

RAG项目实战（使用llamaindex构建自己的知识库）

RAG项目实战（使用llamaindex构建自己的知识库）

- 1.环境配置
- 2.下载Sentence Transformer 模型
- 3.下载InternLM2 1.8B/qwen2.5_0.5B 模型
- 4.创建知识库
- 5.创建web应用

1.环境配置

使用 `conda` 配置项目python环境

```
# 创建环境
conda create -n rag python=3.10

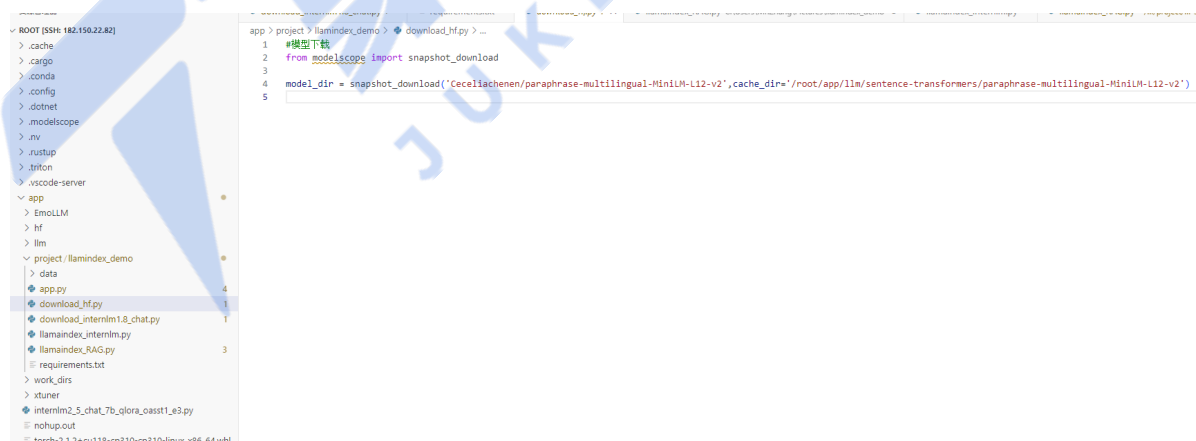
# 激活环境
conda activate rag

# 安装一些必要的库(requirements.txt在项目代码包中)
pip install -r requirements.txt
```

2.下载Sentence Transformer 模型

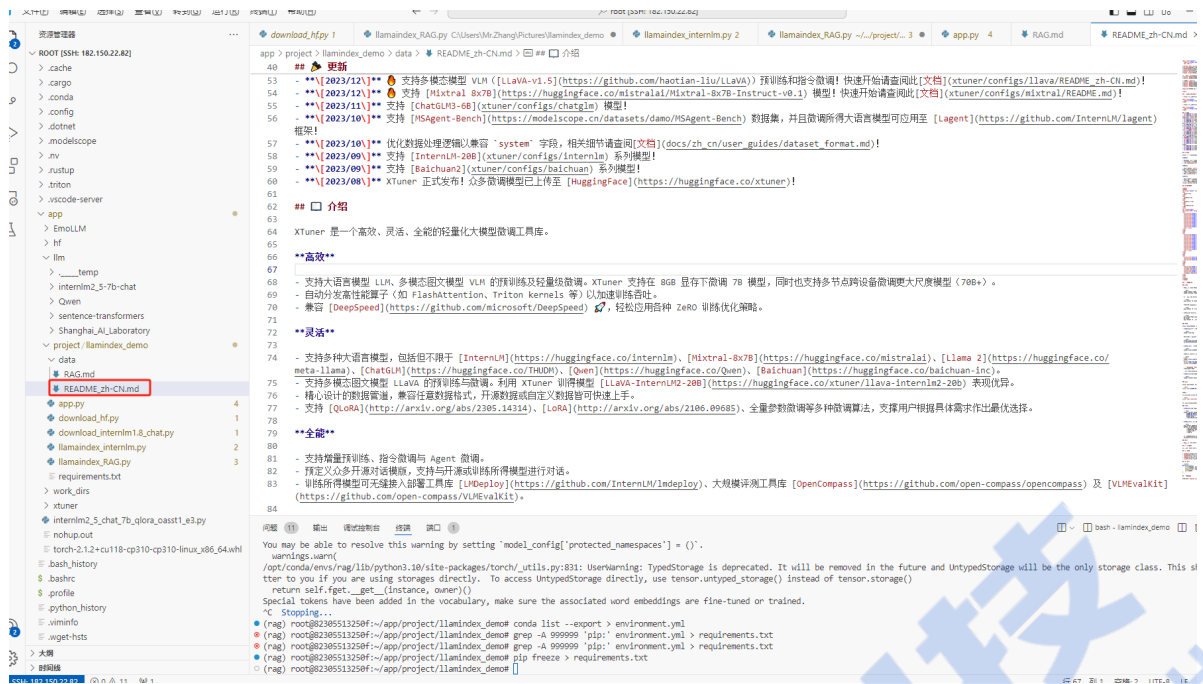
在进行RAG之前，需要使用词向量模型进行Embedding，将文本进行向量化处理，此处选择 Sentence Transformer 模型。执行

llamaindex_demo/download_hf.py下载。



3.下载InternLM2 1.8B/qwen2.5_0.5B 模型

执行llamaindex_demo/download_internlm1.8_chat.py下载。



执行llamaindex_demo/llamaindex_RAG.py，运行次测试后，可以看到可以正确回答，并且可以给出回答的出处：



5.创建web应用

执行llamaindex_demo/app.py，

```
streamlit run app.py
```

运行后可以打开网页端，可以进行提问：

