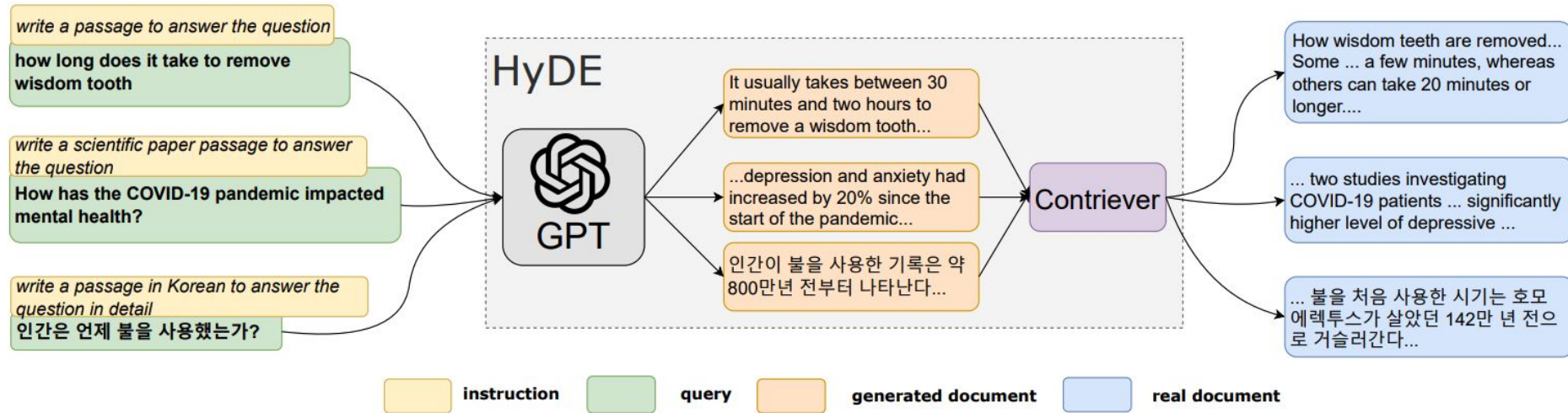- Reranking
- Recursive retrieval
- Embedded tables
- …

# Small to Big Retrieval

- Using smaller units for embedding but using expanded text for LLM inference
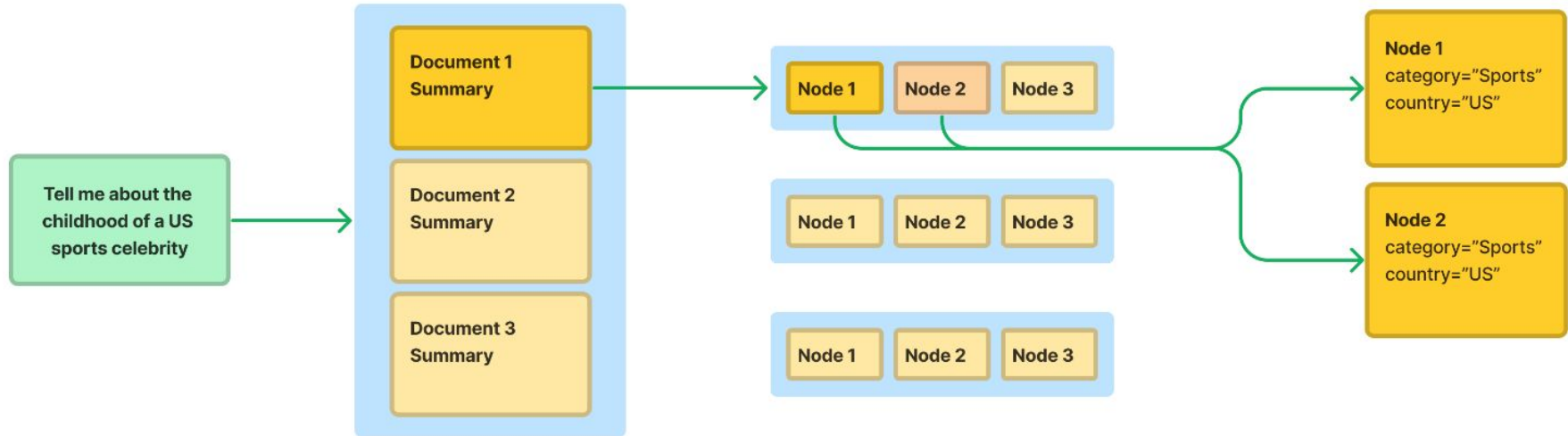- Refer to test_10_pinecone_with_small_to_big.py

# Query Transformation - HyDE

- HyDE: Hypothetical Document Embeddings
- Refer to test_11_pinecone_with_hyde.py

# Recursive Retrieval

**Document Hierarchies (Summaries + Raw Chunks) + Recursive Retrieval**

# Reranking

- A step to rerank the top n results returned from the embedding search
- Why reranking
  - Embedding model might be not strong enough
  - To use context information or user-specific information
- Choosing a reranking model
  - LLMRerank
  - Cohere Rerank
  - SentenceTransformerRerank
  - …

# Homework

- Implement a RAG-based Chatbot using Langchain
    - A console program to continuously receive user input and respond
    - Try smaller dataset first and then larger dataset
    - Refer to the examples in https://github.com/hzeng-otterai/chatbot-example/tree/main/backend_langchain

# Questions?