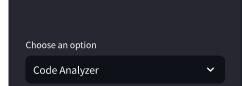
1/28/25. 11:09 PM app



Code Quality Analyzer by DeepSeek

Code Analyzer

Paste your code here:

import re
import streamlit as st
from langchain_community.llms import Ollama

Initialize the Ollama client
ollama_client = Ollama()

def chat(model, messages, stream):
 """Call the Ollama chat function."""
 return ollama_client.chat(model=model, messages=messages, stream=stream)

Set Streamlit page configuration (optional)
st set_page_config(page_title="DeepSeek Chat" layout="centered")

Analyze Code

Analysis Complete!

<think> Okay, I need to analyze this code provided by the user. Let's go through it step by step.
First, I see that they're importing some libraries: re, streamlit as st, and langchain_community.llms like
Ollama. That makes sense because they're using Streamlit with the Llama3.2 model via Ollama.

They initialized the Ollama client, which is a good start for streaming responses. Then, there's a function called chat that takes a model, messages, and stream. This function calls the Ollama client's chat method, so I'm guessing they want the assistant to think and process the user input.

Looking at display functions: display_message seems to show either "user" or "assistant" based on the role. That's standard. Then there are helper functions for formatting reasoning and displaying messages in the chat interface. These would be useful for organizing thoughts properly.

The format_reasoning_response function removes think tags, which is helpful for avoiding confusion. The display_assistant_message function handles formatting by replacing thought content with formatted text. This makes the output cleaner and more readable.

display_chat_history shows all previous messages, keeping only system messages. That's handy for tracking conversation flow.

process_thinking_phase processes the user input during the thinking phase, generating a response that includes thoughts before finalizing it. The process_response_phase handles the natural response after processing.

Handle_user_input is where the actual interaction happens between the user and assistant. They type messages, which are then passed to the chat function. Each message is handled by display_message or display_assistant_message based on role.

The main function sets up the Streamlit interface with some initializations, calls display_chat_history for setup, handles user input, and runs the handle_user_input function.

localhost:8501 1/3

1/28/25, 11:09 PM app

Now, thinking about possible areas for improvement. I notice that Ollama might not have access to the full API version used in the example. That's a common issue because sometimes you can't use older models on newer libraries. So suggesting to upgrade or adjust settings would be good.

Another point is how the user input is handled. The code shows typing messages, which works well but doesn't specify what gets stored or if there are any validations. Maybe adding validation for correct inputs could prevent errors in real-time.

Error handling and loading states are areas where improvement could happen. For example, using a loading state during thinking ensures users aren't stuck waiting. Also, ensuring that the response is handled properly can improve user experience by showing when the system is ready to respond.

I should also think about how messages are stored or if there's any concurrency issues with multiple threads or processes. Streamlit runs on the server, so perhaps using the worker thread for processing would be better but might complicate things.

Testing and documentation could be important too. The code as written works in isolation, so adding error handling, loading states, and documentation would make it more robust and easier to understand by others.

In summary, while this code is functional for a Streamlit app, there are areas like Ollama version compatibility, user input validation, error handling, message storage concurrency, and documentation that could be improved. Addressing these would make the implementation more solid.

The provided code sets up an interactive Streamlit application using Llama 3.2 via the Ollama library, allowing for streaming responses. Here's a summary of key features and areas for improvement:

Summary:

1. Setup:

- Uses streamlit as the frontend.
- Implements an Ollama client for streaming responses (ollama_client.chat().
- Handles user input directly in the chat interface.

2. Key Features:

- Processes user messages during the thinking phase and displays formatted responses.
- Includes functions to format reasoning steps and display messages with formatted content.
- Shows a centered layout initially but can be flipped for vertical mode (streaming).
- Supports message roles (user and assistant) in the chat interface.

3. Areas for Improvement:

- Ollama Version Compatibility: The code uses Ollama 1.x which might not support newer versions like Llama3.2.
 - Suggestion: Upgrade to a newer version (e.g., LlamaLlama) by adjusting API settings if necessary.

User Input Validation:

- The current setup does not validate user input, allowing any text entry without processing.
- Consider adding validation for the input using Streamlit's validation feature or error handling.

• Error Handling and Loading States:

- When processing user messages (e.g., thinking phase), the response might be pending.
 Adding a loading state during this phase would improve user experience.
- Use st.stop to stop streaming after responses are received.
- Message Storage and Concurrency:

1/28/25, 11:09 PM app

- The code processes all incoming messages in memory, which could cause issues with multiple threads or processes on Streamlit's server.
- Consider using the worker thread for more efficient streaming, though this might complicate implementation.
- Documentation and Testing:
 - Add detailed docstrings explaining API calls and error handling.
 - Implement proper testing to validate each function and ensure that responses are handled correctly.

Suggested Improvements:

1. Upgrade Ollama Version:

```
ollama --version
```

Replace the version in your code, e.g., change ollama_client = Ollama() to an upgraded version like ollama --version=latest.

- 2. **Implement Loading States During Thinking:** Modify display_message and other functions to include a loading state when processing user input.
- 3. Add Validation for User Input:
 - Use Streamlit's validation feature (if available) or add error checks in your code.
 - You could also implement message validation before sending the response.
- 4. **Improve Message Storage with Loading States**: Ensure that the application processes messages in a timely manner, avoiding long waiting periods.
- 5. **Enhance Documentation and Tests**: Add detailed docstrings for each API call and test each function to ensure it behaves as expected.

By addressing these areas, the implementation would be more robust, user-friendly, and maintainable.

localhost:8501 3/3