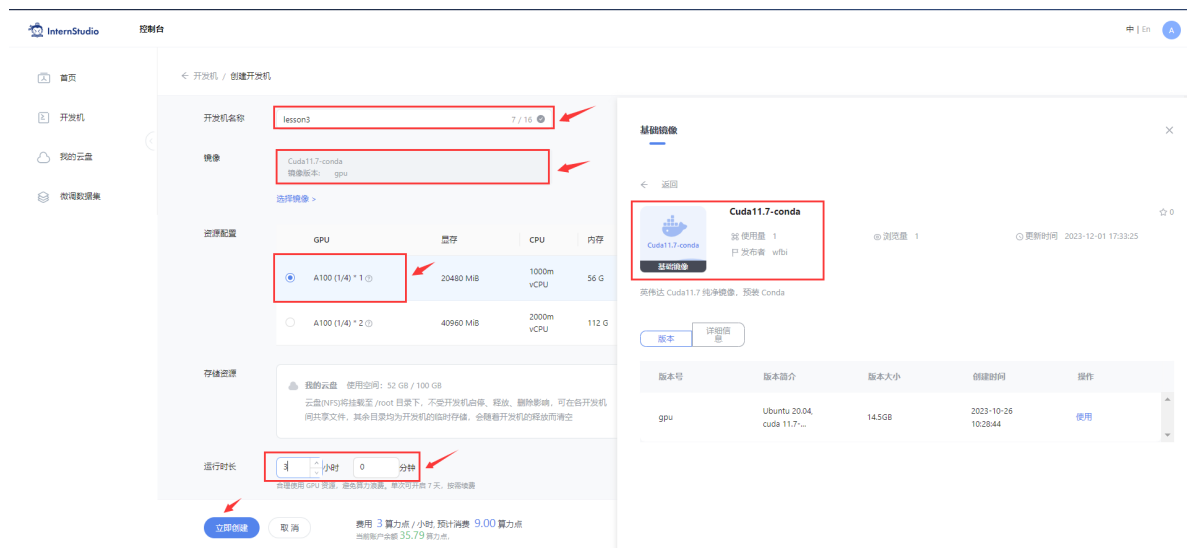


# 1 环境配置

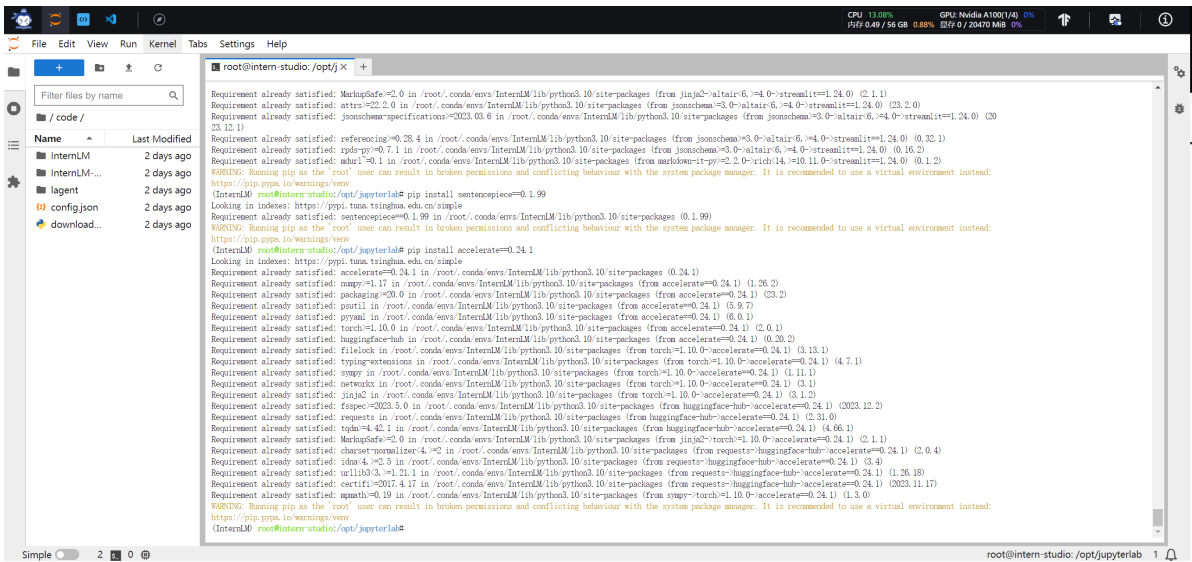
## 1.1 创建开发机

【开发机名称：lesson3，镜像：Cuda11.7-conda，资源配置：A100(1/4)，运行时长：3个小时】



## 1.2 创建Conda环境

```
1 # 进入Conda环境
2 bash
3
4 # 克隆Conda环境
5 conda create -n InternLM --clone /share/conda_envs/internlm-base
6
7 # 激活Conda环境
8 conda activate InternLM
9
10 # 升级pip
11 python -m pip install --upgrade pip
12
13 # 安装相关依赖
14 pip install modelscope==1.9.5
15 pip install transformers==4.35.2
16 pip install streamlit==1.24.0
17 pip install sentencepiece==0.1.99
18 pip install accelerate==0.24.1
```

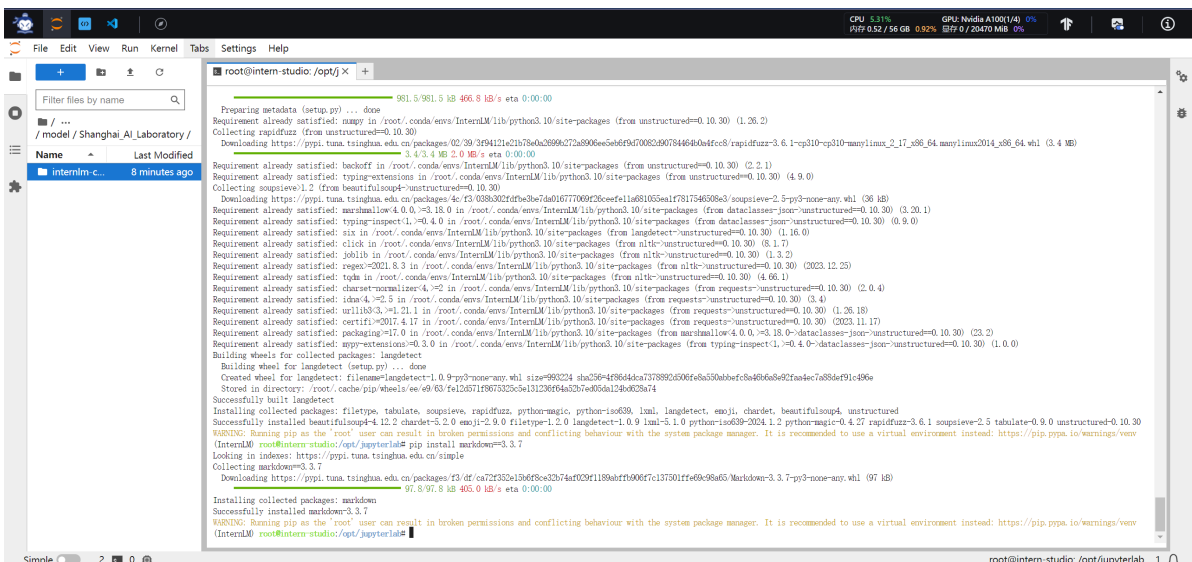


## 1.3 下载模型

- 1 `mkdir -p /root/data/model/Shanghai_AI_Laboratory`
- 2 `cp -r /root/share/temp/model_repos/internlm-chat-7b /root/data/model/Shanghai_AI_Laboratory/internlm-chat-7b`

## 1.4 安装LangChain依赖

- 1 `pip install langchain==0.0.292`
- 2 `pip install gradio==4.4.0`
- 3 `pip install chromadb==0.4.15`
- 4 `pip install sentence-transformers==2.2.2`
- 5 `pip install unstructured==0.10.30`
- 6 `pip install markdown==3.3.7`



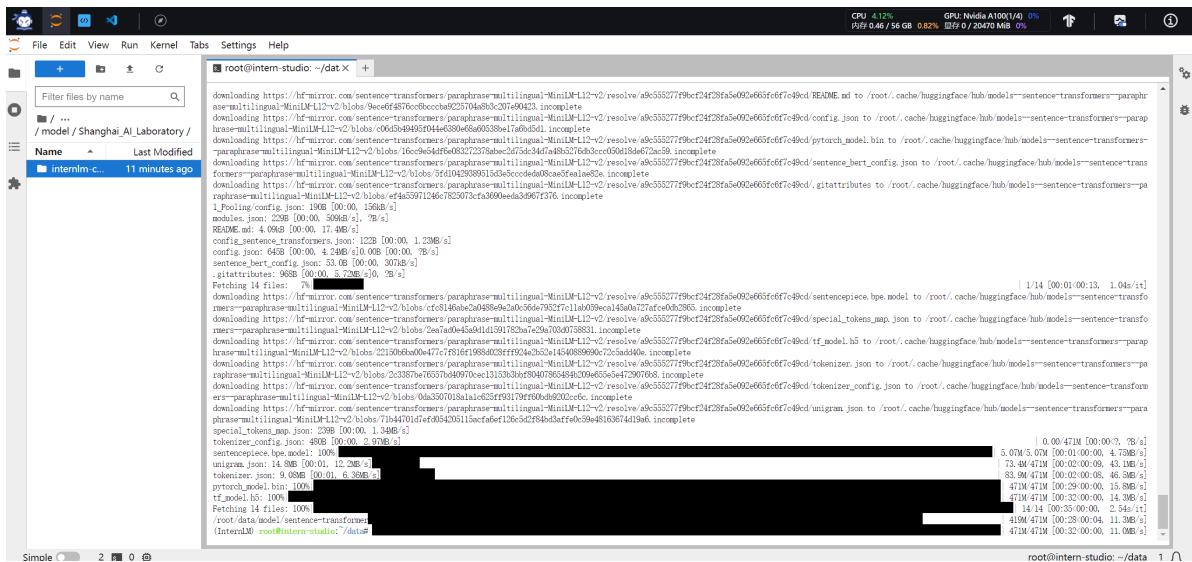
## 1.5 下载开源词向量模型

- 1 `# 安装huggingface_hub`
- 2 `pip install -U huggingface_hub`
- 3
- 4 `# 切换路径`
- 5 `cd /root/data`

```

6
7 # 创建下载脚本
8 vim download_hf.py
9
10 # 填入以下代码
11 import os
12
13 # 设置环境变量
14 os.environ['HF_ENDPOINT'] = 'https://hf-mirror.com'
15
16 # 下载模型
17 os.system('huggingface-cli download --resume-download sentence-
transformers/paraphrase-multilingual-MiniLM-L12-v2 --local-dir
/root/data/model/sentence-transformer')
18
19 # 执行下载脚本
20 python download_hf.py

```



## 1.6 下载NLTK相关资源

```

1 cd /root
2 git clone https://gitee.com/zyy0612/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

```



## 2.2 构建向量数据库

```
1 cd /root/data
2 mkdir demo && cd demo
3 touch create_db.py
4 python create_db.py
```

create\_db.py内容如下:

```
1 # 首先导入所需第三方库
2 from langchain.document_loaders import UnstructuredFileLoader
3 from langchain.document_loaders import UnstructuredMarkdownLoader
4 from langchain.text_splitter import RecursiveCharacterTextSplitter
5 from langchain.vectorstores import Chroma
6 from langchain.embeddings.huggingface import HuggingFaceEmbeddings
7 from tqdm import tqdm
8 import os
9
10 # 获取文件路径函数
11 def get_files(dir_path):
12     # args: dir_path, 目标文件夹路径
13     file_list = []
14     for filepath, dirnames, filenames in os.walk(dir_path):
15         # os.walk 函数将递归遍历指定文件夹
16         for filename in filenames:
17             # 通过后缀名判断文件类型是否满足要求
18             if filename.endswith(".md"):
19                 # 如果满足要求, 将其绝对路径加入到结果列表
20                 file_list.append(os.path.join(filepath, filename))
21             elif filename.endswith(".txt"):
22                 file_list.append(os.path.join(filepath, filename))
23     return file_list
24
25 # 加载文件函数
26 def get_text(dir_path):
27     # args: dir_path, 目标文件夹路径
28     # 首先调用上文定义的函数得到目标文件路径列表
29     file_lst = get_files(dir_path)
30     # docs 存放加载之后的纯文本对象
31     docs = []
32     # 遍历所有目标文件
33     for one_file in tqdm(file_lst):
34         file_type = one_file.split('.')[-1]
35         if file_type == 'md':
36             loader = UnstructuredMarkdownLoader(one_file)
37         elif file_type == 'txt':
38             loader = UnstructuredFileLoader(one_file)
39         else:
40             # 如果是不符合条件的文件, 直接跳过
41             continue
42         docs.extend(loader.load())
43     return docs
44
45 # 目标文件夹
```

```

46 tar_dir = [
47     "/root/data/InternLM",
48     "/root/data/InternLM-XComposer",
49     "/root/data/lagent",
50     "/root/data/lmdeploy",
51     "/root/data/opencompass",
52     "/root/data/xtuner"
53 ]
54
55 # 加载目标文件
56 docs = []
57 for dir_path in tar_dir:
58     docs.extend(get_text(dir_path))
59
60 # 对文本进行分块
61 text_splitter = RecursiveCharacterTextSplitter(
62     chunk_size=500, chunk_overlap=150)
63 split_docs = text_splitter.split_documents(docs)
64
65 # 加载开源词向量模型
66 embeddings = HuggingFaceEmbeddings(model_name="/root/data/model/sentence-
transformer")
67
68 # 构建向量数据库
69 # 定义持久化路径
70 persist_directory = 'data_base/vector_db/chroma'
71 # 加载数据库
72 vectordb = Chroma.from_documents(
73     documents=split_docs,
74     embedding=embeddings,
75     persist_directory=persist_directory # 允许我们将persist_directory目录保存
到磁盘上
76 )
77 # 将加载的向量数据库持久化到磁盘上
78 vectordb.persist()

```

## 2.3 InternLM接入LangChain

为便捷构建 LLM 应用，需要基于本地部署的 InternLM，继承 LangChain 的 LLM 类自定义一个 InternLM LLM 子类，从而实现将 InternLM 接入到 LangChain 框架中。

```

1 | cd /root/data/demo
2 | vim LLM.py

```

LLM.py代码内容如下：

```

1 | from langchain.llms.base import LLM
2 | from typing import Any, List, Optional
3 | from langchain.callbacks.manager import CallbackManagerForLLMRun
4 | from transformers import AutoTokenizer, AutoModelForCausalLM
5 | import torch
6 |

```

```

7 class InternLM_LLM(LLM):
8     # 基于本地 InternLM 自定义 LLM 类
9     tokenizer : AutoTokenizer = None
10    model: AutoModelForCausalLM = None
11
12    def __init__(self, model_path :str):
13        # model_path: InternLM 模型路径
14        # 从本地初始化模型
15        super().__init__()
16        print("正在从本地加载模型...")
17        self.tokenizer = AutoTokenizer.from_pretrained(model_path,
18            trust_remote_code=True)
19        self.model = AutoModelForCausalLM.from_pretrained(model_path,
20            trust_remote_code=True).to(torch.bfloat16).cuda()
21        self.model = self.model.eval()
22        print("完成本地模型的加载")
23
24    def _call(self, prompt : str, stop: Optional[List[str]] = None,
25        run_manager: Optional[CallbackManagerForLLMRun] = None,
26        **kwargs: Any):
27        # 重写调用函数
28        system_prompt = """"You are an AI assistant whose name is InternLM
29        (书生·浦语).
30        - InternLM (书生·浦语) is a conversational language model that is
31        developed by Shanghai AI Laboratory (上海人工智能实验室). It is designed to be
32        helpful, honest, and harmless.
33        - InternLM (书生·浦语) can understand and communicate fluently in the
34        language chosen by the user such as English and 中文.
35        """
36        messages = [(system_prompt, '')]
37        response, history = self.model.chat(self.tokenizer, prompt,
38            history=messages)
39        return response
40
41    @property
42    def _llm_type(self) -> str:
43        return "InternLM"

```

## 2.4 构建检索问答链

```

1 import gradio as gr
2
3 # 实例化核心功能对象
4 model_center = Model_center()
5 # 创建一个 web 界面
6 block = gr.Blocks()
7 with block as demo:
8     with gr.Row(equal_height=True):
9         with gr.Column(scale=15):
10             # 展示的页面标题
11             gr.Markdown("""<h1><center>InternLM</center></h1>
12                 <center>书生浦语</center>
13                 """)
14

```

```

15     with gr.Row():
16         with gr.Column(scale=4):
17             # 创建一个聊天机器人对象
18             chatbot = gr.Chatbot(height=450, show_copy_button=True)
19             # 创建一个文本框组件，用于输入 prompt。
20             msg = gr.Textbox(label="Prompt/问题")
21
22             with gr.Row():
23                 # 创建提交按钮。
24                 db_wo_his_btn = gr.Button("Chat")
25             with gr.Row():
26                 # 创建一个清除按钮，用于清除聊天机器人组件的内容。
27                 clear = gr.ClearButton(
28                     components=[chatbot], value="Clear console")
29
30             # 设置按钮的点击事件。当点击时，调用上面定义的 qa_chain_self_answer 函数，并
            传入用户的消息和聊天历史记录，然后更新文本框和聊天机器人组件。
31             db_wo_his_btn.click(model_center.qa_chain_self_answer, inputs=[
32                 msg, chatbot], outputs=[msg, chatbot])
33
34             gr.Markdown("""提醒: <br>
35             1. 初始化数据库时间可能较长，请耐心等待。
36             2. 使用中如果出现异常，将会在文本输入框进行展示，请不要惊慌。 <br>
37             """)
38 gr.close_all()
39 # 直接启动
40 demo.launch()

```

## 2.5 部署Web Demo

定义一个类，该类负责加载并存储检索问答链，并响应 Web 界面里调用检索问答链进行回答的动作：

```

1 class Model_center():
2     """
3     存储检索问答链的对象
4     """
5     def __init__(self):
6         # 构造函数，加载检索问答链
7         self.chain = load_chain()
8
9     def qa_chain_self_answer(self, question: str, chat_history: list = []):
10        """
11        调用问答链进行回答
12        """
13        if question == None or len(question) < 1:
14            return "", chat_history
15        try:
16            chat_history.append(
17                (question, self.chain({"query": question})["result"]))
18            # 将问答结果直接附加到问答历史中，Gradio 会将其展示出来
19            return "", chat_history
20        except Exception as e:
21            return e, chat_history
22

```



```

1  import gradio as gr
2
3  # 实例化核心功能对象
4  model_center = Model_center()
5  # 创建一个 web 界面
6  block = gr.Blocks()
7  with block as demo:
8      with gr.Row(equal_height=True):
9          with gr.Column(scale=15):
10             # 展示的页面标题
11             gr.Markdown("""<h1><center>InternLM</center></h1>
12                 <center>书生浦语</center>
13                 """)
14
15             with gr.Row():
16                 with gr.Column(scale=4):
17                     # 创建一个聊天机器人对象
18                     chatbot = gr.Chatbot(height=450, show_copy_button=True)
19                     # 创建一个文本框组件，用于输入 prompt。
20                     msg = gr.Textbox(label="Prompt/问题")
21
22                     with gr.Row():
23                         # 创建提交按钮。
24                         db_wo_his_btn = gr.Button("Chat")
25                     with gr.Row():
26                         # 创建一个清除按钮，用于清除聊天机器人组件的内容。
27                         clear = gr.ClearButton(
28                             components=[chatbot], value="Clear console")
29
30             # 设置按钮的点击事件。当点击时，调用上面定义的 qa_chain_self_answer 函数，并
31             # 传入用户的消息和聊天历史记录，然后更新文本框和聊天机器人组件。
32             db_wo_his_btn.click(model_center.qa_chain_self_answer, inputs=[
33                 msg, chatbot], outputs=[msg, chatbot])
34
35             gr.Markdown("""提醒: <br>
36                 1. 初始化数据库时间可能较长，请耐心等待。
37                 2. 使用中如果出现异常，将会在文本输入框进行展示，请不要惊慌。 <br>
38                 """)
39 gr.close_all()
40 # 直接启动
41 demo.launch()

```

执行run\_gradio.py

```
1 | python run_gradio.py
```

### 3 远程演示

```
1 | ssh -CNg -L 7860:127.0.0.1:7860 root@ssh.intern-ai.org.cn -p 【开发机端口号】
```

