# 基础岛任务笔记4-InternLM + LlamaIndex RAG 实践
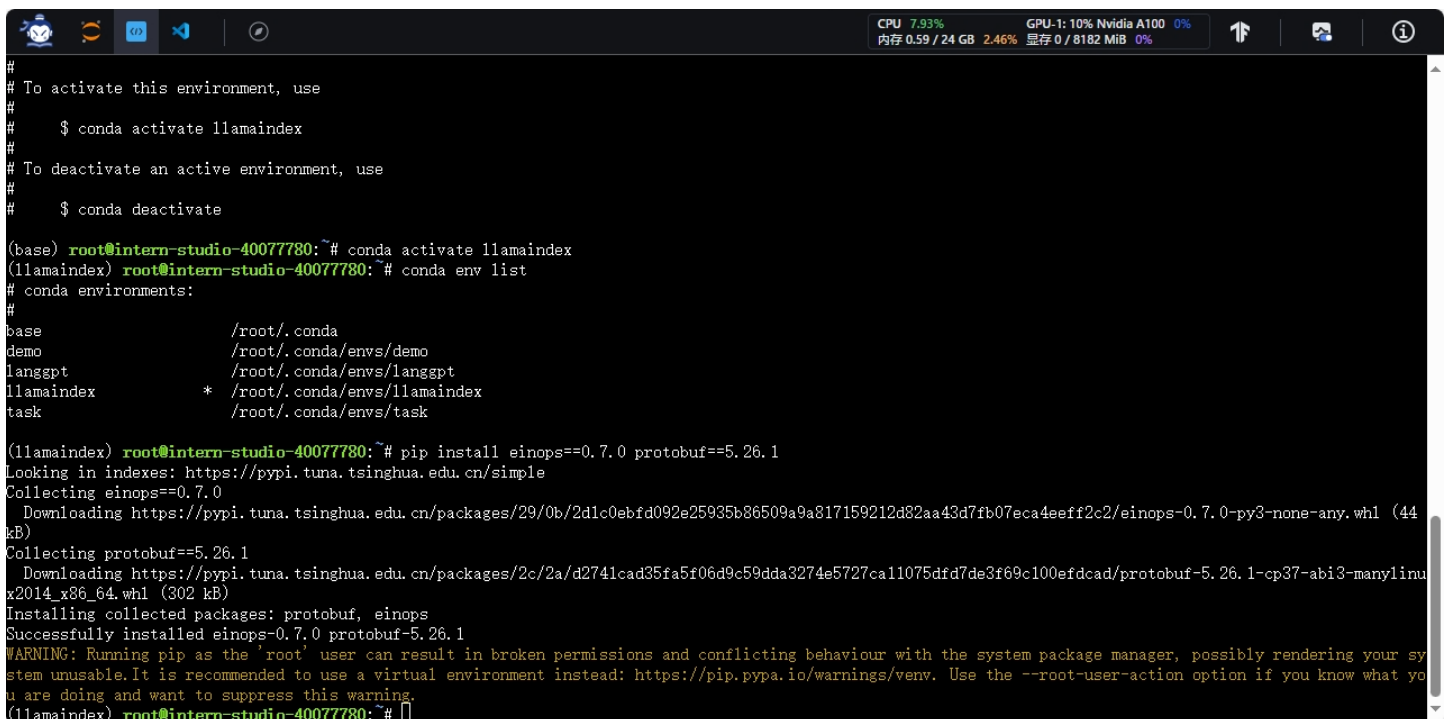
- **任务要求**：基于 LlamaIndex 构建自己的 RAG 知识库，寻找一个问题 A,这个问题在使用 LlamaIndex 之前InternLM2-Chat-1.8B模型不会回答，借助 LlamaIndex 后 InternLM2-Chat-1.8B 模型具备回答 A 的能力，截图保存。

## 环境，模型准备

## 配置基础环境

```
1  #创建新的conda环境，命名为 llamaindex，在命令行模式下运行：
2  conda create -n llamaindex python=3.10
3
4  #运行 conda 命令，激活 llamaindex 然后安装相关基础依赖 python 虚拟环境：
5  conda activate llamaindex
6  conda install pytorch==2.0.1 torchvision==0.15.2 torchaudio==2.0.2 pytorch-cuda=11.7 -c pytorch -c nvidia
7
8  #安装python 依赖包
9  pip install einops==0.7.0 protobuf==5.26.1
```

# 安装 llamaindex

```
1  pip install llama-index==0.10.38 llama-index-llms-huggingface==0.2.0
   "transformers[torch]==4.41.1" "huggingface_hub[inference]==0.23.1"
   huggingface_hub==0.23.1 sentence-transformers==2.7.0 sentencepiece==0.2.0
```

```
(llamaindex) root@intern-studio-40077780:~# pip install llama-index==0.10.38 llama-index-llms-huggingface==0.2.0 "transformers[torch]==4.41.1" "huggingface_hu
b[inference]==0.23.1" huggingface_hub==0.23.1 sentence-transformers==2.7.0 sentencepiece==0.2.0
```

# 下载Sentence Transformer模型

```
1  cd ~
2  mkdir llamaindex_demo
3  mkdir model
4  cd ~/llamaindex_demo
5  touch download_hf.py
```

打开 `download_hf.py` 贴入以下代码

```
1  import os
2  # 设置环境变量
3  os.environ['HF_ENDPOINT'] = 'https://hf-mirror.com'
4  # 下载模型
5  os.system('huggingface-cli download --resume-download sentence-
   transformers/paraphrase-multilingual-MiniLM-L12-v2 --local-dir
   /root/model/sentence-transformer')
```

然后，在 `/root/llamaindex_demo` 目录下执行该脚本：

```
1  cd /root/llamaindex_demo
2  conda activate llamaindex
3  python download_hf.py
```

```
(llamaindex) root@intern-studio-40077780:~# mkdir llamaindex_demo
(llamaindex) root@intern-studio-40077780:~# mkdir model
(llamaindex) root@intern-studio-40077780:~# cd llamaindex_demo
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# touch download_hf.py
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# ls
download_hf.py
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# ls
download_hf.py
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# cd
(llamaindex) root@intern-studio-40077780:~# ;s
bash: syntax error near unexpected token `;'
(llamaindex) root@intern-studio-40077780:~# ls
Tutorial  chat_history  demo  demo2  file  langgpt  llamaindex_demo  model  share
(llamaindex) root@intern-studio-40077780:~# vim download_hf.py
(llamaindex) root@intern-studio-40077780:~# ls
Tutorial  chat_history  demo  demo2  file  langgpt  llamaindex_demo  model  share
(llamaindex) root@intern-studio-40077780:~# cd llamaindex_demo/
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# ls
download_hf.py
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# vim download_hf.py
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# python download_hf.py
Fetching 15 files:   0%|                                                          |
/root/.conda/envs/llamaindex/lib/python3.10/site-packages/huggingface_hub/file_download.py:1132: FutureWarning: `resume_download` is depi
moved in version 1.0.0. Downloads always resume when possible. If you want to force a new download, use `force_download=True`.
  warnings.warn(
Downloading 'model.safetensors' to '/root/model/sentence-transformer/.huggingface/download/model.safetensors.eaa086f0ffee582aeb45b36e34c
a1bbc9bd955917b.incomplete'
Downloading 'modules.json' to '/root/model/sentence-transformer/.huggingface/download/modules.json.f7640f94e81bb7f4f04daf1668850b38763a1
```

# 下载NLTK相关资源

```
1  cd /root
2  git clone https://gitee.com/yzy0612/nltk_data.git  --branch gh-pages
3  cd nltk_data
4  mv packages/*  ./
5  cd tokenizers
6  unzip punkt.zip
7  cd ../taggers
8  unzip averaged_perceptron_tagger.zip
```

```
File "/root/.conda/envs/llamaindex/lib/python3.10/site-packages/requests/adapters.py", line 713, in send
    raise ReadTimeout(e, request=request)
requests.exceptions.ReadTimeout: (ReadTimeoutError("HTTPSConnectionPool(host='cdn-lfs.hf-mirror.com', port=443):
t ID: 20aa6663-c23f-4435-81db-ac0ac82f2a18)')
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# cd /root
(llamaindex) root@intern-studio-40077780:~# git clone https://gitee.com/yzy0612/nltk_data.git --branch gh-pages
Cloning into 'nltk_data'...
remote: Enumerating objects: 1692, done.
remote: Total 1692 (delta 0), reused 0 (delta 0), pack-reused 1692
Receiving objects: 100% (1692/1692), 952.80 MiB | 44.64 MiB/s, done.
Resolving deltas: 100% (909/909), done.
Updating files: 100% (244/244), done.
(llamaindex) root@intern-studio-40077780:~# cd nltk_data
(llamaindex) root@intern-studio-40077780:~/nltk_data# mv packages/* ./
(llamaindex) root@intern-studio-40077780:~/nltk_data# cd tokenizers
(llamaindex) root@intern-studio-40077780:~/nltk_data/tokenizers# unzip punkt.zip
Archive:  punkt.zip
   creating: punkt/
  inflating: punkt/greek.pickle
  inflating: punkt/estonian.pickle
  inflating: punkt/turkish.pickle
  inflating: punkt/polish.pickle
   creating: punkt/PY3/
  inflating: punkt/PY3/greek.pickle
  inflating: punkt/PY3/estonian.pickle
  inflating: punkt/PY3/turkish.pickle
  inflating: punkt/PY3/polish.pickle
  inflating: punkt/PY3/russian.pickle
  inflating: punkt/PY3/czech.pickle
  inflating: punkt/PY3/portuguese.pickle
  inflating: punkt/PY3/README
  inflating: punkt/PY3/dutch.pickle
  inflating: punkt/PY3/norwegian.pickle
```

```
(llamaindex) root@intern-studio-40077780:~/nltk_data/tokenizers# cd ../taggers
(llamaindex) root@intern-studio-40077780:~/nltk_data/taggers# unzip averaged_perceptron_tagger.zip
Archive:  averaged_perceptron_tagger.zip
   creating: averaged_perceptron_tagger/
  inflating: averaged_perceptron_tagger/averaged_perceptron_tagger.pickle
(llamaindex) root@intern-studio-40077780:~/nltk_data/taggers#
```

# LlamaIndex HuggingFaceLLM

```
1  #运行以下指令，把 InternLM2 1.8B 软连接出来
2  cd ~/model
3  ln -s /root/share/new_models/Shanghai_AI_Laboratory/internlm2-chat-1_8b/ ./
4
5  #运行以下指令，新建一个python文件
6  cd ~/llamaindex_demo
7  touch llamaindex_internlm.py
```

```
1  #打开llamaindex_internlm.py 贴入以下代码
2  from llama_index.llms.huggingface import HuggingFaceLLMfrom
   llama_index.core.llms import ChatMessagellm = HuggingFaceLLM(
3      model_name="/root/model/internlm2-chat-1_8b",
```

```
4     tokenizer_name="/root/model/internlm2-chat-1_8b",
5     model_kwargs={"trust_remote_code":True},
6     tokenizer_kwargs={"trust_remote_code":True}
7 )
8
9 rsp = llm.chat(messages=[ChatMessage(content="xtuner是什么？")])
10 print(rsp)
```

```
1 #运行
2 conda activate llamaindex
3 cd ~/llamaindex_demo/
4 python llamaindex_internlm.py
```

问了一个"xtuner是什么？"的问题，回答的不咋地



很无奈，过载了，开发机10%的cpu果然有点不够

# LlamaIndex RAG

```
1  #安装 LlamaIndex 词嵌入向量依赖
2  conda activate llamaindex
3  pip install llama-index-embeddings-huggingface==0.2.0 llama-index-embeddings-
   instructor==0.1.3
4
5  #运行以下命令，获取知识库
6  cd ~/llamaindex_demo
7  mkdir data
8  cd data
9  git clone https://github.com/InternLM/xtuner.git
10 mv xtuner/README_zh-CN.md ./
11
12 #运行以下指令，新建一个python文件
13 cd ~/llamaindex_demo
14 touch llamaindex_RAG.py
```



llamaIndex_RAG.py写入如下代码，使大模型可以从词嵌入向量获得信息

```
1
2  from llama_index.core import VectorStoreIndex, SimpleDirectoryReader, Settings
```

```
 3
 4  from llama_index.embeddings.huggingface import HuggingFaceEmbedding
 5  from llama_index.llms.huggingface import HuggingFaceLLM
 6
 7  #初始化一个HuggingFaceEmbedding对象，用于将文本转换为向量表示
 8  embed_model = HuggingFaceEmbedding(
 9      #指定了一个预训练的sentence-transformer模型的路径
10      model_name="/root/model/sentence-transformer"
11  )
12  #将创建的嵌入模型赋值给全局设置的embed_model属性，
13  #这样在后续的索引构建过程中就会使用这个模型。
14  Settings.embed_model = embed_model
15
16  llm = HuggingFaceLLM(
17      model_name="/root/model/internlm2-chat-1_8b",
18      tokenizer_name="/root/model/internlm2-chat-1_8b",
19      model_kwargs={"trust_remote_code":True},
20      tokenizer_kwargs={"trust_remote_code":True}
21  )
22  #设置全局的llm属性，这样在索引查询时会使用这个模型。
23  Settings.llm = llm
24
25  #从指定目录读取所有文档，并加载数据到内存中
26  documents = SimpleDirectoryReader("/root/llamaindex_demo/data").load_data()
27  #创建一个VectorStoreIndex，并使用之前加载的文档来构建索引。
28  # 此索引将文档转换为向量，并存储这些向量以便于快速检索。
29  index = VectorStoreIndex.from_documents(documents)
30  # 创建一个查询引擎，这个引擎可以接收查询并返回相关文档的响应。
31  query_engine = index.as_query_engine()
32  response = query_engine.query("xtuner是什么?")
33
34  print(response)
```

再次过载，无奈

```
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo/data# cd ..
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# touch llamaindex_RAG.py
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# ls
data  download_hf.py  llamaindex_RAG.py  llamaindex_internlm.py
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# vim llamaindex_RAG.py
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# ls
data  download_hf.py  llamaindex_RAG.py  llamaindex_internlm.py
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# python llamaindex_RAG.py
/root/.conda/envs/llamaindex/lib/python3.10/site-packages/pydantic/_internal/_fields.py:161: UserWarning: Field "model_id" h
ce "model_".

You may be able to resolve this warning by setting `model_config['protected_namespaces'] = ()`.
  warnings.warn(
Traceback (most recent call last):
  File "/root/llamaindex_demo/llamaindex_RAG.py", line 8, in <module>
    embed_model = HuggingFaceEmbedding(
  File "/root/.conda/envs/llamaindex/lib/python3.10/site-packages/llama_index/embeddings/huggingface/base.py", line 86, in _
    self._model = SentenceTransformer(
  File "/root/.conda/envs/llamaindex/lib/python3.10/site-packages/sentence_transformers/SentenceTransformer.py", line 197, i
    modules = self._load_sbert_model(
  File "/root/.conda/envs/llamaindex/lib/python3.10/site-packages/sentence_transformers/SentenceTransformer.py", line 1296,
    module = Transformer(model_name_or_path, cache_dir=cache_folder, **kwargs)
  File "/root/.conda/envs/llamaindex/lib/python3.10/site-packages/sentence_transformers/models/Transformer.py", line 36, in
    self._load_model(model_name_or_path, config, cache_dir, **model_args)
  File "/root/.conda/envs/llamaindex/lib/python3.10/site-packages/sentence_transformers/models/Transformer.py", line 65, in
    self.auto_model = AutoModel.from_pretrained(
  File "/root/.conda/envs/llamaindex/lib/python3.10/site-packages/transformers/models/auto/auto_factory.py", line 563, in fr
    return model_class.from_pretrained(
  File "/root/.conda/envs/llamaindex/lib/python3.10/site-packages/transformers/modeling_utils.py", line 3305, in from_pretra
    raise EnvironmentError(
OSError: Error no file named pytorch_model.bin, model.safetensors, tf_model.h5, model.ckpt.index or flax_model.msgpack found
-transformer.
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# []
```

# LlamaIndex web

```
1  #运行之前首先安装依赖
2  pip install streamlit==1.36.0
3
4  #运行以下指令，新建一个python文件
5  cd ~/llamaindex_demo
6  touch app.py
```

```
-transformer.
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# touch app.py
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# vim app.py
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# streamlit run app.py
bash: streamlit: command not found
(llamaindex) root@intern-studio-40077780:~/llamaindex_demo# cd ..
(llamaindex) root@intern-studio-40077780:~# pip install streamlit==1.36.0
Looking in indexes: https://pypi.tuna.tsinghua.edu.cn/simple
Collecting streamlit==1.36.0
  Downloading https://pypi.tuna.tsinghua.edu.cn/packages/c6/51/f140402202af6ce1bf747243f66415c5eb2f43ba2e2ac419a
whl (8.6 MB)
                                        8.6/8.6 MB 31.7 MB/s eta 0:00:00
Collecting altair<6,>=4.0 (from streamlit==1.36.0)
  Using cached https://pypi.tuna.tsinghua.edu.cn/packages/9b/52/4a86a4fa1cc2aae79137cc9510b7080c3e5aede2310d14fa
8 kB)
Collecting blinker<2,>=1.0.0 (from streamlit==1.36.0)
  Using cached https://pypi.tuna.tsinghua.edu.cn/packages/bb/2a/10164ed1f31196a2f7f3799368a821765c62851ead0e630a
.5 kB)
Collecting cachetools<6,>=4.0 (from streamlit==1.36.0)
  Using cached https://pypi.tuna.tsinghua.edu.cn/packages/a4/07/14f8ad37f2d12a5ce41206c21820d8cb6561b728e51fad45
 (9.5 kB)
Requirement already satisfied: click<9,>=7.0 in ./.conda/envs/llamaindex/lib/python3.10/site-packages (from stre
Requirement already satisfied: numpy<3,>=1.20 in ./.conda/envs/llamaindex/lib/python3.10/site-packages (from str
```

```python
#app.py文件
import streamlit as st
from llama_index.core import VectorStoreIndex, SimpleDirectoryReader, Settings
from llama_index.embeddings.huggingface import HuggingFaceEmbedding
from llama_index.llms.huggingface import HuggingFaceLLM

st.set_page_config(page_title="llama_index_demo", page_icon="🦜🔗")
st.title("llama_index_demo")

# 初始化模型
@st.cache_resource
def init_models():
    embed_model = HuggingFaceEmbedding(
        model_name="/root/model/sentence-transformer"
    )
    Settings.embed_model = embed_model

    llm = HuggingFaceLLM(
        model_name="/root/model/internlm2-chat-1_8b",
        tokenizer_name="/root/model/internlm2-chat-1_8b",
        model_kwargs={"trust_remote_code": True},
        tokenizer_kwargs={"trust_remote_code": True}
    )
    Settings.llm = llm

    documents = SimpleDirectoryReader("/root/llamaindex_demo/data").load_data()
    index = VectorStoreIndex.from_documents(documents)
    query_engine = index.as_query_engine()
```

```python
29
30          return query_engine
31
32      # 检查是否需要初始化模型
33      if 'query_engine' not in st.session_state:
34          st.session_state['query_engine'] = init_models()
35
36      def greet2(question):
37          response = st.session_state['query_engine'].query(question)
38          return response
39
40
41      # Store LLM generated responses
42      if "messages" not in st.session_state.keys():
43          st.session_state.messages = [{"role": "assistant", "content": "你好，我是你的
    助手，有什么我可以帮助你的吗？"}]
44
45          # Display or clear chat messages
46      for message in st.session_state.messages:
47          with st.chat_message(message["role"]):
48              st.write(message["content"])
49
50      def clear_chat_history():
51          st.session_state.messages = [{"role": "assistant", "content": "你好，我是你的
    助手，有什么我可以帮助你的吗？"}]
52
53      st.sidebar.button('Clear Chat History', on_click=clear_chat_history)
54
55      # Function for generating LLaMA2 response
56      def generate_llama_index_response(prompt_input):
57          return greet2(prompt_input)
58
59      # User-provided prompt
60      if prompt := st.chat_input():
61          st.session_state.messages.append({"role": "user", "content": prompt})
62          with st.chat_message("user"):
63              st.write(prompt)
64
65      # Gegenerate_llama_index_response last message is not from assistant
66      if st.session_state.messages[-1]["role"] != "assistant":
67          with st.chat_message("assistant"):
68              with st.spinner("Thinking..."):
69                  response = generate_llama_index_response(prompt)
70                  placeholder = st.empty()
71                  placeholder.markdown(response)
72          message = {"role": "assistant", "content": response}
73          st.session_state.messages.append(message)
```

运行 `streamlit run app.py`



```
Collecting rpds-py>=0.7.1 (from jsonschema>=3.0->altair<6,>=4.0->streamlit==1.36.0)
  Using cached https://pypi.tuna.tsinghua.edu.cn/packages/4a/6d/1166a157b227f2333f8e8ae320b6b7ea2a6a38fbe7a3563ad76dffc8608d/rpds_py-0.20.0-cp310-cp310-manyli
nux_2_17_x86_64.manylinux2014_x86_64.whl (354 kB)
Collecting mdurl~=0.1 (from markdown-it-py>=2.2.0->rich<14,>=10.14.0->streamlit==1.36.0)
  Using cached https://pypi.tuna.tsinghua.edu.cn/packages/b3/38/89ba8ad64ae25be8de66a6d463314cf1eb366222074cfda9ee839c56a4b4/mdurl-0.1.2-py3-none-any.whl (10.
0 kB)
Requirement already satisfied: six>=1.5 in ./.conda/envs/llamaindex/lib/python3.10/site-packages (from python-dateutil>=2.8.2->pandas<3,>=1.3.0->streamlit==1.
36.0) (1.16.0)
Installing collected packages: watchdog, tornado, toml, smmap, rpds-py, pygments, pyarrow, narwhals, mdurl, cachetools, blinker, referencing, pydeck, markdown
-it-py, gitdb, rich, jsonschema-specifications, gitpython, jsonschema, altair, streamlit
Successfully installed altair-5.4.1 blinker-1.8.2 cachetools-5.5.0 gitdb-4.0.11 gitpython-3.1.43 jsonschema-4.23.0 jsonschema-specifications-2023.12.1 markdow
n-it-py-3.0.0 mdurl-0.1.2 narwhals-1.6.2 pyarrow-17.0.0 pydeck-0.9.1 pygments-2.18.0 referencing-0.35.1 rich-13.8.0 rpds-py-0.20.0 smmap-5.0.1 streamlit-1.36.
0 toml-0.10.2 tornado-6.4.1 watchdog-4.0.2
WARNING: Running pip as the 'root' user can result in broken permissions and conflicting behaviour with the system package manager, possibly rendering your sy
stem unusable.It is recommended to use a virtual environment instead: https://pip.pypa.io/warnings/venv. Use the --root-user-action option if you know what yo
u are doing and want to suppress this warning.
(llamaindex) root@intern-studio-40077780:~# streamlit run app.py
Usage: streamlit run [OPTIONS] TARGET [ARGS]...
Try 'streamlit run --help' for help.

Error: Invalid value: File does not exist: app.py
(llamaindex) root@intern-studio-40077780:~# streamlit run llamaindex_demo/app.py

Collecting usage statistics. To deactivate, set browser.gatherUsageStats to false.


  You can now view your Streamlit app in your browser.

  Local URL: http://localhost:8501
  Network URL: http://192.168.236.122:8501
  External URL: http://192.168.236.122:8501
```