# 第一章、任务描述

Retrieval Augmented Generation（RAG）with LLM是目前比较热门的应用之一，实现并不难，但提取内容的准确度是目前普遍存在的问题。想要提高准确度，需要考虑多个细节，例如：

- 如何保证文档切片不会造成相关内容丢失
- 切片大小如何控制
- 如何保证召回内容跟问题是相关的等等。

请提供相关的代码实现，尽可能的解决RAG准确度低的问题。

**Key Function**

- Langchain已经提供了一些api接口，可以调用，但需要写明白解决了哪方面的问题，同时也应该有自己的改进
- 提供一个demo，去展示该方案使用前后的效果对比，给出准确度定量的估计，不少于5个例子
- 加分：Tree-of-Thought，Graph-of-Thought，Knowledge-Graph

# 第二章、方案汇总

本方案基于LLM+RAG构建了一个金融领域的QA Demo，然后通过bad case分析RAG目前存在的一些问题，然后通过数据，召回和生成等多个方面逐步对基础方案进行改进和优化，从而达到更好的效果（**准确率从15%提升到75+%**）。具体地：

**第三章：基础方案** 根据Langchain官方示例，基于ChatGPT构造了一个金融领域的QA Demo，使用了如下数据：

- 汇丰2022年度报告
- 汇丰官网的FAQ问答对，共300多个
- 各大证券公司关于最近中央金融会议的观点文章
- 金融相关的LLM文章

**第四章：评估方法**

定义测试和评估RAG的数据集和指标。

**第五章：改进方案**

整体方案采用Mixture of Expert (MoE)架构，先使用LLM进行意图分类，然后分别调用领域专家RAG。我们对bad case进行分析，针对性的对领域专家RAG进行改进。

**第六章：总结**

对整个项目进行总结。

注：由于时间有限，本项目Demo没有进行充分的Train/Test验证，而从bad case的角度出发，探索和验证若干提升RAG的方法。

## ▾ 第三章、基础方案

按照LangChain教程的默认设置：https://python.langchain.com/docs/use_cases/question_answering/

```
!pip install -q openai==0.28.1
!pip install -q langchain==0.0.330
!pip install -q -U langchainhub
!pip install -q -U chromadb tiktoken pypdf pymupdf lark
!pip install -q -U FlagEmbedding sentence_transformers
!pip install -q -U transformers
!pip install -q rank_bm25 cohere
!pip install -q -U evaluate rouge_score
!pip install -q -U unstructured[all-docs] pydantic lxml
```

```
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 2.7/2.7 MB 18.7 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 84.1/84.1 kB 2.7 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
  Building wheel for rouge_score (setup.py) ... done
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 1.7/1.7 MB 13.3 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 395.8/395.8 kB 20.2 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 358.9/358.9 kB 22.4 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 275.1/275.1 kB 22.7 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 981.5/981.5 kB 32.9 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 3.3/3.3 MB 49.5 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 10.1/10.1 MB 80.1 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 15.7/15.7 MB 91.7 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 239.6/239.6 kB 29.9 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 5.6/5.6 MB 116.7 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 62.0/62.0 kB 8.8 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 101.8/101.8 kB 14.0 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 19.2/19.2 MB 89.5 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 45.7/45.7 kB 6.0 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 5.9/5.9 MB 117.0 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 2.0/2.0 MB 101.3 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 154.8/154.8 kB 20.2 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 112.2/112.2 kB 15.0 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 42.2/42.2 kB 1.7 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 49.0/49.0 kB 5.5 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 112.5/112.5 kB 15.4 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 2.2/2.2 MB 99.1 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 79.5/79.5 kB 11.8 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 3.0/3.0 MB 108.9 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 117.0/117.0 kB 16.8 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
  Building wheel for python-pptx (setup.py) ... done
  Building wheel for langdetect (setup.py) ... done
  Building wheel for olefile (setup.py) ... done
  Building wheel for iopath (setup.py) ... done
  Building wheel for antlr4-python3-runtime (setup.py) ... done
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviou
lida 0.0.10 requires kaleido, which is not installed.
```

```python
import os
import evaluate
import openai
import pandas as pd

from copy import deepcopy

from langchain import hub
from langchain.chains import RetrievalQA
from langchain.chat_models import AzureChatOpenAI
from langchain.document_loaders import DirectoryLoader, PyPDFLoader, TextLoader
from langchain.embeddings import HuggingFaceBgeEmbeddings, OpenAIEmbeddings
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain.vectorstores import Chroma


os.environ["OPENAI_API_TYPE"] = "OPENAI_API_TYPE"
os.environ["OPENAI_API_VERSION"] = "OPENAI_API_VERSION"
os.environ["OPENAI_API_BASE"] = "OPENAI_API_BASE"
os.environ["OPENAI_API_KEY"] = "OPENAI_API_KEY"


OPENAI_API_KEY = "OPENAI_API_KEY"
OPENAI_DEPLOYMENT_NAME = "gpt-35-turbo"
MODEL_NAME = "gpt-35-turbo"
```

## ▾ 1、加载数据

```python
from google.colab import drive
drive.mount('/content/drive')
PROJ_DIR = "/content/drive/My Drive/Colab Notebooks/HSBCRAG"
DATA_DIR = f"{PROJ_DIR}/data"
```

```
 ⤓  Mounted at /content/drive
```

```python
faq_loader = DirectoryLoader(f'{DATA_DIR}/hsbc_faqs/', glob="./*.txt", loader_cls=TextLoader)
faqs = faq_loader.load()
```

```
annual_report_loader = DirectoryLoader(f'{DATA_DIR}/hsbc_annual_reports/', glob="./*.pdf", loader_cls=PyPDFLoa
annual_reports = annual_report_loader.load()


article_loader = DirectoryLoader(f'{DATA_DIR}/jinronghuiyi_articles/', glob="./*.pdf", loader_cls=PyPDFLoader)
articles = article_loader.load()


paper_loader = DirectoryLoader(f'{DATA_DIR}/llm_papers/', glob="./*.pdf", loader_cls=PyPDFLoader)
papers = paper_loader.load()


documents = faqs + annual_reports + articles + papers


print(f"""faq={len(faqs)}
annual_reports={len(annual_reports)}
articles={len(articles)}
papers={len(papers)}
documents={len(documents)}""")
```

```
faq=6
annual_reports=172
articles=141
papers=230
documents=549
```

## ▾ 2、切分文档

```
text_splitter = RecursiveCharacterTextSplitter(chunk_size=500, chunk_overlap=0)
texts = text_splitter.split_documents(documents)
```

```
len(texts)
```

```
2743
```

## ▾ 3、构建索引

OpenAIEmbeddings速度较慢，并且容易碰到RateLimitError问题，我们使用BGE来作为baseline

```
# get around RateLimitError by increasing max_retries
# https://github.com/langchain-ai/langchain/issues/2493
# embedding = OpenAIEmbeddings(
#     deployment="text-embedding-ada-002",
#     show_progress_bar=True,
#     maxConcurrency=5,
#     # chunk_size=1,
#     disallowed_special=(),
#     max_retries=100,
# )

model_name = "BAAI/bge-small-zh"
model_kwargs = {'device': 'cuda'}
encode_kwargs = {'normalize_embeddings': True}
embedding = HuggingFaceBgeEmbeddings(
    model_name=model_name,
    model_kwargs=model_kwargs,
    encode_kwargs=encode_kwargs,
    query_instruction="为这个句子生成表示以用于检索相关文章："
)
```

```
len(embedding.embed_query("text"))
```

    512

```
Chroma().delete_collection()
baseline_vectordb = Chroma.from_documents(
    documents=texts,
    embedding=embedding,
    collection_name="baseline"
)
```

```
baseline_vectordb._collection.count()
```

    2743

```
baseline_retriever = baseline_vectordb.as_retriever()
```

## 4、Prompt

```
default_rag_prompt = hub.pull("rlm/rag-prompt")
```

## 5、配置LLM

```
chatgpt = AzureChatOpenAI(
    openai_api_key=OPENAI_API_KEY,
    deployment_name=OPENAI_DEPLOYMENT_NAME,
    model_name=MODEL_NAME,
    temperature=0
)
```

```
chatgpt.predict("申请汇丰中国信用卡有哪些步骤？ ")
```

    '申请汇丰中国信用卡的步骤如下：\n\n1. 在汇丰中国官网或手机APP上选择信用卡产品，填写个人信
    息并提交申请。\n\n2. 提交申请后，等待汇丰银行的审核，通常需要1-2个工作日。\n\n3. 审核通过
    后，汇丰银行会联系您确认申请信息，并告知您信用卡的额度和发卡时间。\n\n4. 在收到信用卡后，需
    要激活并设置密码，可以通过汇丰中国官网或手机APP完成。\n\n5. 使用信用卡时，需要注意还款日期

```
chatgpt.predict("汇丰中国信用卡有哪些密码？ ")
```

    '汇丰中国信用卡有以下几种密码：\n\n1. 信用卡密码：用于在ATM机上取现金或进行其他操作时输入
    的密码，一般为6位数字。\n\n2. 网上银行密码：用于登录汇丰中国网上银行进行账户管理和交易的密
    码，一般为8-30位数字、字母或符号的组合。\n\n3. 手机银行密码：用于登录汇丰中国手机银行进行
    账户管理和交易的密码、一般为6-8位数字、字母或符号的组合。\n\n4. 短信验证码：用于在进行某些

## 6、构建RAG

```
def initialize_rag(llm, retriever, prompt=default_rag_prompt):
    chain_type_kwargs = {"prompt": prompt}
    rag = RetrievalQA.from_chain_type(
        llm=llm,
        chain_type="stuff",
        retriever=retriever,
        return_source_documents=True,
        chain_type_kwargs=chain_type_kwargs
    )
    return rag
```

```
baseline_rag = initialize_rag(llm=chatgpt, retriever=baseline_retriever)
```

## 第四章、评估方法

```python
class Evaluator:
    @staticmethod
    def score_any(predict_str, answers):
        for answer in answers:
            if answer in predict_str:
                return 1
        return 0

    @staticmethod
    def score_ratio(predict_str, answers):
        s = 0
        for answer in answers:
            if answer in predict_str:
                s += 1
        return s / len(answers)

    @staticmethod
    def score_string(predict_str, answer):
        rouge = evaluate.load("rouge")
        predict_str = predict_str.replace("\n", "").replace("+s", " ")
        answer = answer.replace("\n", "").replace("+s", " ")
        rouge_results = rouge.compute(predictions=[predict_str], references=[[answer]])
        return rouge_results["rougeL"]

    @staticmethod
    def score_metadata(source_documents, answer, field):
        for source in source_documents:
            if source.metadata.get(field) != answer:
                return 0
        return 1

    @staticmethod
    def test_rag(rag, query, verbose=0, return_source=False):
        response = rag(query)
        if verbose > 0:
            print("-"*200)
            print(f"Question: {query}")
            print(f"Response: {response['result']}")
            if verbose > 1 and "source_documents" in response:
                print('\n\nSources:')
                for source in response["source_documents"]:
                    print(f"page={source.metadata.get('page')}, source={source.metadata.get('source')}")
                    if verbose > 2:
                        print(source.page_content)
        if return_source:
            return response.get("result",""), response.get("source_documents",[])
        else:
            return response.get("result","")

    @staticmethod
    def test_rag_all(rag, question_answer_pairs, verbose=0):
        results = deepcopy(question_answer_pairs)
        for category,qas in results.items():
            for qa in qas:
                qa["response"], qa["source_documents"] = Evaluator.test_rag(
                    rag=rag, query=qa["question"], verbose=verbose, return_source=True
                    )
        return results

    @staticmethod
    def evaluate_rag_results(results):
        results = deepcopy(results)
        for category,qas in results.items():
            for qa in qas:
                w = qa.get("weight", 1.0)
                if qa["type"] == "any":
                    qa["score"] = w * Evaluator.score_any(qa["response"], qa["answer"])
                elif qa["type"] == "ratio":
                    qa["score"] = w * Evaluator.score_ratio(qa["response"], qa["answer"])
                elif qa["type"] == "string":
                    qa["score"] = w * Evaluator.score_string(qa["response"], qa["answer"])
                elif qa["type"].startswith("metadata"):
                    qa["score"] = w * Evaluator.score_metadata(qa["source_documents"], qa["answer"], qa["type'
```

```python
        return results

    @staticmethod
    def compute_rag_score(results):
        s = 0
        cnt = 0
        for category,qas in results.items():
            for qa in qas:
                s += qa["score"]
                cnt += 1
        return s/cnt

    @staticmethod
    def inspect_rag_results(results):
        for category,qas in results.items():
            for qa in qas:
                if qa["score"] != 1:
                    print(f"category={category}, question={qa['question']}, score={qa['score']}, response={qa


question_answer_pairs = {
    "huifeng_faq": [
        {
            "question": "如何环球转账？ ",
            "answer": """汇丰环球转账功能，可通过登录网上银行\n\n\n在"我的银行"－"转账及货币兑换"－"转账及货币兑换"。\n在新
            "type": "string"
        },
        {
            "question": "如何查看汇丰环球转账历史记录？ ",
            "answer": "可以。您能查看过去 12 个月以内的环球转账记录。请在汇丰环球网上银行页面点击"环球转账历史记录"，然后选择
            "type": "string"
        },
        {
            "question": "汇丰信用卡有几个密码？ ",
            "answer": ["查询密码","交易密码"],
            "type": "ratio",
        },
        {
            "question": "如果赎回申请成功，我多久才可以取回资金？ ",
            "answer": "若您投资的是代客境外理财计划－开放式海外基金型产品，在海外基金管理人接受银行的赎回要求后，银行将在从海外基
            "type": "string"
        },
        {
            "question": "卓越理财客户服务月费是多少？ ",
            "answer": "如您在汇丰中国的同一个卓越理财客户号码下的所有账户之月内日均总余额低于500,000元人民币/等值外币，本行将
            "type": "string"
        }

    ],
    "huifeng_annual_report": [
        {
            # P7
            "question": "汇丰董事会下属哪些委员会？ ",
            "answer": ["审计委员会","风险及消费者权益保护委员会","关联交易控制委员会","薪酬委员会","提名委员会"],
            "type": "ratio"
        },
        {
            "question": "2022年度，汇丰中国获得《财富管理》杂志什么称号？ ",
            "answer": ["2022年度最佳中国外资私人银行"],
            "type": "ratio"
        },
        {
            "question": "汇丰在什么时候成为首家协助QFI完成北交所交易的外资托管行？ ",
            "answer": ["2022年1月"],
            "type": "any"
        },
        {
            "question": "汇丰的WPB是哪个部门？ ",
            "answer": ["财富管理及个人银行业务部"],
            "type": "any"
        },
        {
            "question": "财富管理及个人银行业务部总监是哪位？ ",
            "answer": ["孙丹莹"],
            "type": "any"
```

```
    },
    {
        # P27~P29
        "question": "与汇丰银行业务相关的主要风险有哪些? ",
        "answer": ["信用风险", "市场风险", "财资风险", "操作风险", "抗逆力风险", "监管合规风险", "金融犯罪风险", "声誉
        "type": "ratio",
    },
    {
        "question": "截至2022年末，汇丰银行资产总计人民币多少亿元? ",
        "answer": ["5,968.5","5,968.5亿元","人民币5,968.5亿元","596,845","人民币596,845百万元","596,845人民币百
        "type": "any",
    },
    {
        "question": "截至2022年末，汇丰银行负债合计人民币多少亿元? ",
        "answer": ["5,386.7","5,386.7亿元","人民币5,386.7亿元","538,674","人民币538,674百万元","538,674人民币百
        "type": "any",
        },
    {
        "question": "汇丰银行在2022会计年度的营业收入为人民币多少亿元? ",
        "answer": ["149.4","149.4亿元","人民币149.4亿元"],
        "type": "any",
    },
    {
        "question": "汇丰银行在2022会计年度的营业支出为人民币多少亿元? ",
        "answer": ["80.7","80.7亿元","人民币80.7亿元"],
        "type": "any",
    },
    {
        "question": "汇丰银行在2022会计年度的净利润为人民币多少亿元? ",
        "answer": ["60.4","60.4亿元","人民币60.4亿元"],
        "type": "any",
    },
    {
        "question": "2022年，汇丰银行不良贷款余额是人民币多少亿元? ",
        "answer": ["5.1","5.1亿元","人民币5.1亿元"],
        "type": "any",
    },
    {
        "question": "2022年，汇丰银行不良贷款率是多少? ",
        "answer": ["0.21%"],
        "type": "any",
    },
    {
        "question": "2022年吸收个人活期存款多少? ",
        "answer": ["47,200,715","47,200,715千元","人民币47,200,715千元","47,200,715人民币千元"],
        "type": "any",
    },
    {
        "question": "2022年吸收个人定期存款多少? ",
        "answer": ["37,742,065","37,742,065千元","人民币37,742,065千元","37,742,065人民币千元"],
        "type": "any",
    },
        {
        "question": "2021年吸收个人活期存款多少? ",
        "answer": ["44,839,631","44,839,631千元","人民币44,839,631千元","44,839,631人民币千元"],
        "type": "any",
    },
    {
        "question": "2021年吸收个人定期存款多少? ",
        "answer": ["26,185,953","26,185,953千元","人民币26,185,953千元","26,185,953人民币千元"],
        "type": "any",
    },
],
"jinronghuiyi_article": [
    {
        "question": "中金对中央金融工作会议的观点有哪些? ",
        "answer": "中金公司",
        "type": "metadata.公司",
    },
    {
        "question": "广发证券对中央金融工作会议的观点有哪些? ",
        "answer": "广发证券",
        "type": "metadata.公司",
    },
```

```
        {
            "question": "平安证券对中央金融工作会议的观点有哪些? ",
            "answer": "平安证券",
            "type": "metadata.公司",
        }
    ],
    "llm_paper": [
        {
            "question": "What is the number of training tokens for Llama 2?",
            "answer": ["2.0T","2 trillion"],
            "type": "any",
        },
        {
            "question": "What is the author of Llama?",
            "answer": ['Hugo Touvron','Thibaut Lavril','Gautier Izacard','Xavier Martinet','Marie-Anne Lachaux
            "type": "ratio",
        },
        {
            "question": "What is the affiliation of the first author of Llama?",
            "answer": ["Meta AI","GenAI, Meta"],
            "type": "any",
        },
        {
            "question": "What is the affiliation of the first author of DISC-FinLLM?",
            "answer": ["Fudan University and Huazhong University of Science and Technology","Huazhong Univers:
            "type": "any",
        },
        {
            "question": "What are the common authors of Llama and Llama 2?",
            "answer": ['Aurelien Rodriguez','Hugo Touvron','Marie-Anne Lachaux','Naman Goyal','Thibaut Lavril
            "type": "ratio",
        },
        {
            "question": "Is the first author of Llama and Llama 2 the same? If yes, please output <<YES>>, oth
            "answer": ["<<Yes>>"],
            "type": "any",
        },
        {
            "question": "Is there any LLM in financial area? What are they?",
            "answer": ["BloombergGPT","FinGPT","DISC-FinLLM","ConFIRM","FinVis-GPT"],
            "type": "ratio"
        },
    ]
}
```

```
sum([len(qs) for qs in question_answer_pairs.values()])
```

```
    32
```

```
evaluator = Evaluator()
```

# 第五章、改进方案

在本章，我们从具体的case出发，依次从数据、召回和生成等三个主要方面来改进RAG。在这之后，我们将给出完整的方案实现。最后，我们将探索COT，KG和Tool在改进RAG上面的应用。

## 1、数据层

### a、【done】利用好Meta数据

```
a = evaluator.test_rag(baseline_rag, "平安证券对中央金融工作会议的观点有哪些? ", verbose=2)
```

```
    --------------------------------------------------------------------------------------------------------
    Question: 平安证券对中央金融工作会议的观点有哪些?
    Response: 平安证券对中央金融工作会议的观点没有在提供的文本中提到。


    Sources:
    page=0, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/jinronghuiyi_articles/国联证券-中央金融工作会议点评: 中央金
    page=1, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/jinronghuiyi_articles/国联证券-中央金融工作会议点评: 中央金
```

```
        page=0, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/jinronghuiyi_articles/国信证券-中央金融工作会议解读：推动
        page=0, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/jinronghuiyi_articles/东莞证券-中央金融工作会议点评：加快
```

以上的case，从召回来源来看，跟问题不相关。考虑到文件名存在结构化信息，譬如"XX证券-YY.pdf"，可以将该信息添加到meta信息里面，通过meta数据进行过滤，提供召回的相关性。

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


metadata_field_info = [
    AttributeInfo(
        name="公司",
        description="公司",
        type="string",
    )
]
content_description = "资料"


def augment_metadata():
    for article in articles:
        title = article.metadata["source"].split("/")[-1].split(".pdf")[0]
        article.metadata["公司"] = title.split("-")[0]

    for faq in faqs:
        faq.metadata["公司"] = "汇丰银行"

    for report in annual_reports:
        report.metadata["公司"] = "汇丰银行"

    for paper in papers:
        paper.metadata["公司"] = ""

augment_metadata()


documents = faqs + annual_reports + articles + papers

text_splitter = RecursiveCharacterTextSplitter(chunk_size=500, chunk_overlap=0)
texts = text_splitter.split_documents(documents)

Chroma().delete_collection()
metadata_vectordb = Chroma.from_documents(
    documents=texts,
    embedding=embedding,
    collection_name="add_metadata"
)
metadata_vectordb._collection.count()

metadata_retriever = SelfQueryRetriever.from_llm(
    chatgpt,
    metadata_vectordb,
    content_description,
    metadata_field_info,
    use_original_query=True,
    verbose=True,
)

metadata_rag = initialize_rag(llm=chatgpt, retriever=metadata_retriever)


a = evaluator.test_rag(metadata_rag, query="平安证券对中央金融工作会议的观点有哪些？", verbose=2)
```

```
    ------------------------------------------------------------------------------------------------------------------
    Question: 平安证券对中央金融工作会议的观点有哪些？
    Response: 平安证券认为，中央金融工作会议肯定了金融回归本源的成绩，但也指出了金融服务实体经济质效不高等问题。未来一段时间，强监管和防风险仍是重

    Sources:
    page=0, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/jinronghuiyi_articles/平安证券-首席宏观报告：中央金融工作
    page=3, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/jinronghuiyi_articles/平安证券-首席宏观报告：中央金融工作
    page=3, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/jinronghuiyi_articles/平安证券-首席宏观报告：中央金融工作
    page=1, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/jinronghuiyi_articles/平安证券-首席宏观报告：中央金融工作
```

可以看到，使用meta进行过滤后，可以得到比较相关的召回结果。

## ▾ b、【done】更好地处理PDF数据

从下面的case可以看到，常用的PyPDFLoader对于年报的解析不够理想（主要是格式和表格数据），出现了严重的叠词现象，会对RAG的召回和生成产生影响。我们可以使用其他的PDF解析工具来避免此类问题，也能提高RAG的效果。

```
for report in annual_reports:
    if "风风险险" in report.page_content:
        print(report)

    page_content='- 21 -汇丰银行（中国)有限公司  \n \n董事会报告 -薪酬报告（续）\n \n - 21 - \nRESTRICTED    \n薪薪酬酬报报告告(续)  \n
    page_content='- 24 -汇丰银行(中国)有限公司  \n \n董事会报告 -风险管理  \n \n - 24 - \nRESTRICTED    \n风风险险管管理理  \n \n我我们
    page_content='- 25 -汇丰银行（中国)有限公司  \n \n董事会报告 -风险管理（续）\n \n - 25 - \nRESTRICTED    \n风风险险管管理理(续)  \
    page_content='- 26 -汇丰银行（中国)有限公司  \n \n董事会报告 -风险管理（续）\n \n - 26 - \nRESTRICTED    \n风风险险管管理理(续)  \
    page_content='- 27 -汇丰银行（中国)有限公司  \n \n董事会报告 -风险管理（续）\n \n - 27 - \nRESTRICTED    \n风风险险管管理理(续)  \
    page_content='- 28 -汇丰银行（中国)有限公司  \n \n董事会报告 -风险管理（续）\n \n - 28 - \nRESTRICTED    \n风风险险管管理理(续)  \
    page_content='- 29 -汇丰银行（中国)有限公司  \n \n董事会报告 -风险管理（续）\n \n - 29 - \nRESTRICTED    \n风风险险管管理理(续)  \
    page_content='- 30 -汇丰银行（中国)有限公司  \n \n董事会报告 -风险管理（续）\n \n - 30 - \nRESTRICTED    \n风风险险管管理理(续)  \
    page_content='- 31 -汇丰银行（中国)有限公司  \n \n董事会报告 -风险管理（续）\n \n - 31 - \nRESTRICTED    \n风风险险管管理理(续)  \
    page_content='- 32 -汇丰银行（中国)有限公司  \n \n董事会报告 -风险管理（续）\n \n - 32 - \nRESTRICTED    \n风风险险管管理理(续)  \
    page_content='- 33 -汇丰银行（中国)有限公司  \n \n董事会报告 -风险管理（续）\n \n - 33 - \nRESTRICTED    \n风风险险管管理理(续)  \
    page_content='- 34 -汇丰银行（中国)有限公司  \n \n董事会报告 -风险管理（续）\n \n - 34 - \nRESTRICTED    \n风风险险管管理理(续)  \
    page_content='- 35 -汇丰银行（中国)有限公司  \n \n董事会报告 -风险管理（续）\n \n - 35 - \nRESTRICTED    \n风风险险管管理理(续)  \
```

```
from langchain.document_loaders import PyMuPDFLoader

annual_report_loader = DirectoryLoader(f'{DATA_DIR}/hsbc_annual_reports/', glob="./*.pdf", loader_cls=PyMuPDFl
annual_reports = annual_report_loader.load()

article_loader = DirectoryLoader(f'{DATA_DIR}/jinronghuiyi_articles/', glob="./*.pdf", loader_cls=PyMuPDFLoade
articles = article_loader.load()

paper_loader = DirectoryLoader(f'{DATA_DIR}/llm_papers/', glob="./*.pdf", loader_cls=PyMuPDFLoader)
papers = paper_loader.load()

for report in annual_reports:
    if "风风险险" in report.page_content:
        print(report)
```

可以看到，叠词的问题解决了。

```
documents = faqs + annual_reports + articles + papers

text_splitter = RecursiveCharacterTextSplitter(chunk_size=500, chunk_overlap=0)
texts = text_splitter.split_documents(documents)

Chroma().delete_collection()
better_pdf_vectordb = Chroma.from_documents(
    documents=texts,
    embedding=embedding,
    collection_name="better_pdf"
)

better_pdf_retriever = better_pdf_vectordb.as_retriever()

better_pdf_rag = initialize_rag(llm=chatgpt, retriever=better_pdf_retriever)
```

## ▾ 2、召回层

## ▾ a、【done】更好的Embedding

我们使用更好的Embedding可以获得更好的召回效果。参考C-MTEB上面的任务，我们选取了BGE。
（https://github.com/FlagOpen/FlagEmbedding）

```
model_name = "BAAI/bge-large-zh-v1.5"
model_kwargs = {'device': 'cuda'}
```

```python
encode_kwargs = {'normalize_embeddings': True}
embedding = HuggingFaceBgeEmbeddings(
    model_name=model_name,
    model_kwargs=model_kwargs,
    encode_kwargs=encode_kwargs,
    query_instruction="为这个句子生成表示以用于检索相关文章: "
)
```

| Downloading (…)3d520/.gitattributes: 100% | 1.52k/1.52k [00:00<00:00, 149kB/s] |
| Downloading (…)_Pooling/config.json: 100% | 191/191 [00:00<00:00, 17.7kB/s] |
| Downloading (…)39c423d520/README.md: 100% | 27.7k/27.7k [00:00<00:00, 2.63MB/s] |
| Downloading (…)c423d520/config.json: 100% | 1.00k/1.00k [00:00<00:00, 92.5kB/s] |
| Downloading (…)ce_transformers.json: 100% | 124/124 [00:00<00:00, 11.4kB/s] |
| Downloading model.safetensors: 100% | 1.30G/1.30G [00:03<00:00, 342MB/s] |
| Downloading pytorch_model.bin: 100% | 1.30G/1.30G [00:03<00:00, 402MB/s] |

```python
Chroma().delete_collection()
better_emb_vectordb = Chroma.from_documents(
    documents=texts,
    embedding=embedding,
    collection_name="better_emb"
)

better_emb_retriever = better_emb_vectordb.as_retriever()

better_emb_rag = initialize_rag(llm=chatgpt, retriever=better_emb_retriever)

a = evaluator.test_rag(baseline_rag, "2022年，汇丰银行不良贷款率是多少? ", verbose=2)
```

```
-------------------------------------------------------------------------------------------------------
Question: 2022年，汇丰银行不良贷款率是多少?
Response: I'm sorry, there is no information provided in the given context about the non-performing loan ratio of HSBC B


Sources:
page=164, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/hsbc_annual_reports/汇丰银行(中国)有限公司2022年度报告.
page=24, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/hsbc_annual_reports/汇丰银行(中国)有限公司2022年度报告.p
page=105, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/hsbc_annual_reports/汇丰银行(中国)有限公司2022年度报告.
page=71, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/hsbc_annual_reports/汇丰银行(中国)有限公司2022年度报告.p
```

```python
a = evaluator.test_rag(better_emb_rag, "2022年，汇丰银行不良贷款率是多少? ", verbose=2)
```

```
-------------------------------------------------------------------------------------------------------
Question: 2022年，汇丰银行不良贷款率是多少?
Response: 汇丰银行2022年不良贷款率为0.21%。


Sources:
page=23, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/hsbc_annual_reports/汇丰银行(中国)有限公司2022年度报告.p
page=152, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/hsbc_annual_reports/汇丰银行(中国)有限公司2022年度报告.
page=153, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/hsbc_annual_reports/汇丰银行(中国)有限公司2022年度报告.
page=24, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/hsbc_annual_reports/汇丰银行(中国)有限公司2022年度报告.p
```

从上面的结果可以看到，使用更好的Embedding模型，可以提高召回的相关性，从而提升RAG的效果。（Chroma似乎有随机性，结果每次跑有所不同，https://github.com/langchain-ai/langchain/issues/1946）

## b、【done】优化Chunk粒度

Chunk的粒度需要根据语料的粒度，Embedding模型的效果，以及LLM context长度等因素来决定。

- 如果多是FAQ等短的问答对，可以选择较小的chunk size;

- 譬如语录多是短的FAQ，可以使用较短的窗口，如果是横跨多页的篇章（譬如汇丰年报里面，关于"与银行业务相关的主要风险"的部分，横跨了P27～P29三页），可以选择较大的chunk size；
- 为了兼容两种情况，可以使用Langchain里面的ParentDocumentRetriever，先用较小的child chunk来检索，然后返回较大的parent chunk来进行后续的生成；
- 由于LLM一般有context长度的限制，所以并非chunk size越大越好（有可能被截断），一般需要根据应用场景来调整到合适的chunk size。

```
# small chunk
text_splitter = RecursiveCharacterTextSplitter(chunk_size=200, chunk_overlap=50)
texts = text_splitter.split_documents(documents)

small_chunks_vectordb = Chroma.from_documents(
    documents=texts,
    embedding=embedding,
    collection_name="small_chunks"
)

small_chunks_retriever = small_chunks_vectordb.as_retriever()

small_chunks_rag = initialize_rag(llm=chatgpt, retriever=small_chunks_retriever)


# big chunk
text_splitter = RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=50)
texts = text_splitter.split_documents(documents)

big_chunks_vectordb = Chroma.from_documents(
    documents=texts,
    embedding=embedding,
    collection_name="big_chunks"
)

big_chunks_retriever = big_chunks_vectordb.as_retriever()

big_chunks_rag = initialize_rag(llm=chatgpt, retriever=big_chunks_retriever)
```

## c、【done】多级召回保证召回粒度

```
from langchain.retrievers import ParentDocumentRetriever
from langchain.storage import InMemoryStore

# text splitter for big chunks
parent_splitter = RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=50)

# text splitter for small chunks
child_splitter = RecursiveCharacterTextSplitter(chunk_size=200, chunk_overlap=50)

# vectorstore for small chunks
vectorstore = Chroma(collection_name="parent_chunks", embedding_function=embedding)

# storage for big chunks
store = InMemoryStore()

parent_chunks_retriever = ParentDocumentRetriever(
    vectorstore=vectorstore,
    docstore=store,
    child_splitter=child_splitter,
    parent_splitter=parent_splitter,
)

parent_chunks_retriever.add_documents(documents)

len(list(parent_chunks_retriever.docstore.yield_keys()))
```

```
    1459
```

```
parent_chunks_rag = initialize_rag(llm=chatgpt, retriever=parent_chunks_retriever)
```

## d、【wip】多路召回提升召回率

```
# branches = [
#     "上海自贸试验区支行",
#     "深圳华侨城支行",
#     "深圳海天路支行",
#     "深圳华强北路支行",
#     "苏州玄妙广场支行",
#     "天津国际大厦支行",
#     "阳江支行",
# ]
# for branch in branches:
#     query = f"汇丰银行{branch}的地址和电话是多少? "
#     for rag in [baseline_rag, small_chunks_rag, big_chunks_rag, parent_chunks_rag]:
#         a = test_rag(rag=rag, query=query, verbose=1)


# from langchain.retrievers import BM25Retriever, EnsembleRetriever

# bm25_retriever = BM25Retriever.from_documents(texts)
# bm25_retriever.k = 2

# bm25_retriever.get_relevant_documents("阳江支行")

# ensemble_retriever = EnsembleRetriever(
#     retrievers=[bm25_retriever, parent_chunks_retriever],
#     weights=[0.5, 0.5]
# )

# ensemble_rag = initialize_rag(llm=chatgpt, retriever=ensemble_retriever)


# a = test_rag(rag=ensemble_rag, query=query, verbose=1)
```

## e、【wip】MultiQuery提升召回多样性

```
# from langchain.retrievers.multi_query import MultiQueryRetriever

# retriever_from_llm = MultiQueryRetriever.from_llm(
#     retriever=parent_chunks_retriever, llm=chatgpt
# )

# llm_rag = initialize_rag(llm=chatgpt, retriever=retriever_from_llm)

# query = "2022年的营收额是多少? 相比2021增长了多少? "
# a = test_rag(rag=llm_rag, query=query, verbose=1)


# from typing import List
# from langchain.chains import LLMChain
# from pydantic import BaseModel, Field
# from langchain.prompts import PromptTemplate
# from langchain.output_parsers import PydanticOutputParser


# # Output parser will split the LLM result into a list of queries
# class LineList(BaseModel):
#     # "lines" is the key (attribute name) of the parsed output
#     lines: List[str] = Field(description="Lines of text")


# class LineListOutputParser(PydanticOutputParser):
#     def __init__(self) -> None:
#         super().__init__(pydantic_object=LineList)

#     def parse(self, text: str) -> LineList:
#         lines = text.strip().split("\n")
#         return LineList(lines=lines)


# output_parser = LineListOutputParser()

# QUERY_PROMPT = PromptTemplate(
#     input_variables=["question"],
#     template="""You are an AI language model assistant. Your task is to extract the keywords
```

```
#     temptate-    You are an AI tanguage modet assistant. Your task is to extract the keywords
#     in the given user question.
#     Provide these keywords separated by newlines.

#     Examples:
#     Original question: 汇丰银行阳江支行的地址和电话是多少?
#     Output:
#     汇丰银行
#     阳江支行

#     Original question: {question}""",
# )


# # Chain
# llm_chain = LLMChain(llm=chatgpt, prompt=QUERY_PROMPT, output_parser=output_parser)


# llm_retriever = MultiQueryRetriever(
#     retriever=parent_chunks_retriever, llm_chain=llm_chain, parser_key="lines"
# )


# llm_rag = initialize_rag(llm=chatgpt, retriever=llm_retriever)


# query = "汇丰银行阳江支行的地址和电话是多少? "
# a = test_rag(rag=llm_rag, query=query, verbose=1)


# llm_chain(query)
```

## f、【wip】MultiVector提高召回

```
# !apt install tesseract-ocr
# !apt-get install poppler-utils
# !pip install -q pytesseract


# from lxml import html
# from pydantic import BaseModel
# from typing import Any, Optional
# from unstructured.partition.pdf import partition_pdf


# # Get elements
# raw_pdf_elements = partition_pdf(filename=f"{DATA_DIR}/llm_papers/Llama 2- Open Foundation and Fine-Tuned Ch
#                                 # Unstructured first finds embedded image blocks
#                                 extract_images_in_pdf=False,
#                                 # Use layout model (YOLOX) to get bounding boxes (for tables) and find tit]
#                                 # Titles are any sub-section of the document
#                                 infer_table_structure=True,
#                                 # Post processing to aggregate text once we have the title
#                                 chunking_strategy="by_title",
#                                 # Chunking params to aggregate text blocks
#                                 # Attempt to create a new chunk 3800 chars
#                                 # Attempt to keep chunks > 2000 chars
#                                 max_characters=4000,
#                                 new_after_n_chars=3800,
#                                 combine_text_under_n_chars=2000,
#                                 image_output_dir_path=f"{DATA_DIR}/llm_papers")


# class Element(BaseModel):
#     type: str
#     text: Any


# # Categorize by type
# categorized_elements = []
# for element in raw_pdf_elements:
#     if "unstructured.documents.elements.Table" in str(type(element)):
#         categorized_elements.append(Element(type="table", text=str(element)))
#     elif "unstructured.documents.elements.CompositeElement" in str(type(element)):
#         categorized_elements.append(Element(type="text", text=str(element)))


# # Tables
# table_elements = [e for e in categorized_elements if e.type == "table"]
# print(len(table_elements))
```

```
# # Text
# text_elements = [e for e in categorized_elements if e.type == "text"]
# print(len(text_elements))


# from langchain.docstore.document import Document
# from langchain.prompts import ChatPromptTemplate
# from langchain.schema.output_parser import StrOutputParser
# from langchain.chains.summarize import load_summarize_chain

# # Prompt
# prompt_text="""You are an assistant tasked with summarizing tables and text. \
# Give a concise summary of the table or text. Table or text chunk: {text} """
# prompt = ChatPromptTemplate.from_template(prompt_text)

# # Summary chain
# summarize_chain = load_summarize_chain(chatgpt, chain_type="stuff")

# # Apply to tables
# tables = [Document(page_content=i.text) for i in table_elements]
# table_summaries = [summarize_chain.run([table]) for table in tables]

# # Apply to texts
# texts = [Document(page_content=i.text) for i in text_elements]
# text_summaries = [summarize_chain.run([text]) for text in texts]


# len(tables), len(table_summaries), len(texts), len(text_summaries)


# import uuid
# from langchain.vectorstores import Chroma
# from langchain.storage import InMemoryStore
# from langchain.schema.document import Document
# from langchain.embeddings import OpenAIEmbeddings
# from langchain.retrievers.multi_vector import MultiVectorRetriever

# # The vectorstore to use to index the child chunks
# vectorstore = Chroma(
#     collection_name="summaries",
#     embedding_function=embedding
# )

# # The storage layer for the parent documents
# store = InMemoryStore()
# id_key = "doc_id"

# # The retriever (empty to start)
# multi_vector_retriever = MultiVectorRetriever(
#     vectorstore=vectorstore,
#     docstore=store,
#     id_key=id_key,
# )

# # Add texts
# doc_ids = [str(uuid.uuid4()) for _ in texts]
# summary_texts = [Document(page_content=s,metadata={id_key: doc_ids[i]}) for i, s in enumerate(text_summaries
# multi_vector_retriever.vectorstore.add_documents(summary_texts)
# multi_vector_retriever.docstore.mset(list(zip(doc_ids, texts)))

# # Add tables
# table_ids = [str(uuid.uuid4()) for _ in tables]
# summary_tables = [Document(page_content=s,metadata={id_key: table_ids[i]}) for i, s in enumerate(table_summa
# multi_vector_retriever.vectorstore.add_documents(summary_tables)
# multi_vector_retriever.docstore.mset(list(zip(table_ids, tables)))


# multi_vector_rag = initialize_rag(llm=chatgpt, retriever=multi_vector_retriever)


# query = "What is the number of training tokens for Llama 2?"
# for rag in [baseline_rag, big_chunks_rag, multi_vector_rag]:
#     test_rag(rag, query, verbose=1)
```

▼ g、【done】MultiEmbedding

考虑到LLM Paper相关的数据是英文的，所以这里我们采用英文Embedding模型对其进行编码能得到更好的效果。参考MTEB，我们采用了
BAAI/bge-large-en-v1.5（https://github.com/FlagOpen/FlagEmbedding）。

```
model_name = "BAAI/bge-large-en-v1.5"
model_kwargs = {'device': 'cuda'}
encode_kwargs = {'normalize_embeddings': True}
en_embedding = HuggingFaceBgeEmbeddings(
    model_name=model_name,
    model_kwargs=model_kwargs,
    encode_kwargs=encode_kwargs,
    query_instruction="Represent this sentence for searching relevant passages: "
)
```

| Downloading (…)5e2c6/.gitattributes: 100% | 1.52k/1.52k [00:00<00:00, 145kB/s] |
| Downloading (…)_Pooling/config.json: 100% | 191/191 [00:00<00:00, 16.6kB/s] |
| Downloading (…)ba76d5e2c6/README.md: 100% | 90.3k/90.3k [00:00<00:00, 1.11MB/s] |
| Downloading (…)76d5e2c6/config.json: 100% | 779/779 [00:00<00:00, 60.8kB/s] |
| Downloading (…)ce_transformers.json: 100% | 124/124 [00:00<00:00, 11.1kB/s] |
| Downloading pytorch_model.bin: 100% | 1.34G/1.34G [00:04<00:00, 237MB/s] |
| Downloading | 52.0/52.0 [00:00<00:00, |

```
text_splitter = RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=50)
texts = text_splitter.split_documents(papers)

papers_en_emb_vectordb = Chroma.from_documents(
    documents=texts,
    embedding=en_embedding,
    collection_name="papers_en_emb"
)

papers_en_emb_retriever = papers_en_emb_vectordb.as_retriever()

papers_en_emb_rag = initialize_rag(llm=chatgpt, retriever=papers_en_emb_retriever)

a = evaluator.test_rag(baseline_rag, query="What is the number of training tokens for LLaMA2?", verbose=2)
```

```
    --------------------------------------------------------------------------------------------------------------
    Question: What is the number of training tokens for LLaMA2?
    Response: The number of training tokens for Llama 2 is not explicitly stated in the given context.

    Sources:
    page=1, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/LLaMA- Open and Efficient Foundation Lang
    page=15, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned
    page=5, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned C
    page=1, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/LLaMA- Open and Efficient Foundation Lang
```

```
a = evaluator.test_rag(papers_en_emb_rag, query="What is the number of training tokens for LLaMA2?", verbose=2
```

```
    --------------------------------------------------------------------------------------------------------------
    Question: What is the number of training tokens for LLaMA2?
    Response: The number of training tokens for Llama 2 is 2.0 trillion tokens. This is stated in the table comparing the at

    Sources:
    page=5, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned C
    page=4, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned C
    page=4, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned C
    page=0, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/LLaMA- Open and Efficient Foundation Lang
```

▼ h、【done】Rerank提升召回相关性

使用Rerank的原因：

- 虽然Embedding+向量召回可以很快在大量的文档里面实现召回，但是返回的召回不一定都相关；
- 基于HNSW的召回具有一定的随机性，多次召回结果可能会不一致；
- LLM有context长度限制，如果相关的内容排序不是特别靠前，有可能被截断，影响提取信息的准确度；

可以引入Reranker提升相关性和稳定性，同时减轻输入LLM的context长度限制。

- Reranker：BAAI/bge-reranker-large, https://github.com/FlagOpen/FlagEmbedding
- 参考文章：https://mp.weixin.qq.com/s/4UoRi8VhQjfE7zcpFnre4A

```
from FlagEmbedding import FlagReranker
reranker = FlagReranker('BAAI/bge-reranker-large', use_fp16=True) # Setting use_fp16 to True speeds up computa

score = reranker.compute_score(['query', 'passage'])
print(score)

scores = reranker.compute_score([['what is panda?', 'hi'], ['what is panda?', 'The giant panda (Ailuropoda mel
print(scores)
```

| Downloading | 443/443 [00:00<00:00, |
| (…)okenizer_config.json: 100% | 42.2kB/s] |
| Downloading | 5.07M/5.07M [00:00<00:00, |
| (…)tencepiece.bpe.model: 100% | 22.1MB/s] |
| Downloading | 279/279 [00:00<00:00, |
| (…)cial_tokens_map.json: 100% | 26.0kB/s] |
| Downloading (…)lve/main/config.json: | 801/801 [00:00<00:00, |

```
import asyncio
from typing import Dict, Optional, Sequence

from langchain.callbacks.manager import Callbacks
from langchain.pydantic_v1 import Extra
from langchain.retrievers import ContextualCompressionRetriever
from langchain.retrievers.document_compressors.base import BaseDocumentCompressor
from langchain.schema import Document


class BGERerank(BaseDocumentCompressor):
    """Document compressor that uses BGE Reranker."""

    reranker: FlagReranker

    class Config:
        """Configuration for this pydantic object."""

        extra = Extra.forbid
        arbitrary_types_allowed = True

    def compress_documents(
        self,
        documents: Sequence[Document],
        query: str,
        callbacks: Optional[Callbacks] = None,
    ) -> Sequence[Document]:
        if len(documents) == 0:  # to avoid empty api call
            return []
        doc_list = list(documents)
        query_doc_pairs = [[query, d.page_content] for d in doc_list]
        scores = self.reranker.compute_score(query_doc_pairs)
        for doc,score in zip(doc_list, scores):
            doc.metadata["relevance_score"] = score
        doc_list = sorted(doc_list, key=lambda d: d.metadata["relevance_score"], reverse=True)
        return doc_list

    async def acompress_documents(
        self,
        documents: Sequence[Document],
        query: str,
        callbacks: Optional[Callbacks] = None,
    ) -> Sequence[Document]:
        """Compress retrieved documents given the query context."""
```

```
            return await asyncio.get_running_loop().run_in_executor(
                None, self.compress_documents, documents, query, callbacks
            )


compressor = BGERerank(reranker=reranker)

# 更好的策略是先用base_retriever召回一个比较大的池子，譬如10，然后使用base_compressor重排序选出topk，譬如3
# 这里为了方便，直接使用base_retriever的topk大小作为base_compressor的大小
big_chunks_compression_retriever = ContextualCompressionRetriever(
    base_compressor=compressor, base_retriever=big_chunks_retriever
)


big_chunks_compression_rag = initialize_rag(llm=chatgpt, retriever=big_chunks_compression_retriever)


a = evaluator.test_rag(big_chunks_rag, query="如何环球转账？", verbose=3)
```

Question: 免责声明
Answer: 以上常见问题旨在提供于中国大陆开展金融活动所可能面临的与金融产品、服务或监管要求相关的概略信息。因金融产品/服务及相关监管要求的多样性和
page=None, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/hsbc_faqs/我想要进行国内外转账及汇款.txt
400-820-3090
若您在海外或港澳台地区，请在号码前加拨中国地区码+86

Question: 10. 我是否可在中国境内划转外币？
Answer: 可以。但根据国家外汇管理局相关规定，个人外汇账户内资金境内划转，仅限于本人账户之间、个人与近亲属账户之间。若有外汇境内划转需要，需要您

Question: 11. 通过个人网上银行，如何才能连接不同国家的汇丰账户？
Answer: 首先，将被连接的国家/地区的网上银行账号必须是活跃的。在您选择要添加的国家/地区后，根据屏幕提示输入网上银行用户名、密码、安全密码/安全问

Question: 12. 什么是环球网上银行？我能否由此浏览我在世界各地的汇丰账户？
Answer: "环球网上银行"服务是专门为我客户推出的网上银行增值服务。我行客户可以通过"环球网上银行"服务浏览其拥有的其他国家或地区的汇丰账户信息。

Question: 13. 环球转账
Answer: 环球转账支持币种包括：美元，英镑，欧元，日元，港元，加拿大元，澳元，新加坡元，瑞士法郎，新西兰元。
环球转账仅支持外币现汇储蓄账户，不支持外币现钞储蓄账户。
由于环球转账不能自动加开附属外币账号，您可在网上银行/柜台开立外币附属账号。

Question: 14. 汇丰环球转账的服务时间是怎样的？
Answer: 提供 7*24 小时全天服务，但不包括网上银行系统维护时间（通常为周末凌晨）。 您可以在网上银行登录页面查看具体安排。

Question: 15. 我可以在哪里使用环球转账？
Answer: 通过网上银行和手机银行，您可向您在其他国家/地区所持有的汇丰同名账户进行外币跨境资金划转，目前手机银行率先提供支持客户使用"环球转账亲情

Question: 16. 我是否能进行同币种或跨币种的外币环球转账？
Answer: 可以，两者支持。对于跨外汇币种而言,环球转账检视详情页面会显示当时参考汇率供客户确认，实际执行汇率将以客户确认的汇率为准并显示在交易完成
page=None, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/hsbc_faqs/我想要进行国内外转账及汇款.txt
### 国内外转账及汇款

Question: 我如何使用汇丰手机银行进行人民币转账？
Answer: 汇丰中国手机银行支持实时、普通、稍后、循环四种转账方式。日累计不超过5万元的人民币非同名转账，可实现免密码器实时到账！

手机银行首页选择"人民币转账"，点击"银行账户转账"；
选择"转出账户"，输入"转账金额"，选择"已保存的我的收款人"和"转账方式"后继续；
输入手机银行登入密码和短信验证码。确认交易信息后，即完成转账。

Question: 如何添加"我的收款人"？
Answer: 如果您没有已经保存的收款人，请填写或拍照扫描"收款账号"，选择"收款人开户银行"和填写"收款人姓名"，根据提示使用安全密码器生成交易安全码并

Question: 我想进行环球转账，如何操作？
Answer: 汇丰环球转账功能，可通过登录网上银行


在"我的银行" –"转账及货币兑换" –"转账及货币兑换"。
在新交易中，"转出账户"选择外币储蓄账户。转入账户选择"我的账户"，并选择另一个汇丰其他国家 / 地区的外币账户。
您可验证交易详情并确认，即可完成交易。


您还可以通过登录"汇丰银行"手机App进行相关操作。

Question: 我如何使用外币境外第三方转账？
Answer: 登录汇丰银行网上银行，前往"我的银行"–"转账及货币兑换"。


新交易中，"转出账户"选择外币，"转入账户"选择"新的收款人"，输入"收款人银行详细信息"。
添加"收款人个人详细信息"。

```
a = evaluator.test_rag(big_chunks_compression_rag, query="如何环球转账？", verbose=3)
```

------------------------------------------------------------------------------------------------------
Question: 如何环球转账？
Response: 如何环球转账？您可以通过登录网上银行或手机银行进行操作。在网上银行中，选择"我的银行"–"转账及货币兑换"–"转账及货币兑换"，选择转出账户


Sources:
page=None, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/hsbc_faqs/我想要进行国内外转账及汇款.txt
400-820-3090

若您在海外或港澳台地区，请在号码前加拨中国地区码+86

Question: 10. 我是否可在中国境内划转外币？
Answer: 可以。但根据国家外汇管理局相关规定，个人外汇账户内资金境内划转，仅限于本人账户之间、个人与近亲属账户之间。若有外汇境内划转需要，需要您?

Question: 11. 通过个人网上银行，如何才能连接不同国家的汇丰账户？
Answer: 首先,将被连接的国家/地区的网上银行账号必须是活跃的。在您选择要添加的国家/地区后，根据屏幕提示输入网上银行用户名、密码、安全密码/安全问

Question: 12. 什么是环球网上银行？我能否由此浏览我在世界各地的汇丰账户？
Answer: "环球网上银行"服务是专门为我行客户推出的网上银行增值服务。我行客户可以通过"环球网上银行"服务浏览其拥有的其他国家或地区的汇丰账户信息。

Question: 13. 环球转账
Answer: 环球转账支持币种包括：美元，英镑，欧元，日元，港元，加拿大元，澳元，新加坡元，瑞士法郎，新西兰元。
环球转账仅支持外币现汇储蓄账户，不支持外币现钞储蓄账户。
由于环球转账不能自动加开附属外币账号，您可在网上银行/柜台开立外币附属账户。

Question: 14. 汇丰环球转账的服务时间是怎样的？
Answer: 提供 7*24 小时全天服务，但不包括网上银行系统维护时间（通常为周末凌晨）。 您可以在网上银行登录页面查看具体安排。

Question: 15. 我可以在哪里使用环球转账？
Answer: 通过网上银行和手机银行，您可向您在其他国家/地区所持有的汇丰同名账户进行外币跨境资金划转，目前手机银行率先提供支持客户使用"环球转账亲情?

Question: 16. 我是否能进行同币种或跨币种的外币环球转账？
Answer: 可以，两者支持。对于跨外汇币种而言，环球转账检视详情页面会显示当时参考汇率供客户确认，实际执行汇率将以客户确认的汇率为准并显示在交易完成
page=None, source=/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/hsbc_faqs/我想要进行国内外转账及汇款.txt
### 国内外转账及汇款

Question: 我如何使用汇丰手机银行进行人民币转账？
Answer: 汇丰中国手机银行支持实时、普通、稍后、循环四种转账方式。日累计不超过5万元的人民币非同名转账，可实现免密码器实时到账！

手机银行首页选择"人民币转账"，点击"银行账户转账"；
选择"转出账户"，输入"转账金额"，选择"已保存的我的收款人"和"转账方式"后继续；
输入手机银行登入密码和短信验证码。确认交易信息后，即完成转账。

Question: 如何添加"我的收款人"？
Answer: 如果您没有已经保存的收款人，请填写或拍照扫描"收款账号"，选择"收款人开户银行"和填写"收款人姓名"，根据提示使用安全密码器生成交易安全码并

Question: 我想进行环球转账，如何操作？
Answer: 汇丰环球转账功能，可通过登录网上银行

在"我的银行" –"转账及货币兑换" –"转账及货币兑换"。
在新交易中，"转出账户"选择外币储蓄账户。转入账户选择"我的账户"，并选择另一个汇丰其他国家 / 地区的外币账户。
您可验证交易详情并确认，即可完成交易。

您还可以通过登录"汇丰银行"手机App进行相关操作。

Question: 我如何使用外币境外第三方转账？
Answer: 登录汇丰银行网上银行、前往"我的银行"–"转账及货币兑换"。

可以看到经过rerank之后，"环球转账"相关的内容从第3位上升到第2位，从而可以被送入LLM进行生成（LLM有context的限制，相关的内容排序越靠前，越有利于LLM能提取到准确的信息）。

# 3、生成层

## a、【done】根据场景选择更好的LLM

好的LLM模型本身就很大程度决定了生成质量的基础水平；（本Demo没有进行验证，直接使用ChatGPT作为基础模型）

## b、【done】优化Prompt

Prompt Engineering是影响LLM生成质量很重要的环节，其影响着LLM指令跟随的能力，同时在Prompt中加入额外的信息，也能帮助LLM生成更好和更具有事实性的回复。

```
fin_rag_prompt = deepcopy(default_rag_prompt)
fin_rag_prompt.messages[0].prompt.template = """You are an experienced financial analyst for HSBC with an inte
Question: {question}
Context: {context}
Answer:"""


prompt_dict = {
    "default_rag_prompt": default_rag_prompt,
    "fin_rag_prompt": fin_rag_prompt,
    "default_prompt": None,
}
retriever_dict = {
    "baseline_retriever": baseline_retriever,
```

```
        "better_pdf_retriever": better_pdf_retriever,
        "better_emb_retriever": better_emb_retriever,
        "small_chunks_retriever": small_chunks_retriever,
        "big_chunks_retriever": big_chunks_retriever,
        "parent_chunks_retriever": parent_chunks_retriever,
        "big_chunks_compression_retriever": big_chunks_compression_retriever,
}


rag_results = []
for prompt_name, prompt in prompt_dict.items():
    for retriever_name, retriever in retriever_dict.items():
        rag = initialize_rag(llm=chatgpt, retriever=retriever, prompt=prompt)
        res = evaluator.test_rag_all(rag=rag, question_answer_pairs=question_answer_pairs, verbose=0)
        result = {
            "prompt": prompt_name,
            "retriever": retriever_name,
            "results": res
        }
        rag_results.append(result)


for result in rag_results:
    result["results"] = evaluator.evaluate_rag_results(result["results"])
    result["score"] = evaluator.compute_rag_score(result["results"])
```

Downloading builder script:                                     6.27k/6.27k [00:00<00:00,
100%                                                            464kB/s]

```
df = pd.DataFrame(rag_results).sort_values("score", ascending=False)


df = df[["prompt","retriever","score"]]
df["rag"] = "single_rag"
df = df[["prompt","retriever","rag","score"]]


best_prompt = df.iloc[0]["prompt"]
best_retriever = df.iloc[0]["retriever"]
best_score = df.iloc[0]["score"]


best_prompt, best_retriever, best_score
```

    ('default_prompt', 'big_chunks_compression_retriever', 0.772767857142857)

## 4、整体实现

使用LLM进行意图分类，然后分别调用领域专家RAG，类似于Mixture of Expert (MoE)架构。

```
class EnsembleRAG:
    def __init__(self, intent_prompt, tools, llm):
        self.intent_prompt = intent_prompt
        self.tools = tools
        self.llm = llm

    def __call__(self, query):
        intent_prompt = self.intent_prompt.format(question=query)
        category = self.llm.predict(intent_prompt)
        category = category.lower()
        if category == "none" or category not in self.tools:
            print(f"{query}: {category}")
            return {"query": "query", "result": "I don't know."}
        else:
            tool = self.tools[category]
            return tool(query)


# 使用llm进行意图分类，然后分别调用对应的RAG
intent_domain_prompt = """你是一个有用的助手。
以下是用户可能提出问题的四个意图领域的描述。
对于给定的用户问题，请在这些意图领域中进行选择。
请仅返回意图领域。回答后不要返回任何其他内容。
如果您认为没有任何与之相关的领域，请返回 NONE。
```

域名: huifeng_faq
描述: 对于回答与汇丰金融产品相关的问题非常有用，如账户、账单、转账、汇款、支付、数字银行、微信服务、存款、住房抵押贷款、投资、保险、

域名: huifeng_annual_report
描述: 对于回答与汇丰财报和年度报告相关的问题非常有用，例如营业收入、营业支出、成本、吸收存款、公司结构、部门结构、公司治理、薪酬结构

域名: jinronghuiyi_article
说明: 适用于回答中央金融工作会议的问题，比如不同证券公司的观点。

域名: llm_paper
描述: 对于回答与大型语言模型相关的问题非常有用，如LLM、LLaMA、LLaMA2和金融LLM（FinGPT和FinLLM）等。

Question: {question}
Domain:
"""

```python
best_rag = initialize_rag(llm=chatgpt, retriever=retriever_dict[best_retriever], prompt=prompt_dict[best_prompt
paper_rag = initialize_rag(llm=chatgpt, retriever=papers_en_emb_retriever, prompt=prompt_dict[best_prompt])
tools = {
    "huifeng_faq": best_rag, # 优化了pdf+zh embedding+chunk+rerank+prompt
    "huifeng_annual_report": best_rag, # 优化了pdf+zh embedding+chunk+rerank+prompt
    "jinronghuiyi_article": metadata_rag, # 优化了metadata
    "llm_paper": paper_rag # 优化了pdf+en embedding+chunk+prompt
}
```

```python
intent_results = {}
for category,qas in question_answer_pairs.items():
    intent_results[category] = 0
    for qa in qas:
        i = chatgpt.predict(intent_domain_prompt.format(question=qa["question"]))
        intent_results[category] += int(i == category)
    intent_results[category] /= len(qas)
intent_results
```

```
{'huifeng_faq': 1.0,
 'huifeng_annual_report': 0.7647058823529411,
 'jinronghuiyi_article': 1.0,
 'llm_paper': 1.0}
```

```python
ensemble_rag = EnsembleRAG(intent_prompt=intent_domain_prompt, tools=tools, llm=chatgpt)
ensemble_rag_results = evaluator.test_rag_all(rag=ensemble_rag, question_answer_pairs=question_answer_pairs,
ensemble_rag_results = evaluator.evaluate_rag_results(ensemble_rag_results)
evaluator.inspect_rag_results(ensemble_rag_results)
```

category=huifeng_faq, question=如何环球转账？, score=0.4, response=您可以通过以下步骤进行汇丰环球转账：

1. 登录汇丰银行网上银行，进入"我的银行"–"转账及货币兑换"–"转账及货币兑换"页面。

2. 在新交易中，"转出账户"选择外币储蓄账户。转入账户选择"我的账户"，并选择另一个汇丰其他国家/地区的外币账户。

3. 输入转账金额和币种，选择交易编码/原因，然后验证交易详情并确认，即可完成交易。

您还可以通过登录汇丰银行手机App进行相关操作。请注意，环球转账支持的币种包括美元、英镑、欧元、日元、港元、加拿大元、澳元、新加坡元、瑞士法郎和新
category=jinronghuiyi_article, question=广发证券对中央金融工作会议的观点有哪些？, score=0.0, response=广发证券对中央金融工作会议的观
category=llm_paper, question=What is the author of Llama?, score=0.9285714285714286, response=The authors of Llama are H
category=llm_paper, question=What is the affiliation of the first author of DISC-FinLLM?, score=0.0, response=The first
category=llm_paper, question=What are the common authors of Llama and Llama 2?, score=0.0, response=Unfortunately, the p
category=llm_paper, question=Is the first author of Llama and Llama 2 the same? If yes, please output <<YES>>, otherwise
category=llm_paper, question=Is there any LLM in financial area? What are they?, score=0.4, response=Yes, there are seve

```python
score = evaluator.compute_rag_score(ensemble_rag_results)
score
```

```
0.835267857142857
```

```python
df.loc[len(df.index)] = [f"{best_prompt}", f"{best_retriever}+SelfQueryRetriever", "ensemble_rag", score]
```

## ▾ 5、进一步改进

## ▾ a、【done】引入KG结构化知识

```python
paper_rag(question_answer_pairs["llm_paper"][5]["question"])
```

```
{'query': 'Is the first author of Llama and Llama 2 the same? If yes, please output <<YES>>, otherwise output <<NO>>.',
 'result': 'No.',
 'source_documents': [Document(page_content='Llama 2: Open Foundation and Fine-Tuned Chat Models\nHugo Touvron*\nLouis
Martin†\nKevin Stone†\nPeter Albert Amjad Almahairi Yasmine Babaei Nikolay Bashlykov Soumya Batra\nPrajjwal Bhargava
Shruti Bhosale Dan Bikel Lukas Blecher Cristian Canton Ferrer Moya Chen\nGuillem Cucurull David Esiobu Jude Fernandes
Jeremy Fu Wenyin Fu Brian Fuller\nCynthia Gao Vedanuj Goswami Naman Goyal Anthony Hartshorn Saghar Hosseini Rui
Hou\nHakan Inan Marcin Kardas Viktor Kerkez Madian Khabsa Isabel Kloumann Artem Korenev\nPunit Singh Koura Marie-Anne
Lachaux Thibaut Lavril Jenya Lee Diana Liskovich\nYinghai Lu Yuning Mao Xavier Martinet Todor Mihaylov Pushkar
Mishra\nIgor Molybog Yixin Nie Andrew Poulton Jeremy Reizenstein Rashi Rungta Kalyan Saladi\nAlan Schelten Ruan Silva
Eric Michael Smith Ranjan Subramanian Xiaoqing Ellen Tan Binh Tang\nRoss Taylor Adina Williams Jian Xiang Kuan Puxin Xu
Zheng Yan Iliyan Zarov Yuchen Zhang\nAngela Fan Melanie Kambadur Sharan Narang Aurelien Rodriguez Robert
Stojnic\nSergey Edunov', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref',
'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned
Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 0, 'producer': 'pdfTeX-
1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-
Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''}),
  Document(page_content='Figure 3: Safety human evaluation results for Llama 2-Chat compared to other open-source and
closed-\nsource models. Human raters judged model generations for safety violations across ~2,000 adversarial\nprompts
consisting of both single and multi-turn prompts. More details can be found in Section 4.4. It is\nimportant to caveat
these safety results with the inherent bias of LLM evaluations due to limitations of the\nprompt set, subjectivity of
the review guidelines, and subjectivity of individual raters. Additionally, these\nsafety evaluations are performed
using content standards that are likely to be biased towards the Llama\n2-Chat models.\nWe are releasing the following
models to the general public for research and commercial use‡:\n1. Llama 2, an updated version of Llama 1, trained on a
new mix of publicly available data. We also\nincreased the size of the pretraining corpus by 40%, doubled the context
length of the model, and', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with
hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and
Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 3,
'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open
Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''}),
  Document(page_content='• Armand Joulin, Edouard Grave, Guillaume Lample, and Timothee Lacroix, members of the
original\nLlama team who helped get this work started.\n• Drew Hamlin, Chantal Mora, and Aran Mun, who gave us some
design input on the figures in the\npaper.\n• Vijai Mohan for the discussions about RLHF that inspired our Figure 20,
and his contribution to the\ninternal demo.\n• Early reviewers of this paper, who helped us improve its quality,
including Mike Lewis, Joelle Pineau,\nLaurens van der Maaten, Jason Weston, and Omer Levy.\nA.2\nAdditional Details for
Pretraining\nA.2.1\nArchitecture Changes Compared to Llama 1\nContext Length.\nWe expand the context window for Llama 2
from 2048 tokens to 4096 tokens. The longer\ncontext window enables models to process more information, which is
particularly useful for supporting\nlonger histories in chat applications, various summarization tasks, and
understanding longer documents.', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with
hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and
Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 46,
'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open
Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''}),
  Document(page_content='Sergey Edunov\nThomas Scialom*\nGenAI, Meta\nAbstract\nIn this work, we develop and release
Llama 2, a collection of pretrained and fine-tuned\nlarge language models (LLMs) ranging in scale from 7 billion to 70
billion parameters.\nOur fine-tuned LLMs, called Llama 2-Chat, are optimized for dialogue use cases. Our\nmodels
outperform open-source chat models on most benchmarks we tested, and based on\nour human evaluations for helpfulness
and safety, may be a suitable substitute for closed-\nsource models. We provide a detailed description of our approach
to fine-tuning and safety\nimprovements of Llama 2-Chat in order to enable the community to build on our work
and\ncontribute to the responsible development of LLMs.\n*Equal contribution, corresponding authors: {tscialom,
htouvron}@meta.com\n†Second author\nContributions for all the authors can be found in Section A.1.\narXiv:2307.09288v2
[cs.CL]  19 Jul 2023', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref',
'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned
Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 0, 'producer': 'pdfTeX-
1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-
Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''})]}
```

我们通过在Prompt中加入KG的信息来模拟LLM+KG。更具体的方案可以参考:

- https://blog.langchain.dev/using-a-knowledge-graph-to-implement-a-devops-rag-application/
- https://mp.weixin.qq.com/s/VJRG0MUaEGR6iM_xFRroyg

```
paper_rag.combine_documents_chain.llm_chain.prompt = deepcopy(paper_rag.combine_documents_chain.llm_chain.pro
paper_rag.combine_documents_chain.llm_chain.prompt.messages[0].prompt.template = """Use the following pieces
If you don't know the answer, just say that you don't know, don't try to make up an answer.
You are also provided the following Academic Knowledge Graph for checking affiliation. Please only use it whe
----------------
KG:
("Hugo Touvron", "affiliated_with", "Meta AI")
("Wei Chen", "affiliated_with", "Fudan University and Huazhong University of Science and Technology")
----------------
{context}"""

paper_rag.combine_documents_chain.llm_chain.prompt.messages[1].prompt.template = """{question}"""


paper_rag(question_answer_pairs["llm_paper"][2]["question"])

    {'query': 'What is the affiliation of the first author of Llama?',
     'result': 'The first author of Llama is Hugo Touvron and he is affiliated with Meta AI.',
     'source_documents': [Document(page_content='Llama 2: Open Foundation and Fine-Tuned Chat Models\nHugo Touvron*\nLouis
    Martin†\nKevin Stone†\nPeter Albert Amjad Almahairi Yasmine Babaei Nikolay Bashlykov Soumya Batra\nPrajjwal Bhargava
    Shruti Bhosale Dan Bikel Lukas Blecher Cristian Canton Ferrer Moya Chen\nGuillem Cucurull David Esiobu Jude Fernandes
    Jeremy Fu Wenyin Fu Brian Fuller\nCynthia Gao Vedanuj Goswami Naman Goyal Anthony Hartshorn Saghar Hosseini Rui
    Hou\nHakan Inan Marcin Kardas Viktor Kerkez Madian Khabsa Isabel Kloumann Artem Korenev\nPunit Singh Koura Marie-Anne
    Lachaux Thibaut Lavril Jenya Lee Diana Liskovich\nYinghai Lu Yuning Mao Xavier Martinet Todor Mihaylov Pushkar
```

Mishra\nIgor Molybog Yixin Nie Andrew Poulton Jeremy Reizenstein Rashi Rungta Kalyan Saladi\nAlan Schelten Ruan Silva Eric Michael Smith Ranjan Subramanian Xiaoqing Ellen Tan Binh Tang\nRoss Taylor Adina Williams Jian Xiang Kuan Puxin Xu Zheng Yan Iliyan Zarov Yuchen Zhang\nAngela Fan Melanie Kambadur Sharan Narang Aurelien Rodriguez Robert Stojnic\nSergey Edunov', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 0, 'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''}),
  Document(page_content='LLaMA: Open and Efficient Foundation Language Models\nHugo Touvron∗, Thibaut Lavril∗, Gautier Izacard∗, Xavier Martinet\nMarie-Anne Lachaux, Timothee Lacroix, Baptiste Rozière, Naman Goyal\nEric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin\nEdouard Grave∗, Guillaume Lample∗\nMeta AI\nAbstract\nWe introduce LLaMA, a collection of founda-\ntion language models ranging from 7B to 65B\nparameters. We train our models on trillions\nof tokens, and show that it is possible to train\nstate-of-the-art models using publicly avail-\nable datasets exclusively, without resorting\nto proprietary and inaccessible datasets.\nIn\nparticular, LLaMA-13B outperforms GPT-3\n(175B) on most benchmarks, and LLaMA-\n65B is competitive with the best models,\nChinchilla-70B and PaLM-540B. We release\nall our models to the research community1.\n1\nIntroduction\nLarge Languages Models (LLMs) trained on mas-\nsive corpora of texts have shown their ability to per-\nform new tasks from textual instructions or from a', metadata={'author': '', 'creationDate': 'D:20230228015746Z', 'creator': 'LaTeX with hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/LLaMA- Open and Efficient Foundation Language Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230228015746Z', 'page': 0, 'producer': 'pdfTeX-1.40.21', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/LLaMA- Open and Efficient Foundation Language Models.pdf', 'subject': '', 'title': '', 'total_pages': 27, 'trapped': ''}),
  Document(page_content='Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question answering\nchallenge targeting commonsense knowledge. arXiv preprint arXiv:1811.00937, 2018.\nRohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and\nTatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/\ntatsu-lab/stanford_alpaca, 2023.\nRoss Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew\nPoulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for science. arXiv preprint\narXiv:2211.09085, 2022.\n43', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 42, 'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''}),
  Document(page_content='• Armand Joulin, Edouard Grave, Guillaume Lample, and Timothee Lacroix, members of the original\nLlama team who helped get this work started.\n• Drew Hamlin, Chantal Mora, and Aran Mun, who gave us some design input on the figures in the\npaper.\n• Vijai Mohan for the discussions about RLHF that inspired our Figure 20, and his contribution to the\ninternal demo.\n• Early reviewers of this paper, who helped us improve its quality, including Mike Lewis, Joelle Pineau,\nLaurens van der Maaten, Jason Weston, and Omer Levy.\nA.2\nAdditional Details for Pretraining\nA.2.1\nArchitecture Changes Compared to Llama 1\nContext Length.\nWe expand the context window for Llama 2 from 2048 tokens to 4096 tokens. The longer\ncontext window enables models to process more information, which is particularly useful for supporting\nlonger histories in chat applications, various summarization tasks, and understanding longer documents.', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 46, 'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''})]}

```python
ensemble_rag = EnsembleRAG(intent_prompt=intent_domain_prompt, tools=tools, llm=chatgpt)
ensemble_rag_results = evaluator.test_rag_all(rag=ensemble_rag, question_answer_pairs=question_answer_pairs, \
ensemble_rag_results = evaluator.evaluate_rag_results(ensemble_rag_results)
evaluator.inspect_rag_results(ensemble_rag_results)
```

category=huifeng_faq, question=如何环球转账？, score=0.4, response=您可以通过以下步骤进行汇丰环球转账：

1. 登录汇丰银行网上银行，进入"我的银行"-"转账及货币兑换"-"转账及货币兑换"页面。

2. 在新交易中，"转出账户"选择外币储蓄账户。转入账户选择"我的账户"，并选择另一个汇丰其他国家/地区的外币账户。

3. 输入转账金额和币种，选择交易编码/原因，然后验证交易详情并确认，即可完成交易。

您还可以通过登录汇丰银行手机App进行相关操作。请注意，环球转账支持的币种包括美元、英镑、欧元、日元、港元、加拿大元、澳元、新加坡元、瑞士法郎和新
category=jinronghuiyi_article, question=广发证券对中央金融工作会议的观点有哪些？, score=0.0, response=广发证券对中央金融工作会议的观
category=llm_paper, question=What is the author of Llama?, score=0.7857142857142857, response=There are multiple authors
category=llm_paper, question=What are the common authors of Llama and Llama 2?, score=0.16666666666666666, response=Unfo
category=llm_paper, question=Is the first author of Llama and Llama 2 the same? If yes, please output <<YES>>, otherwise
category=llm_paper, question=Is there any LLM in financial area? What are they?, score=0.4, response=Yes, there are seve

```python
score = evaluator.compute_rag_score(ensemble_rag_results)
score
```

0.8672619047619047

```python
df.loc[len(df.index)] = [f"{best_prompt}", f"{best_retriever}+SelfQueryRetriever", "ensemble_rag+KG", score]
```

## ▾ b、【done】Chain-of-Thought处理复杂Query

```python
paper_rag(question_answer_pairs["llm_paper"][5]["question"])
```

{'query': 'Is the first author of Llama and Llama 2 the same? If yes, please output <<YES>>, otherwise output <<NO>>.',
 'result': 'NO',
 'source_documents': [Document(page_content='Llama 2: Open Foundation and Fine-Tuned Chat Models\nHugo Touvron∗\nLouis Martin†\nKevin Stone†\nPeter Albert Amjad Almahairi Yasmine Babaei Nikolay Bashlykov Soumya Batra\nPrajjwal Bhargava Shruti Bhosale Dan Bikel Lukas Blecher Cristian Canton Ferrer Moya Chen\nGuillem Cucurull David Esiobu Jude Fernandes

Jeremy Fu Wenyin Fu Brian Fuller\nCynthia Gao Vedanuj Goswami Naman Goyal Anthony Hartshorn Saghar Hosseini Rui Hou\nHakan Inan Marcin Kardas Viktor Kerkez Madian Khabsa Isabel Kloumann Artem Korenev\nPunit Singh Koura Marie-Anne Lachaux Thibaut Lavril Jenya Lee Diana Liskovich\nYinghai Lu Yuning Mao Xavier Martinet Todor Mihaylov Pushkar Mishra\nIgor Molybog Yixin Nie Andrew Poulton Jeremy Reizenstein Rashi Rungta Kalyan Saladi\nAlan Schelten Ruan Silva Eric Michael Smith Ranjan Subramanian Xiaoqing Ellen Tan Binh Tang\nRoss Taylor Adina Williams Jian Xiang Kuan Puxin Xu Zheng Yan Iliyan Zarov Yuchen Zhang\nAngela Fan Melanie Kambadur Sharan Narang Aurelien Rodriguez Robert Stojnic\nSergey Edunov', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 0, 'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''}),
  Document(page_content='Figure 3: Safety human evaluation results for Llama 2-Chat compared to other open-source and closed-\nsource models. Human raters judged model generations for safety violations across ~2,000 adversarial\nprompts consisting of both single and multi-turn prompts. More details can be found in Section 4.4. It is\nimportant to caveat these safety results with the inherent bias of LLM evaluations due to limitations of the\nprompt set, subjectivity of the review guidelines, and subjectivity of individual raters. Additionally, these\nsafety evaluations are performed using content standards that are likely to be biased towards the Llama\n2-Chat models.\nWe are releasing the following models to the general public for research and commercial use‡:\n1. Llama 2, an updated version of Llama 1, trained on a new mix of publicly available data. We also\nincreased the size of the pretraining corpus by 40%, doubled the context length of the model, and', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 3, 'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''}),
  Document(page_content='• Armand Joulin, Edouard Grave, Guillaume Lample, and Timothee Lacroix, members of the original\nLlama team who helped get this work started.\n• Drew Hamlin, Chantal Mora, and Aran Mun, who gave us some design input on the figures in the\npaper.\n• Vijai Mohan for the discussions about RLHF that inspired our Figure 20, and his contribution to the\ninternal demo.\n• Early reviewers of this paper, who helped us improve its quality, including Mike Lewis, Joelle Pineau,\nLaurens van der Maaten, Jason Weston, and Omer Levy.\nA.2\nAdditional Details for Pretraining\nA.2.1\nArchitecture Changes Compared to Llama 1\nContext Length.\nWe expand the context window for Llama 2 from 2048 tokens to 4096 tokens. The longer\ncontext window enables models to process more information, which is particularly useful for supporting\nlonger histories in chat applications, various summarization tasks, and understanding longer documents.', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 46, 'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''}),
  Document(page_content='Sergey Edunov\nThomas Scialom∗\nGenAI, Meta\nAbstract\nIn this work, we develop and release Llama 2, a collection of pretrained and fine-tuned\nlarge language models (LLMs) ranging in scale from 7 billion to 70 billion parameters.\nOur fine-tuned LLMs, called Llama 2-Chat, are optimized for dialogue use cases. Our\nmodels outperform open-source chat models on most benchmarks we tested, and based on\nour human evaluations for helpfulness and safety, may be a suitable substitute for closed-\nsource models. We provide a detailed description of our approach to fine-tuning and safety\nimprovements of Llama 2-Chat in order to enable the community to build on our work and\ncontribute to the responsible development of LLMs.\n∗Equal contribution, corresponding authors: {tscialom, htouvron}@meta.com\n†Second author\nContributions for all the authors can be found in Section A.1.\narXiv:2307.09288v2 [cs.CL]  19 Jul 2023', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 0, 'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''})]}

```
print(paper_rag.combine_documents_chain.llm_chain.prompt.messages[0].prompt.template)
```

```
Use the following pieces of context to answer the users question.
If you don't know the answer, just say that you don't know, don't try to make up an answer.
You are also provided the following Academic Knowledge Graph for checking affiliation. Please only use it when necessary
----------------
KG:
("Hugo Touvron", "affiliated_with", "Meta AI")
("Wei Chen", "affiliated_with", "Fudan University and Huazhong University of Science and Technology")
----------------
{context}
```

```
paper_rag.combine_documents_chain.llm_chain.prompt = deepcopy(paper_rag.combine_documents_chain.llm_chain.pro
paper_rag.combine_documents_chain.llm_chain.prompt.messages[0].prompt.template = """Use the following pieces o
You are also provided the following Academic Knowledge Graph for checking affiliation. Please only use it whe
----------------
KG:
("Hugo Touvron", "affiliated_with", "Meta AI")
("Wei Chen", "affiliated_with", "Fudan University and Huazhong University of Science and Technology")
----------------
{context}"""

paper_rag.combine_documents_chain.llm_chain.prompt.messages[1].prompt.template = """{question}\nLet's think st
```

```
paper_rag(question_answer_pairs["llm_paper"][5]["question"])
```

```
{'query': 'Is the first author of Llama and Llama 2 the same? If yes, please output <<YES>>, otherwise output <<NO>>.',
 'result': '1. Look for the names of the authors of Llama and Llama 2.\n2. Identify if the first author of Llama is the same as the first author of Llama 2.\n3. Output <<YES>> if they are the same, otherwise output <<NO>.\n\nBased on the given context, the first author of Llama is not explicitly mentioned. Therefore, we cannot determine if the first author of Llama and Llama 2 are the same.',
 'source_documents': [Document(page_content='Llama 2: Open Foundation and Fine-Tuned Chat Models\nHugo Touvron∗\nLouis Martin†\nKevin Stone†\nPeter Albert Amjad Almahairi Yasmine Babaei Nikolay Bashlykov Soumya Batra\nPrajjwal Bhargava Shruti Bhosale Dan Bikel Lukas Blecher Cristian Canton Ferrer Moya Chen\nGuillem Cucurull David Esiobu Jude Fernandes
```

Jeremy Fu Wenyin Fu Brian Fuller\nCynthia Gao Vedanuj Goswami Naman Goyal Anthony Hartshorn Saghar Hosseini Rui Hou\nHakan Inan Marcin Kardas Viktor Kerkez Madian Khabsa Isabel Kloumann Artem Korenev\nPunit Singh Koura Marie-Anne Lachaux Thibaut Lavril Jenya Lee Diana Liskovich\nYinghai Lu Yuning Mao Xavier Martinet Todor Mihaylov Pushkar Mishra\nIgor Molybog Yixin Nie Andrew Poulton Jeremy Reizenstein Rashi Rungta Kalyan Saladi\nAlan Schelten Ruan Silva Eric Michael Smith Ranjan Subramanian Xiaoqing Ellen Tan Binh Tang\nRoss Taylor Adina Williams Jian Xiang Kuan Puxin Xu Zheng Yan Iliyan Zarov Yuchen Zhang\nAngela Fan Melanie Kambadur Sharan Narang Aurelien Rodriguez Robert Stojnic\nSergey Edunov', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 0, 'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''}),
  Document(page_content='Figure 3: Safety human evaluation results for Llama 2-Chat compared to other open-source and closed-\nsource models. Human raters judged model generations for safety violations across ~2,000 adversarial\nprompts consisting of both single and multi-turn prompts. More details can be found in Section 4.4. It is\nimportant to caveat these safety results with the inherent bias of LLM evaluations due to limitations of the\nprompt set, subjectivity of the review guidelines, and subjectivity of individual raters. Additionally, these\nsafety evaluations are performed using content standards that are likely to be biased towards the Llama\n2-Chat models.\nWe are releasing the following models to the general public for research and commercial use‡:\n1. Llama 2, an updated version of Llama 1, trained on a new mix of publicly available data. We also\nincreased the size of the pretraining corpus by 40%, doubled the context length of the model, and', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 3, 'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''}),
  Document(page_content='• Armand Joulin, Edouard Grave, Guillaume Lample, and Timothee Lacroix, members of the original\nLlama team who helped get this work started.\n• Drew Hamlin, Chantal Mora, and Aran Mun, who gave us some design input on the figures in the\npaper.\n• Vijai Mohan for the discussions about RLHF that inspired our Figure 20, and his contribution to the\ninternal demo.\n• Early reviewers of this paper, who helped us improve its quality, including Mike Lewis, Joelle Pineau,\nLaurens van der Maaten, Jason Weston, and Omer Levy.\nA.2\nAdditional Details for Pretraining\nA.2.1\nArchitecture Changes Compared to Llama 1\nContext Length.\nWe expand the context window for Llama 2 from 2048 tokens to 4096 tokens. The longer\ncontext window enables models to process more information, which is particularly useful for supporting\nlonger histories in chat applications, various summarization tasks, and understanding longer documents.', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 46, 'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''}),
  Document(page_content='Sergey Edunov\nThomas Scialom∗\nGenAI, Meta\nAbstract\nIn this work, we develop and release Llama 2, a collection of pretrained and fine-tuned\nlarge language models (LLMs) ranging in scale from 7 billion to 70 billion parameters.\nOur fine-tuned LLMs, called Llama 2-Chat, are optimized for dialogue use cases. Our\nmodels outperform open-source chat models on most benchmarks we tested, and based on\nour human evaluations for helpfulness and safety, may be a suitable substitute for closed-\nsource models. We provide a detailed description of our approach to fine-tuning and safety\nimprovements of Llama 2-Chat in order to enable the community to build on our work and\ncontribute to the responsible development of LLMs.\n∗Equal contribution, corresponding authors: {tscialom, htouvron}@meta.com\n†Second author\nContributions for all the authors can be found in Section A.1.\narXiv:2307.09288v2 [cs.CL]  19 Jul 2023', metadata={'author': '', 'creationDate': 'D:20230720003036Z', 'creator': 'LaTeX with hyperref', 'file_path': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'format': 'PDF 1.5', 'keywords': '', 'modDate': 'D:20230720003036Z', 'page': 0, 'producer': 'pdfTeX-1.40.25', 'source': '/content/drive/My Drive/Colab Notebooks/HSBCRAG/data/llm_papers/Llama 2- Open Foundation and Fine-Tuned Chat Models.pdf', 'subject': '', 'title': '', 'total_pages': 77, 'trapped': ''})]}

```python
ensemble_rag = EnsembleRAG(intent_prompt=intent_domain_prompt, tools=tools, llm=chatgpt)
ensemble_rag_results = evaluator.test_rag_all(rag=ensemble_rag, question_answer_pairs=question_answer_pairs, \
ensemble_rag_results = evaluator.evaluate_rag_results(ensemble_rag_results)
evaluator.inspect_rag_results(ensemble_rag_results)
```

category=huifeng_faq, question=如何环球转账？, score=0.4, response=您可以通过以下步骤进行汇丰环球转账:

1. 登录汇丰银行网上银行，进入"我的银行"-"转账及货币兑换"-"转账及货币兑换"页面。

2. 在新交易中，"转出账户"选择外币储蓄账户。转入账户选择"我的账户"，并选择另一个汇丰其他国家/地区的外币账户。

3. 输入转账金额和币种，选择交易编码/原因，然后验证交易详情并确认，即可完成交易。

您还可以通过登录汇丰银行手机App进行相关操作。请注意，环球转账支持的币种包括美元、英镑、欧元、日元、港元、加拿大元、澳元、新加坡元、瑞士法郎和新西
category=jinronghuiyi_article, question=广发证券对中央金融工作会议的观点有哪些？, score=0.0, response=广发证券对中央金融工作会议的观
category=llm_paper, question=What is the author of Llama?, score=0.07142857142857142, response=1. The paper "LLaMA: Open
2. The first author listed is Hugo Touvron.
3. However, the paper also mentions that Llama was developed by Meta AI.
4. The Academic Knowledge Graph provided shows that Hugo Touvron is affiliated with Meta AI.
5. Therefore, it can be concluded that both Hugo Touvron and Meta AI are authors of Llama.
category=llm_paper, question=Is the first author of Llama and Llama 2 the same? If yes, please output <<YES>>, otherwise
2. Identify if the first author of Llama is the same as the first author of Llama 2.
3. Output <<YES>> if they are the same, otherwise output <<NO>>.

Based on the given context, the first author of Llama is not explicitly mentioned. Therefore, we cannot determine if the
category=llm_paper, question=Is there any LLM in financial area? What are they?, score=0.4, response=Yes, there are LLMs

```python
score = evaluator.compute_rag_score(ensemble_rag_results)
score
```

0.8709821428571428

```python
df.loc[len(df.index)] = [f"{best_prompt}+COT", f"{best_retriever}+SelfQueryRetriever", "ensemble_rag+KG", scor
```

- c、【todo】引入上下文

如果是多轮QA，需要引入Memory来记录上下文信息，来帮助提高生成效果。对于对话，可以使用Coversation Buffer：
https://python.langchain.com/docs/modules/memory/types/buffer

### d、【todo】Self-RAG

paper：https://github.com/AkariAsai/self-rag

# 第六章：总结

本文通过金融LLM+RAG的Demo，探索了RAG优化的一些思路。**从最初15%的准确率，通过一系列的优化，最终达到75+%**（Chroma似乎有随机性，结果每次跑有所不同，https://github.com/langchain-ai/langchain/issues/1946）。有以下一些Insight：

- 数据层：整个RAG的输入
  - 对其进行高质量的预处理，可以帮助后续的召回和生成模块；*（在本Demo有正向效果）*
  - 同时利用meta源数据，可以提高召回的相关性；*（在本Demo有正向效果）*
- 召回层：核心环节，负责检索送进LLM进行生成的原材料，其召回结果的全面性、多样性和相关性对生成质量至关重要。
  - 为了保证召回结果的相关性，可以通过改进Embedding模型和进行Rerank；*（在本Demo有正向效果）*
  - 为了提高召回的全面性，需要调节合适的chunk大小或者使用多粒度召回的方式；*（在本Demo有正向效果）*
  - 为了提高召回的多样性，可以采用MultiQuery和MultiVector等方法*（在本Demo中，没有验证到有效性）*。
- 生成层：RAG的最后一环，负责最终结果的生成。
  - 好的LLM模型本身就很大程度决定了生成质量的基础水平；*（本Demo没有进行验证，直接使用ChatGPT作为基础模型）*
  - 好的Prompt能帮助LLM更好的生成；*（在本Demo有正向效果）*
- 其他
  - 对于复杂的问题，使用COT技术进行问题的分解和步步推理可以提升效果；*（在本Demo有正向效果）*
  - 使用Agent（Intent Classification+EnsembleRAG/Tool）可以使用针对领域优化的RAG来提升整体效果；*（在本Demo有正向效果）*
  - 引入KG等外部信息，可以帮助解决LLM幻觉问题；*（在本Demo有正向效果）*

```
df.sort_values("score", ascending=False)
```

| | prompt | retriever | rag | score |
|---|---|---|---|---|
| 23 | default_prompt+COT | big_chunks_compression_retriever+SelfQueryRetr... | ensemble_rag+KG | 0.870982 |
| 22 | default_prompt | big_chunks_compression_retriever+SelfQueryRetr... | ensemble_rag+KG | 0.867262 |
| 18 | default_prompt | big_chunks_retriever | single_rag | 0.697768 |
| 4 | default_rag_prompt | big_chunks_retriever | single_rag | 0.510268 |
| 13 | fin_rag_prompt | big_chunks_compression_retriever | single_rag | 0.481250 |
| 6 | default_rag_prompt | big_chunks_compression_retriever | single_rag | 0.480941 |
| 19 | default_prompt | parent_chunks_retriever | single_rag | 0.473624 |