

书生·阅读

LLM 的局限性

- 知识时效性受限：如何让LLM能够获取最新的知识
- 专业能力有限：如何打造垂域大模型
- 定制化成本高：如何打造个人专属的LLM应用

两种不同解决策略：

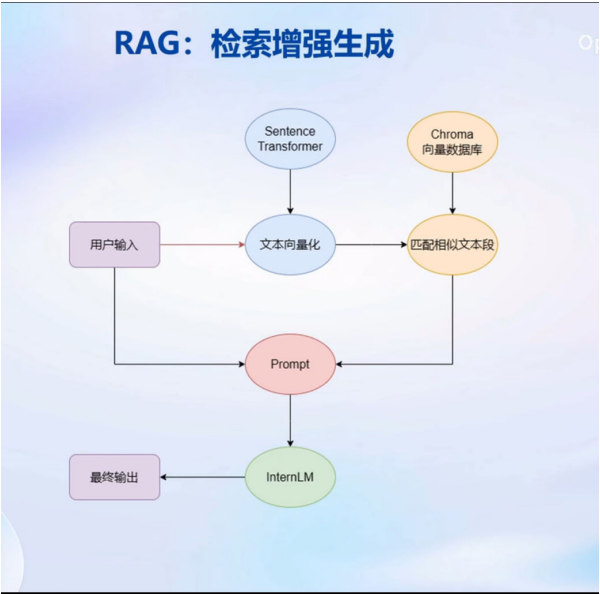
书生·阅读

RAG VS Finetune

- 低成本
- 可实时更新
- 受基座模型影响大
- 单次回答知识有限

- 可个性化微调
- 知识覆盖面广
- 成本高昂
- 无法实时更新

RAG: 外挂知识库，把检索的文章和用户提示词一起交给大模型，占用上下文长度大
RAG过程：



基于LangChain:

OpenMMLab



OpenMMLab 书生·读图

- 分区 cs 的第 2 页


```
root@intern-studio: ~/dat.x | create_db.py | root@intern-studio: ~/dat.x | +
(InternLM root@intern-studio:~/dat.x# git clone https://gitee.com/InternLM/lagent.git
Cloning into 'lagent'...
remote: Enumerating objects: 414, done.
remote: Counting objects: 100% (414/414), done.
remote: Compressing objects: 100% (188/188), done.
remote: Total 414 (delta 197), reused 414 (delta 197), pack-reused 0
Receiving objects: 100% (414/414), 214.97 KIB | 1005.00 KIB/s, done.
Resolving deltas: 100% (197/197), done.
(InternLM root@intern-studio:~/dat.x# git clone https://gitee.com/InternLM/InternLM.git
Cloning into 'InternLM'...
remote: Enumerating objects: 2604, done.
remote: Counting objects: 100% (592/592), done.
remote: Compressing objects: 100% (264/264), done.
remote: Total 2604 (delta 323), reused 581 (delta 318), pack-reused 2012
Receiving objects: 100% (2604/2604), 4.87 MiB | 2.61 MiB/s, done.
Resolving deltas: 100% (1607/1607), done.
Updating files: 100% (203/203), done.
(InternLM root@intern-studio:~/dat.x#
(InternLM root@intern-studio:~/dat.x#
(InternLM root@intern-studio:~/dat.x# cd ~/demo
bash: cd: /root/demo: No such file or directory
(InternLM root@intern-studio:~/dat.x# cd /root/data/demo
(InternLM root@intern-studio:~/dat.x# python create_db
python: can't open file '/root/data/demo/create_db': [Errno 2] No such file or directory
(InternLM root@intern-studio:~/dat.x# python create_db.py
0% | 0/25 [00:00:00, ?it/s]
/root/.conda/envs/InternLM/lib/python3.10/site-packages/unstructured/documents/html.py:498: FutureWarning: The behavior of this method will change in future versions. Use specific 'len(elem)' or 'elem is not None' te
st instead.
rows = body.findall("tr") if body else []
40% | 10/25 [00:18:00:15, 1.01s/it]
/root/.conda/envs/InternLM/lib/python3.10/site-packages/unstructured/documents/html.py:498: FutureWarning: The behavior of this method will change in future versions. Use specific 'len(elem)' or 'elem is not None' te
st instead.
rows = body.findall("tr") if body else []
100% | 25/25 [00:18:00:00, 1.35it/s]
100% | 9/9 [00:00:00:00, 23.57it/s]
100% | 18/18 [00:00:00:00, 50.69it/s]
100% | 72/72 [00:02:00:00, 28.71it/s]
100% | 113/113 [00:04:00:00, 24.15it/s]
100% | 26/26 [00:01:00:00, 23.28it/s]
(InternLM root@intern-studio:~/dat.x#
(InternLM root@intern-studio:~/dat.x#
```

构建检索问答链:

```
File Edit View Run Kernel Tabs Settings Help
+ / data / demo /
Name Last Modified
data_base 3 hours ago
create_db.py 3 hours ago
LLM.py 3 minutes ago
untitled.txt 22 seconds ago

root@x | LLM.py | untitled | root@x | downl | READ | react | downl | web_d | cli_den | downl | +
2 from langchain.vectorstores import Chroma
3 from langchain.embeddings.huggingface import HuggingFaceEmbeddings
4 import os
5 from LLM import InternLM_LLM
6 from langchain.prompts import PromptTemplate
7 from langchain.chains import RetrievalQA
8
9 def load_chain():
10     # 加载问答链
11     # 定义 Embeddings
12     embeddings = HuggingFaceEmbeddings(model_name="/root/data/model/sentence-transformer")
13
14     # 向量数据库持久化路径
15     persist_directory = 'data_base/vector_db/chroma'
16
17     # 加载数据库
18     vectordb = Chroma(
19         persist_directory=persist_directory, # 允许我们将persist_directory目录保存到磁盘上
20         embedding_function=embeddings
21     )
22
23     # 加载自定义 LLM
24     llm = InternLM_LLM(model_path = "/root/data/model/Shanghai_AI_Laboratory/InternLM-chat-7b")
25
26     # 定义一个 Prompt Template
27     template = """使用以下上下文来回答最后的问题。如果你不知道答案，就说你不知道，不要试图编造答
28     案。尽量使答案简明扼要。总是在回答的最后说“谢谢你的提问！”。
29     {context}
30     问题: {question}
31     有用的回答: """
32
33     QA_CHAIN_PROMPT = PromptTemplate(input_variables=["context","question"],template=template)
34
35     # 运行 chain
36     qa_chain = RetrievalQA.from_chain_type(llm,retriever=vectordb.as_retriever(),return_source_documents=True,chain_type_kwargs=
37     {"prompt":QA_CHAIN_PROMPT})
38
39     return qa_chain
```

ssh链接:

```
Windows PowerShell
版权所有 (C) Microsoft Corporation. 保留所有权利。

安装最新的 PowerShell, 了解新功能和改进! https://aka.ms/PSWindows

PS C:\Users\23594> ssh -CNq -L 7860:127.0.0.1:7860 root@ssh.intern-ai.org.cn -p 34795
```

在本地电脑打开web demo:

127.0.0.1:7860

InternLM

书生浦语

Chatbot

Prompt/问题

Chat

Clear console

提醒：

1. 初始化数据库时间可能较长，请耐心等待。

InternLM

书生浦语

Chatbot

InternLM是什么

仕嗨込竿潜能exp1031實驗屋上s海h人r工g智z能n实s验y室s原ori创gin

InternLM是什么

InternLM是一个开源的轻量级训练框架，旨在支持大模型训练，可在拥有数千个GPU的大型集群上进行预训练，并在单个GPU上进行微调，同时实现了卓越的性能优化。该框架在1024个GPU上训练时，可以实现近90%的加速效率。

Prompt/问题

|

Chat

Clear console

提醒：

1. 初始化数据库时间可能较长，请耐心等待。