

Ứng dụng RAG xây dựng ứng dụng hỏi đáp - Project

Ngày 6 tháng 7 năm 2024

Ngày thực hiện:	07/07/2024
Người thực hiện:	Đinh Thị Tâm
Nguồn:	AIO2024 - Project RAG
Nguồn dữ liệu (nếu có):	Link of Data Sources
Từ khóa:	RAG Project
Người tóm tắt:	Đinh Thị Tâm

I. Tự luận

1. Câu 1:

(a) Code

```
1 # -*- coding: utf-8 -*-
2 """M1_RAG_Project_colab.ipynb
3
4 Automatically generated by Colab.
5
6 Original file is located at
7     https://colab.research.google.com/drive/12tbXy-meBuTsiuZ0mhk-t73kYisE-aKA
8 """
9
10 !pip install -q transformers==4.41.2
11 !pip install -q bitsandbytes==0.43.1
12 !pip install -q accelerate==0.31.0
13 !pip install -q langchain==0.2.5
14 !pip install -q langchainhub==0.1.20
15 !pip install -q langchain-chroma==0.1.1
16 !pip install -q langchain-community==0.2.5
17 !pip install -q langchain-openai==0.1.9
18 !pip install -q langchain_huggingface==0.0.3
19 !pip install -q chainlit==1.1.304
20 !pip install -q python-dotenv==1.0.1
21 !pip install -q pypdf==4.2.0
22 !npm install -g localtunnel
23 !pip install -q numpy==1.24.4
24
25
26
```

```
27 # Commented out IPython magic to ensure Python compatibility.
28 %%writefile app.py
29
30 import chainlit as cl
31 import torch
32
33 from chainlit.types import AskFileResponse
34
35 from transformers import BitsAndBytesConfig
36 from transformers import AutoTokenizer, AutoModelForCausalLM, pipeline
37 from langchain_huggingface import HuggingFaceEmbeddings
38 from langchain_huggingface.llms import HuggingFacePipeline
39
40 from langchain_community.chat_message_histories import ChatMessageHistory
41 from langchain_community.document_loaders import PyPDFLoader, TextLoader
42 from langchain.chains import ConversationalRetrievalChain
43 from langchain.memory import ConversationBufferMemory
44 from langchain_chroma import Chroma
45 from langchain_text_splitters import RecursiveCharacterTextSplitter
46 from langchain_core.runnables import RunnablePassthrough
47 from langchain_core.output_parsers import StrOutputParser
48 from langchain import hub
49
50 #text extraction startup
51 text_splitter = RecursiveCharacterTextSplitter ( chunk_size =1000 ,
52         chunk_overlap =100)
53 embedding = HuggingFaceEmbeddings ()
54
55 def process_file ( file : AskFileResponse ):
56     if file.type == "text/plain":
57         Loader = TextLoader
58     elif file.type == "application/pdf":
59         Loader = PyPDFLoader
60
61     loader = Loader (file.path )
62     documents = loader.load ()
63     docs = text_splitter.split_documents ( documents )
64     for i, doc in enumerate (docs):
65         doc.metadata ["source"] = f"source_{i}"
66     return docs
67
68 def get_vector_db ( file : AskFileResponse ):
69     docs = process_file ( file )
70     cl.user_session.set ("docs", docs )
71     vector_db = Chroma.from_documents ( documents =docs , embedding =
72         embedding )
73     return vector_db
74
75 def get_huggingface_llm(model_name : str ="lmsys/vicuna-7b-v1.5",
76         max_new_token : int = 512) :
77     nf4_config = BitsAndBytesConfig (
78         load_in_4bit =True ,
79         bnb_4bit_quant_type ="nf4",
80         bnb_4bit_use_double_quant =True ,
81         bnb_4bit_compute_dtype = torch.bfloat16
82     )
83     model = AutoModelForCausalLM.from_pretrained (
84         model_name,
85         quantization_config = nf4_config ,
86         low_cpu_mem_usage = True
```

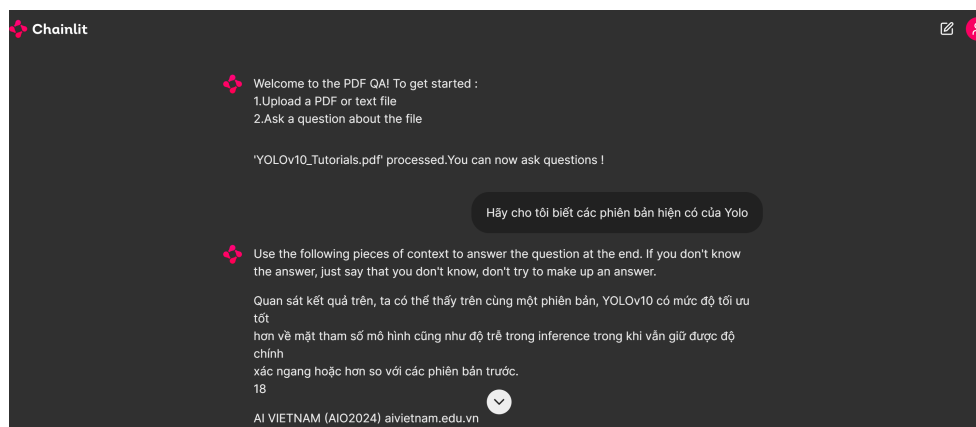
```

85     )
86     tokenizer = AutoTokenizer.from_pretrained ( model_name )
87     model_pipeline = pipeline (
88         "text-generation",
89         model =model ,
90         tokenizer = tokenizer ,
91         max_new_tokens = max_new_token ,
92         pad_token_id = tokenizer.eos_token_id ,
93         device_map ="auto"
94     )
95     llm = HuggingFacePipeline (
96         pipeline = model_pipeline ,
97     )
98     return llm
99 LLM = get_huggingface_llm ()
100
101 welcome_message = """Welcome to the PDF QA! To get started :
102 1.Upload a PDF or text file
103 2.Ask a question about the file
104 """
105
106 @cl.on_chat_start
107 async def on_chat_start ():
108     files = None
109     while files is None :
110         files = await cl.AskFileMessage (
111             content = welcome_message ,
112             accept=["text/plain","application/pdf"],
113             max_size_mb =20 ,
114             timeout =180 ,
115         ).send()
116         file = files [0]
117
118     msg = cl.Message ( content =f"Processing '{ file.name }'...",
119         disable_feedback = True )
120     await msg.send ()
121     vector_db = await cl.make_async ( get_vector_db )( file )
122
123     message_history = ChatMessageHistory ()
124     memory = ConversationBufferMemory (
125         memory_key ="chat_history",
126         output_key ="answer",
127         chat_memory = message_history ,
128         return_messages =True ,
129     )
130     retriever = vector_db.as_retriever ( search_type ="mmr",
131         search_kwargs ={ 'k': 3})
132
133     chain = ConversationalRetrievalChain.from_llm (
134         llm =LLM ,
135         chain_type ="stuff",
136         retriever = retriever ,
137         memory =memory ,
138         return_source_documents = True
139     )
140
141     msg.content = f"'{ file.name}' processed.You can now ask questions !"
142     await msg.update ()
143     cl.user_session.set ("chain", chain )
144

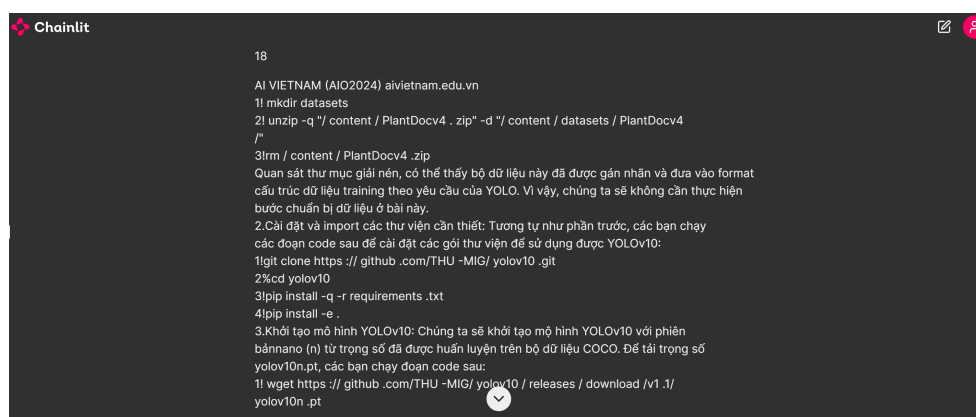
```

```
145 #function message
146 @cl.on_message
147 async def on_message(message: cl.Message):
148     chain = cl.user_session.get("chain")
149     cb = cl.AsyncLangchainCallbackHandler()
150     res = await chain.ainvoke(message.content, callbacks=[cb])
151     answer = res["answer"]
152     source_documents = res["source_documents"]
153     text_elements = []
154
155     if source_documents:
156         for source_idx, source_doc in enumerate(source_documents):
157             source_name = f"source_{source_idx}"
158             text_elements.append(
159                 cl.Text(content=source_doc.page_content,
160                        name=source_name)
161             )
162             source_names = [text_el.name for text_el in text_elements]
163
164             if source_names:
165                 answer += f"\nSources: {'', '.join(source_names)}"
166             else:
167                 answer += "\nNo sources found"
168
169     await cl.Message(content=answer, elements=text_elements).send()
170
171
172
173 #install pyngrok
174 !pip install pyngrok -q
175
176 from pyngrok import ngrok
177 !ngrok config add-authtoken 2ipnt7AAxR7r6ODdAbLJJJEZJkTV_6JHLnEYagYP72wkNiKRfL
178
179 public_url = ngrok.connect(8000).public_url
180 print(public_url)
181
182 #run app.py
183 !chainlit run app.py
184
```

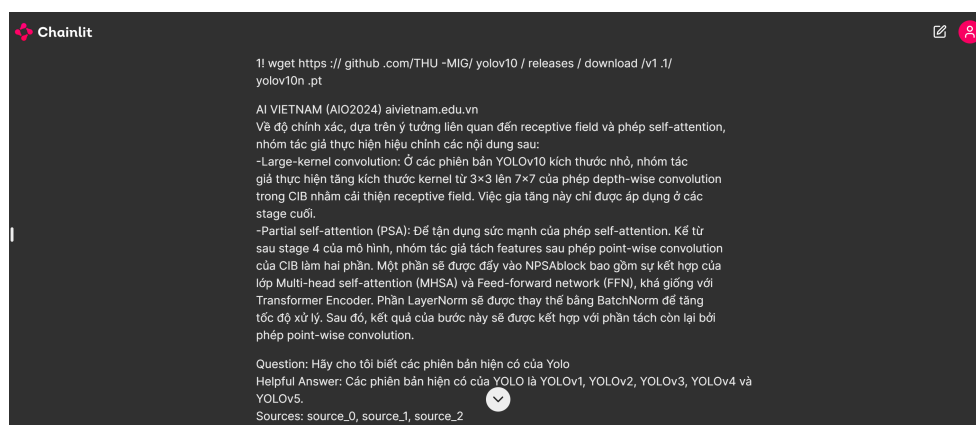
(b) Kết quả thực thi



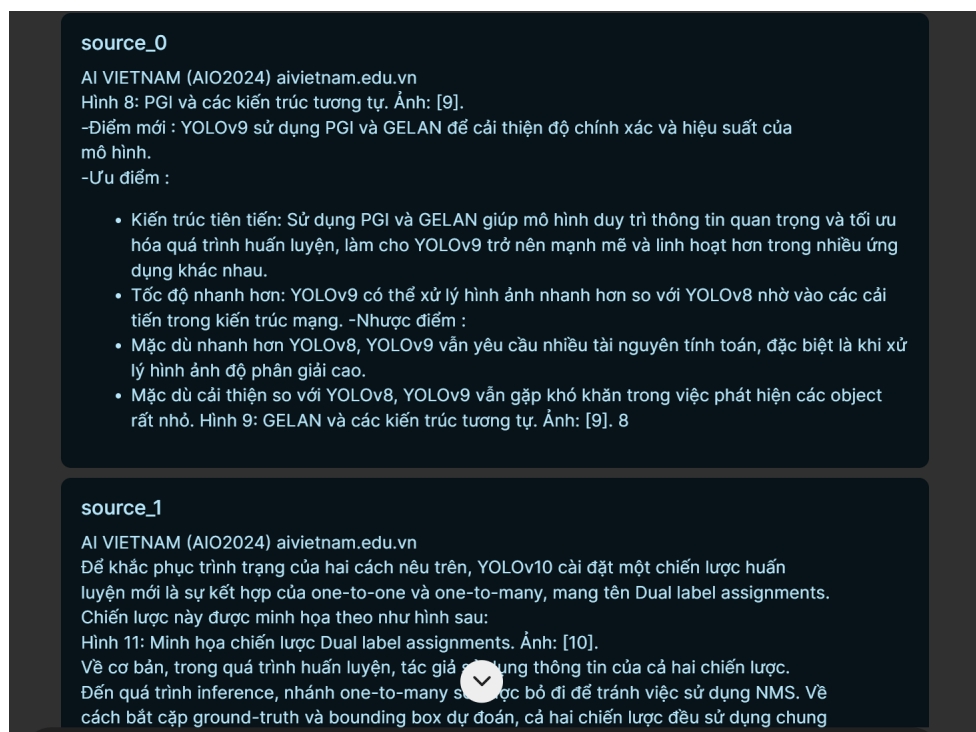
Hình 1



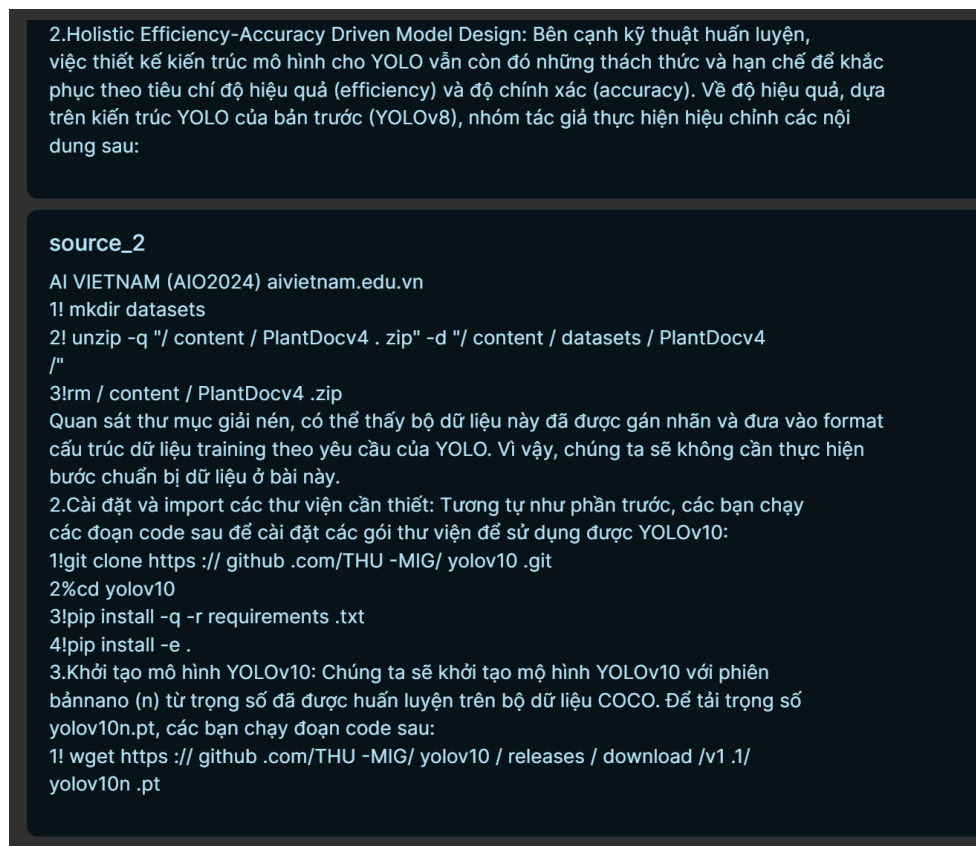
Hình 2



Hình 3



Hình 4



Hình 5

Hình ảnh thực thi chương trình