Why we used the design we used

The given code demonstrates a design for retrieving, processing, and reranking responses to queries about cricket matches using web scraping, natural language processing (NLP), and retrieval-augmented generation (RAG) techniques. Here are the key design choices:

1. Web Scraping: The code uses the requests library to scrape data from ESPNcricinfo for match details like scorecards, reports, and commentaries.

2. Document Loading: It loads the scraped content into memory using the WebBaseLoader from langchain_community to process the data for further analysis.

3. Text Splitting:  The code splits the loaded documents into smaller chunks using RecursiveCharacterTextSplitter to efficiently handle large text data.

4.  Embedding:  It generates embeddings for the text chunks using GPT4AllEmbeddings to convert the text into numerical representations suitable for machine learning models.

5. Retrieval and Generation:  The code uses a retrieval model (RetrievalQA) with a prompt to find relevant information from the embeddings and generate responses to user queries.

6.  Reranking:  It implements a simple reranking logic based on response length to reorder the initial responses and assesses the accuracy of the reranking process.

Overall, the design focuses on leveraging web scraping, NLP techniques, and a retrieval-based approach to provide accurate and relevant responses to user queries about cricket matches.


After evaluating the responses (either right or wrong) and computing the accuracy we arrived at the following results
Accuracy = number_of_correct_responses / total_queries
           = 29/50
            =0.58

Conclusion:

Based on the implemented design and the evaluation results, we can draw the following conclusions:

1. Design Effectiveness: The design effectively leveraged web scraping, natural language processing (NLP), and retrieval-augmented generation (RAG) techniques to retrieve, process, and rerank responses to queries about cricket matches.

2. Accuracy Assessment: The accuracy of the implemented RAG model, measured as the proportion of correct responses to total queries, was 58%. This means that out of 50 queries evaluated, 29 were correctly reranked by the model.

3. Room for Improvement: While a 58% accuracy rate is reasonable, there is room for improvement. Fine-tuning the retrieval and generation models, enhancing the reranking logic, and incorporating more diverse training data could potentially improve the accuracy of the system.

4. Real-World Applicability: The design demonstrates a practical application of RAG models for information retrieval and generation tasks. With further refinement, such systems could be deployed in real-world scenarios to provide accurate and relevant responses to user queries.

5. Scalability and Generalization: The design can be scaled to handle a larger number of queries and can be generalized to other domains beyond cricket matches, showcasing its versatility and applicability to diverse use cases.

In conclusion, while the implemented RAG model (both naive and advanced ) showed promising results, there is scope for enhancement to further improve its accuracy and applicability in real-world scenarios.