

Multi-agent LLMs

Dinh-Thang Duong, Nguyen-Thuan Duong và Quang-Vinh Dinh

Ngày 5 tháng 5 năm 2024

Phần I: Giới thiệu

Multi-agent LLMs là một hướng phát triển liên quan đến mô hình ngôn ngữ lớn. Trong đó, ta ứng dụng LLMs kết hợp với các **agents** có khả năng sử dụng một công cụ hoặc một chức năng bên ngoài nào đó, nhằm giải quyết bài toán của chúng ta. Ví dụ, ta có thể cài đặt cho LLMs sử dụng calculator để giải quyết các phép tính, từ đó sẽ đảm bảo kết quả đầu ra của mô hình.



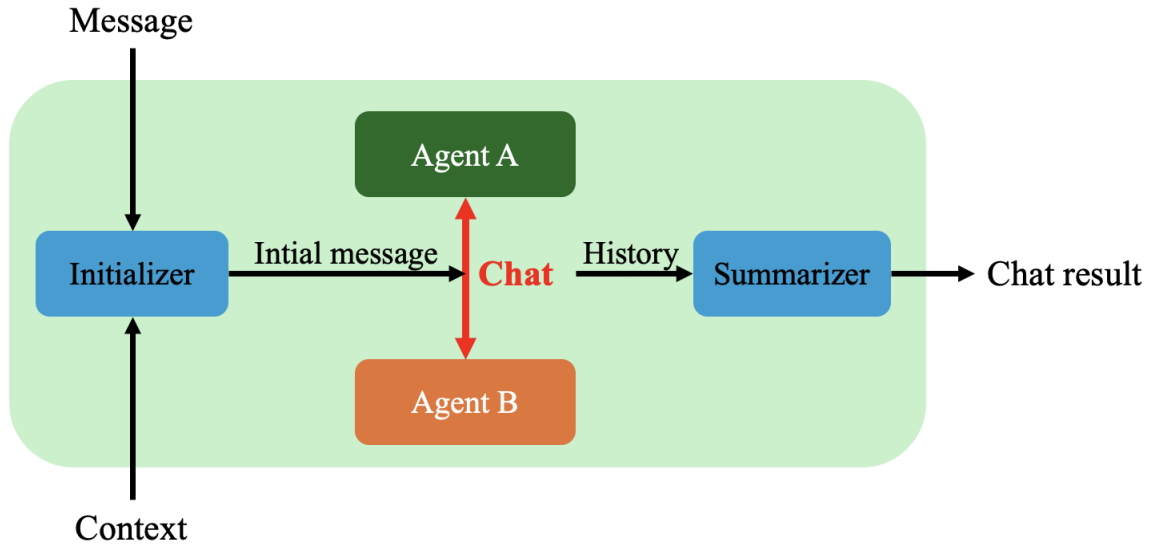
Hình 1: Minh họa về các agent trong môi trường ChatDev. Ảnh: [link](#).

Trong bài viết này, chúng ta sẽ xây dựng một ứng dụng về **Multi-agent LLMs** tạo sinh code đồng thời kiểm tra code thông qua trình thực thi. Mã nguồn được xây dựng bằng thư viện **Autogen**. Theo đó, ta sẽ xây dựng một ứng dụng có dạng như **Chatbot**. Ở đó, sẽ có **hai agent** tự giao tiếp với nhau gồm **AssistantAgent** và **UserProxyAgent**:

- **AssistantAgent** là một **LLM**, có vai trò tạo sinh code theo yêu cầu và kiểm tra lại code nếu cần.

- **UserProxyAgent** có vai trò gửi yêu cầu tạo sinh đến AssistantAgent và thực thi code nhận được để kiểm tra code đã chạy được hay chưa.

Thông qua hội thoại giữa hai agent này, ta sẽ có được một đoạn code theo yêu cầu đã được kiểm tra thông qua trình thực thi. Tổng quan, pipeline của project như sau:



Hình 2: Pipeline của project.

1. Từ yêu cầu về tạo sinh một đoạn code với chức năng nào đó từ người dùng, ta kết hợp với context có sẵn trong Autogen để hình thành nội dung khởi đầu của cuộc hội thoại giữa hai agent.
2. Với nội dung trên, ta tiến hành cho hai agent giao tiếp với nhau.
3. Từ nội dung chat giữa hai agent, ta trích ra nội dung cần thiết để làm kết quả của chương trình (ở đây là đoạn code được tạo sinh theo yêu cầu và kết quả thực thi).

Phần II: Cài đặt chương trình

Trong phần này, chúng ta sẽ tiến hành cài đặt nội dung của project. Mã nguồn được xây dựng trên hệ điều hành Ubuntu với GPU 24GB. Các bước thực hiện như sau:

1. **Tổ chức thư mục code:** Để mã nguồn trở nên rõ ràng nhằm phục vụ cho mục đích đọc hiểu code, chúng ta sẽ có ảnh minh họa tổ chức cây thư mục như sau:

```
multiagent_llms/
├── agent.py
├── app.py
├── chat.py
├── utils.py
├── OAI_CONFIG_LIST.json
└── requirements.txt
```

2. **Cập nhật file requirements.txt:** Ta liệt kê một số thư viện cần thiết để chạy được source code thông qua file requirements như sau:

```
1 pyautogen==0.2.27
2 gradio==4.29.0
```

3. **Cập nhật file agent.py:** Trong file này, ta định nghĩa các hàm khởi tạo các agents. Nội dung như sau:

- (a) **Import các thư viện và modules cần thiết:**

```
1 import autogen
2 from autogen import AssistantAgent, UserProxyAgent
3 from autogen.code_utils import extract_code
4
5 TIMEOUT = 60
```

- (b) **Khai báo biến config cho LLM (AssistantAgent):** AssistAgent của chúng ta là một LLM. Theo đó, đối với việc sử dụng GPT, AutoGen yêu cầu ta cung cấp một file .json có chứa API KEY và model version. Các bạn tạo một file tên **OAI_CONFIG_LIST.json** có nội dung sau:

```
1 [
2     {
3         "model": "gpt-3.5-turbo",
4         "api_key": "Type your API_KEY here"
5     }
6 ]
```

Quay lại file code agent.py, các bạn đọc file và import config như sau:

```
1 config_list = autogen.config_list_from_json("./OAI_CONFIG_LIST.json")
```

- (c) **Xây dựng hàm check điều kiện dừng hội thoại cho UserProxyAgent:** UserProxyAgent cần nhận được đoạn code được tạo sinh từ Agent còn lại để thực thi. Vì vậy, ta sẽ khai báo một hàm kết thúc lượt nhấn của UserProxyAgent nếu chưa nhận được code và ngược lại:

```

1 def _is_termination_msg(message):
2     if isinstance(message, dict):
3         message = message.get("content")
4         if message is None:
5             return False
6     cb = extract_code(message)
7     contain_code = False
8     for c in cb:
9         if c[0] == "python":
10             contain_code = True
11             break
12     return not contain_code

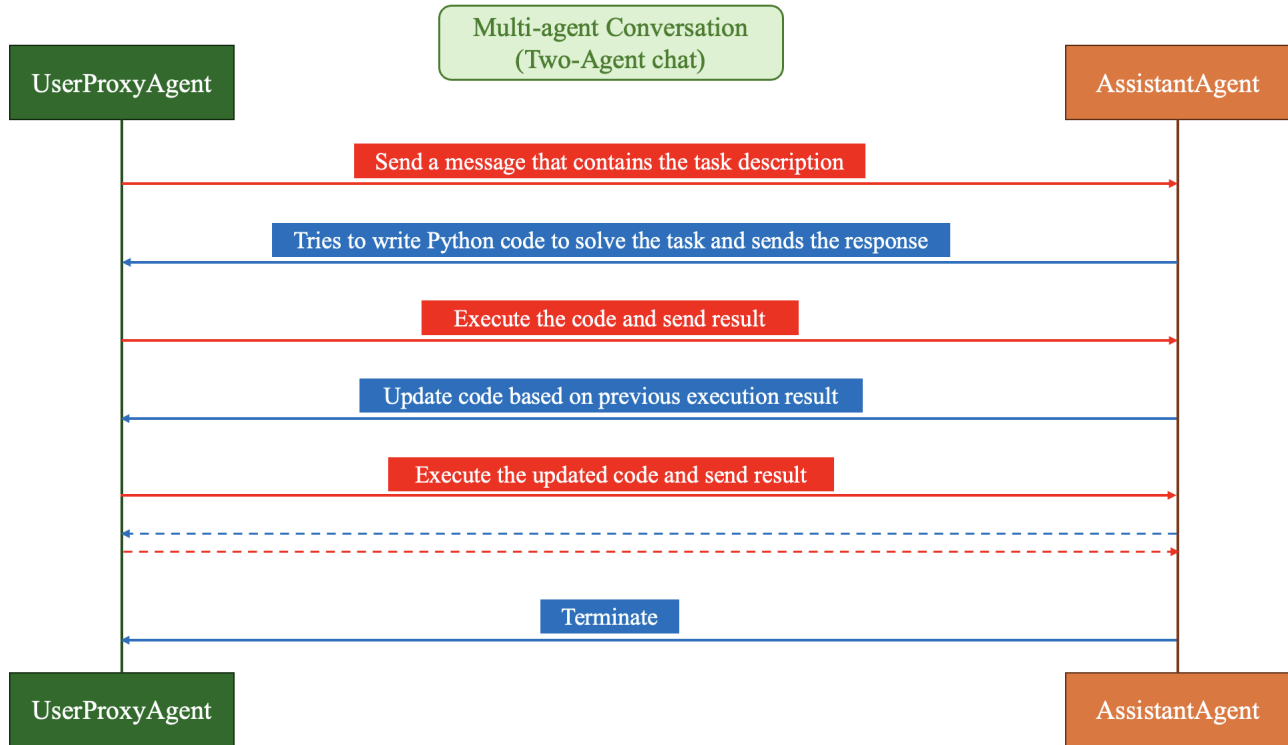
```

- (d) **Xây dựng hàm khởi tạo các agents:** Để thuận tiện trong việc triển khai code, ta xây dựng hàm initialize_agents() để tạo hai agent, sử dụng trong phần code của file khác trong source code:

```

1 def initialize_agents():
2     assistant = AssistantAgent(
3         name="assistant",
4         max_consecutive_auto_reply=5,
5         llm_config={
6             "timeout": TIMEOUT,
7             "config_list": config_list,
8         },
9     )
10
11     userproxy = UserProxyAgent(
12         name="userproxy",
13         human_input_mode="NEVER",
14         is_termination_msg=_is_termination_msg,
15         max_consecutive_auto_reply=5,
16         code_execution_config={
17             "work_dir": "coding",
18             "use_docker": False,
19         },
20     )
21
22     return assistant, userproxy

```



Hình 3: Minh họa về tương tác giữa hai agent - UserProxyAgent (Trình thực thi code) và AssistantAgent (LLM tạo sinh code).

4. **Cập nhật file `utils.py`:** Trong file này, ta định một class liên quan đến threading, sẽ được sử dụng ở file `chat.py`. Nội dung file `utils.py` như sau:

```

1 import sys
2 import threading
3
4 LOG_LEVEL = "INFO"
5 TIMEOUT = 60
6
7 class thread_with_trace(threading.Thread):
8     def __init__(self, *args, **keywords):
9         threading.Thread.__init__(self, *args, **keywords)
10        self.killed = False
11        self._return = None
12
13    def start(self):
14        self.__run_backup = self.run
15        self.run = self.__run
16        threading.Thread.start(self)
17
18    def __run(self):
19        sys.settrace(self.globaltrace)
20        self.__run_backup()
21        self.run = self.__run_backup
22
23    def run(self):
24        if self._target is not None:
25            self._return = self._target(*self._args, **self._kwargs)
26
  
```

```

27 def globaltrace(self, frame, event, arg):
28     if event == "call":
29         return self.localtrace
30     else:
31         return None
32
33 def localtrace(self, frame, event, arg):
34     if self.killed:
35         if event == "line":
36             raise SystemExit()
37     return self.localtrace
38
39 def kill(self):
40     self.killed = True
41
42 def join(self, timeout=0):
43     threading.Thread.join(self, timeout)
44     return self._return

```

5. **Cập nhật file chat.py:** Trong file này, ta định nghĩa các hàm xây dựng hộp hội thoại (chatbox) giữa hai agent. Từ đó, cho chúng tương tác qua lại và lấy về kết quả cuối cùng - lịch sử cuộc trò chuyện giữa hai agent, từ đó trích xuất được đoạn code và kết quả thực thi. Theo đó, ta triển khai như sau:

- (a) **Import các thư viện và module cần thiết:**

```

1 import os
2 from autogen import OpenAIWrapper
3 from agent import initialize_agents, config_list
4 from utils import *

```

- (b) **Khởi tạo agents:** Ta khởi tạo hai agent, UserProxyAgent và AssistantAgent. Đây là hai agent sẽ tương tác qua lại với nhau trong chatbox:

```

1 assistant, userproxy = initialize_agents()

```

- (c) **Xây dựng hàm chuyển đổi format lịch sử hội thoại:** Để đưa thông tin các cuộc hội thoại giữa 2 agent tại quá khứ trong chatbox vào model GPT, ta sẽ xây dựng hàm chuyển đổi thông tin lịch sử ta ghi nhận được sang format của OpenAI thông qua hàm `chat_to_oai_message()`:

```

1 def chat_to_oai_message(chat_history):
2     """Convert chat history to OpenAI message format."""
3     messages = []
4     if LOG_LEVEL == "DEBUG":
5         print(f"chat_to_oai_message: {chat_history}")
6     for msg in chat_history:
7         messages.append(
8             {
9                 "content": msg[0].split()[0] if msg[0].startswith("exitcode")
10                else msg[0],
11                 "role": "user",
12             }
13         )
14         messages.append({"content": msg[1], "role": "assistant"})
15     return messages

```

- (d) **Xây dựng hàm chuyển đổi ngược format lịch sử hội thoại:** Để hiển thị thông tin lịch sử lên giao diện của chúng ta sau khi chạy xong, ta cần lấy lại thông tin lịch sử từ GPT. Vì thông tin đang ở format của OpenAI, ta cần có hàm chuyển đổi ngược lại thành một list như ban đầu. Theo đó, ta khai báo hàm `oai_message_to_chat()`:

```

1 def oai_message_to_chat(oai_messages, sender):
2     """Convert OpenAI message format to chat history."""
3     chat_history = []
4     messages = oai_messages[sender]
5     if LOG_LEVEL == "DEBUG":
6         print(f"oai_message_to_chat: {messages}")
7     for i in range(0, len(messages), 2):
8         chat_history.append(
9             [
10                 messages[i]["content"],
11                 messages[i + 1]["content"] if i + 1 < len(messages) else "",
12             ]
13         )
14     return chat_history

```

- (e) **Xây dựng hàm khởi tạo đoạn hội thoại ban đầu:** Khi bắt đầu chương trình, ta cần khởi tạo môi trường chat cho hai agent. Ta có hàm `initiate_chat()` như sau:

```

1 def initiate_chat(config_list, user_message, chat_history):
2     if LOG_LEVEL == "DEBUG":
3         print(f"chat_history_init: {chat_history}")
4
5     if len(config_list[0].get("api_key", "")) < 2:
6         chat_history.append(
7             [
8                 user_message,
9             ]
10        )
11        return chat_history
12    else:
13        llm_config = {
14            "timeout": TIMEOUT,
15            "config_list": config_list,
16        }
17        assistant.llm_config.update(llm_config)
18        assistant.client = OpenAIWrapper(**assistant.llm_config)
19
20        if user_message.strip().lower().startswith("show file:"):
21            filename = user_message.strip().lower().replace("show file:", "").strip()
22            filepath = os.path.join("coding", filename)
23            if os.path.exists(filepath):
24                chat_history.append([user_message, (filepath,)])
25            else:
26                chat_history.append([user_message, f"File {filename} not found."])
27        ])
28
29        return chat_history
30
31    assistant.reset()
32    oai_messages = chat_to_oai_message(chat_history)
33    assistant._oai_system_message_origin = assistant._oai_system_message.copy()
34    assistant._oai_system_message += oai_messages
35
36    try:
37        userproxy.initiate_chat(assistant, message=user_message)
38        messages = userproxy.chat_messages
39        chat_history += oai_message_to_chat(messages, assistant)
40    except Exception as e:
41        chat_history.append([user_message, str(e)])

```

```

41     assistant._oai_system_message = assistant._oai_system_message_origin.copy
42     ()
43     if LOG_LEVEL == "DEBUG":
44         print(f"chat_history: {chat_history}")
45     return chat_history

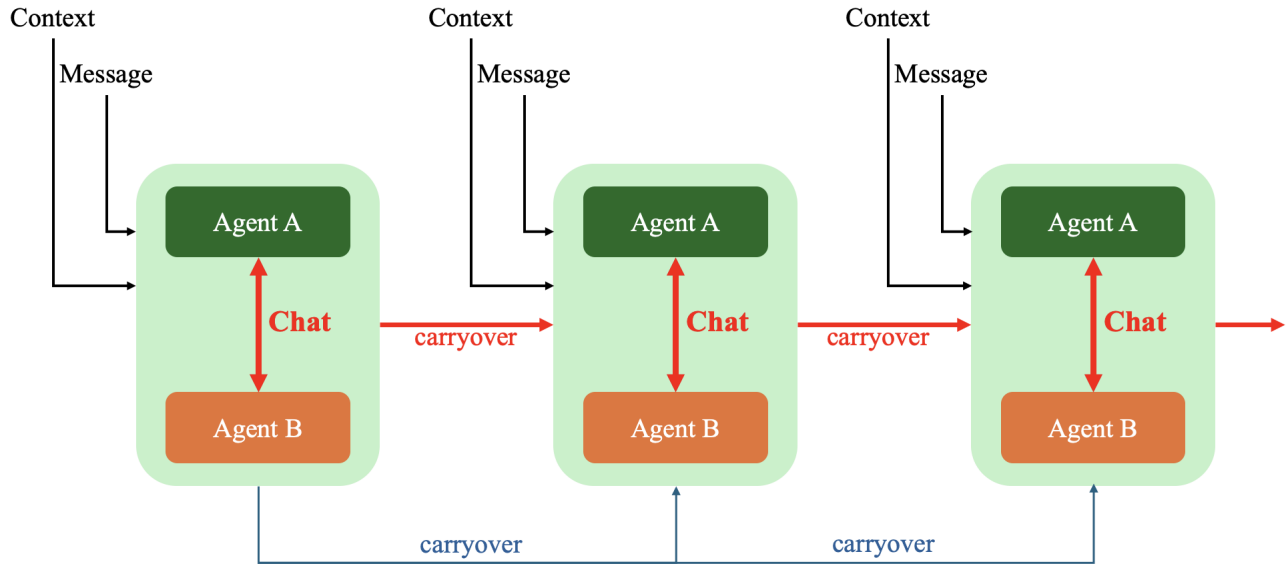
```

- (f) **Xây dựng hàm trả về kết quả đoạn hội thoại:** Ta xây dựng hàm trả về kết quả của chương trình. Kết quả chương trình của chúng ta là lịch sử đoạn hội thoại giữa hai agent (code của AssistantAgent và kết quả thực thi code của UserProxyAgent). Theo đó, hàm trả về kết quả được triển khai như sau:

```

1  def chatbot_reply_thread(input_text, chat_history, config_list):
2      """Chat with the agent through terminal."""
3      thread = thread_with_trace(target=initiate_chat, args=(config_list,
4      input_text, chat_history))
5      thread.start()
6      try:
7          messages = thread.join(timeout=TIMEOUT)
8          if thread.is_alive():
9              thread.kill()
10             thread.join()
11             messages = [
12                 input_text,
13                 "Timeout Error: Please check your API keys and try again
14                 later.",
15             ]
16         except Exception as e:
17             messages = [
18                 input_text,
19                 str(e) if len(str(e)) > 0 else "Invalid Request to OpenAI,
20                 please check your API keys.",
21             ]
22         return messages
23
24 def chatbot_reply(input_text, chat_history, config_list):
25     """Chat with the agent through terminal."""
26     return chatbot_reply_thread(input_text, chat_history, config_list)
27
28 def chat_respond(message, chat_history, model, oai_key, aoai_key, aoai_base):
29     chat_history[:] = chatbot_reply(message, chat_history, config_list)
30     if LOG_LEVEL == "DEBUG":
31         print(f"return chat_history: {chat_history}")
32     return ""

```

Hình 4: Minh họa việc trả về nội dung đoạn hội thoại (carryover) giữa hai agent. Các carryover theo thời gian trở thành lịch sử đoạn hội thoại làm context cho các lần chat khác.

6. **Cập nhật file app.py:** Cuối cùng, ta xây dựng một giao diện đơn giản để tương tác với chương trình, trực quan hóa hội thoại giữa hai agent cũng như kết quả của chương trình mà ta mong muốn. Theo đó, ta sẽ xây dựng giao diện bằng thư viện Gradio. Triển khai như sau:

- (a) **Import các thư viện và modules cần thiết:** Tại đây, ta sẽ import thư viện Gradio. Ta sẽ sử dụng giao diện chat cơ bản của Gradio là `ChatInterface`:

```
1 import gradio as gr
2 from gradio import ChatInterface
3
4 from chat import chat_respond
5
6 LOG_LEVEL = "INFO"
```

- (b) **Xây dựng block Gradio:** Để triển khai giao diện, ta xây dựng một block gradio với các thành phần trên giao diện như sau:

```
1 with gr.Blocks() as demo:
2
3     description = gr.Markdown("""
4         # Microsoft AutoGen
5         ## Multi-agent Conversation
6         """)
7
8     chatbot = gr.Chatbot(
9         [],
10        elem_id="chatbot",
11        bubble_full_width=False,
12        avatar_images=(
13            "../images/human.png",
14            "../images/autogen.png",
15        ),
16        render=False,
17        height=600,
```

```

18     )
19
20     txt_input = gr.Textbox(
21         scale=4,
22         show_label=False,
23         placeholder="Enter text and press enter",
24         container=False,
25         render=False,
26         autofocus=True,
27     )
28
29     chatiface = ChatInterface(
30         chat_respond,
31         chatbot=chatbot,
32         textbox=txt_input,
33         examples=[
34             ["write a python function to count the sum of two numbers?"],
35             ["what if the production of two numbers?"],
36             ["Plot a chart of the last year's stock prices of Microsoft,
37             Google and Apple and save to stock_price.png."],
38             ["show file: stock_price.png"],
39         ],
40     )

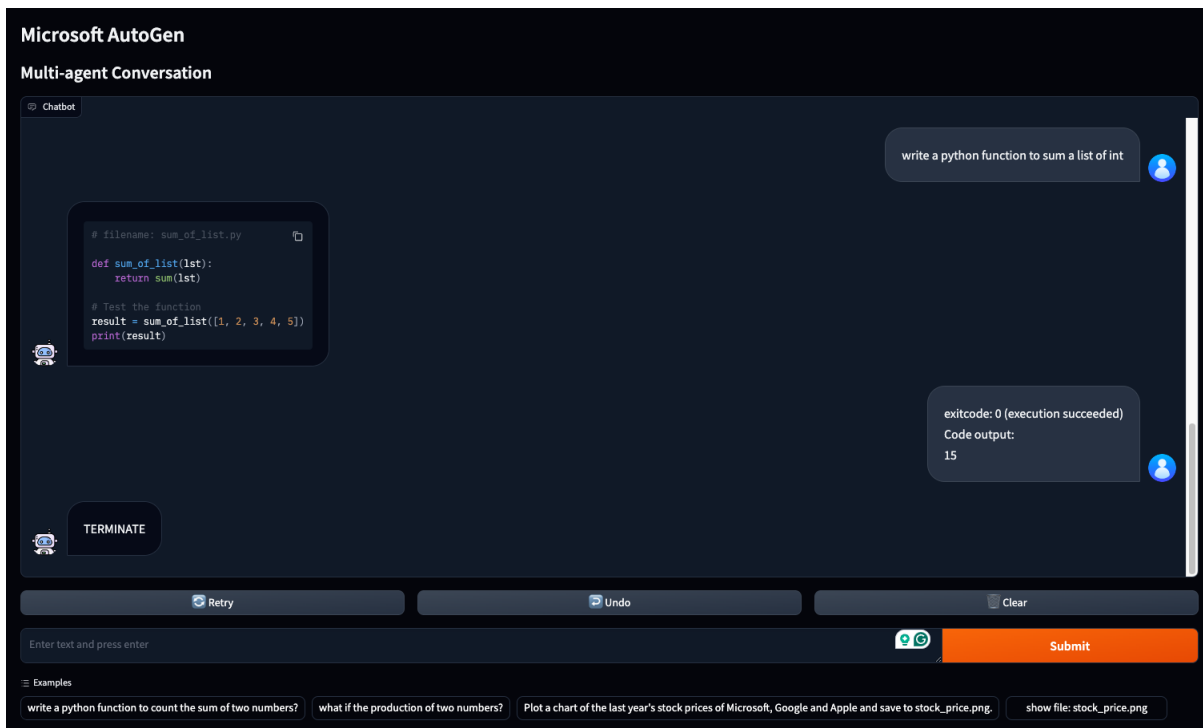
```

- (c) **Khởi động giao diện:** Cuối cùng, ta cài đặt lệnh khởi động giao diện khi chạy file app.py như sau:

```

1 if __name__ == "__main__":
2     demo.launch(share=False, server_name="0.0.0.0", server_port=7868)

```



Hình 5: Giao diện ứng dụng Multi-agent LLMs sau khi được triển khai và chạy thử một ví dụ.

Phần III: Câu hỏi trắc nghiệm

1. Trong Multi-agent LLMs, khái niệm agent có thể được hiểu như là?
 - (a) Là một Docker Container.
 - (b) Là một máy ảo (Virtual machine).
 - (c) Là một mô hình LLM.
 - (d) Là một thực thể (Entity) có khả năng nhận và gửi thông điệp (Message) đến các agent khác.
2. Đây là ưu điểm của việc sử dụng phối hợp của nhiều agent (Multi-agent)?
 - (a) Giải quyết một vấn đề nhanh hơn.
 - (b) Có khả năng giải quyết đa dạng các bài toán hơn.
 - (c) Giải quyết được các vấn đề, bài toán khó và phức tạp.
 - (d) Dễ dàng cài đặt.

3. Xét đoạn code dưới đây:

```
1 import os
2 from autogen import ConversableAgent
3
4 agent = ConversableAgent(
5     "chatbot",
6     llm_config={"config_list": [{
7         "model": "gpt-3.5-turbo",
8         "api_key": os.environ.get("OPENAI_API_KEY")
9     }]}
10 )
```

Agent trên được khởi tạo như là:

- (a) Một mô hình LLM.
 - (b) Một công cụ thực thi code.
 - (c) Một máy ảo.
 - (d) Công cụ lưu trữ lịch sử hội thoại.
4. Một agent có khả năng thực thi script (Code Executors) (ví dụ như python code) sẽ tồn tại rủi ro tiềm ẩn nào?
 - (a) Giống như chúng ta chạy code python nên không có rủi ro.
 - (b) LLM có thể sinh ra các script nhằm đánh cắp dữ liệu của máy HOST.
 - (c) Khiến máy HOST bị chậm do tiêu tốn tài nguyên CPU/GPU.
 - (d) Tiêu tốn không gian lưu trữ dữ liệu mà agent tải về.
 5. Đây **KHÔNG** là một giải pháp để giảm thiểu rủi ro cho một Code Executors Agent?
 - (a) Luôn luôn yêu cầu sự kiểm tra của con người trước khi cho agent thực thi code.
 - (b) Sử dụng máy ảo để thực thi code.
 - (c) Sử dụng công cụ Containerization (ví dụ Docker Container).
 - (d) Sử dụng Jupyter Notebook để thực thi code.
 6. Xét đoạn code dưới đây:

```

1 executor = DockerCommandLineCodeExecutor(
2     image="python:3.11-slim",
3     timeout=10,
4     work_dir=temp_dir.name,
5 )

```

Ý nghĩa của dòng thứ 2 là gì?

- (a) Khai báo tên của docker container.
- (b) Khai báo tên của docker image.
- (c) Khai báo tên docker image được dùng để thực thi.
- (d) Khai báo tên của Code Executor Agent.

7. Xét đoạn code dưới đây:

```

1 human_proxy = ConversableAgent(
2     "human_proxy",
3     human_input_mode="ALWAYS",
4 )

```

Ý nghĩa của dòng thứ 3 là gì?

- (a) Là message từ user truyền cho agent tại thời điểm khởi tạo.
- (b) Cho phép agent có khả năng thực thi code.
- (c) Luôn luôn yêu cầu message của con người trước khi agent thực thi.
- (d) Cho phép con người truyền message bất cứ khi nào.

8. Trong Autogen framework, đâu là cách để chấm dứt cuộc hội thoại, trò chuyện?

- (a) Dùng **max_turns** parameter trong hàm **initiate_chat**.
- (b) Dùng **max_consecutive_auto_reply** parameter khi khởi tạo agent.
- (c) Dùng **is_termination_msg** parameter khi khởi tạo agent.
- (d) Cả 3 cách trên.

9. Xét đoạn code dưới đây:

```

1 executor = DockerCommandLineCodeExecutor(
2     image="python:3.11-slim",
3     timeout=10,
4     work_dir=temp_dir.name,
5 )
6
7 userproxy = UserProxyAgent(
8     name="userproxy",
9     human_input_mode="NEVER",
10    is_termination_msg=lambda msg: "TERMINATE" in msg["content"],
11    max_consecutive_auto_reply=5,
12    code_execution_config={"executor": executor},
13 )
14
15 executor.stop()

```

Ý nghĩa của line thứ 15 là gì?

- (a) Chấm dứt cuộc hội thoại.
- (b) Ngắt hội thoại của **userproxy** agent.

- (c) Dùng docker container được dùng để thực thi code.
 - (d) Dùng việc thực thi code.
10. Vai trò của **Group Chat Manager** trong group chat là gì?
- (a) Chọn ra một Agent duy nhất để thực hiện task.
 - (b) Nhận task từ Boss Agent, sau đó phân tán task cho tất cả các Agent trong group.
 - (c) Thu thập kết quả từ các Agent sau đó đưa ra câu trả lời.
 - (d) Nhận task từ Boss Agent, chọn một Agent thực hiện task, lặp lại cho tới khi hoàn thành.

- *Hết* -