

Optimizing a Retrieval Augmented Generation (RAG) model can significantly enhance its performance, relevance, and accuracy. Here are two innovative techniques for optimizing the RAG model developed in Task 1:

1. Contextual Query Expansion

Overview: Contextual Query Expansion (CQE) involves enhancing the original query with additional context to improve the retrieval of relevant documents. This can be achieved using pre-trained language models to generate expansions or by utilizing historical query logs to identify common expansions.

Implementation:

- **Embedding-based Expansion:** Use embeddings from a pre-trained model to find semantically similar terms or phrases that can be appended to the original query.
- **Historical Data Analysis:** Analyze past queries and corresponding successful document retrievals to identify patterns and common expansions that can be applied to new queries.

Steps:

1. **Generate Query Embeddings:**
 - Use OpenAI's language model to create embeddings for the original query.
2. **Find Related Terms:**
 - Search a large corpus or use a similarity search on the Pinecone index to find terms or phrases with embeddings similar to the original query.
3. **Expand Query:**
 - Append these related terms or phrases to the original query to form an expanded query.
4. **Retrieve Documents:**
 - Use the expanded query to perform the retrieval in Pinecone.

2. Dynamic Retrieval Ranking with Feedback Loop

Overview: Dynamic Retrieval Ranking (DRR) involves adjusting the relevance ranking of retrieved documents based on user feedback or interaction history. This can be implemented through a feedback loop that continuously updates the retrieval model.

Implementation:

- **User Feedback:** Collect user feedback on the relevance of retrieved documents and use this data to update the ranking model.
- **Interaction History:** Analyze interaction logs to understand which documents are frequently accessed or considered relevant for similar queries.

Steps:

1. **Collect Feedback:**
 - Implement a mechanism to collect feedback from users regarding the relevance of retrieved documents.
2. **Update Retrieval Model:**
 - Use the feedback data to fine-tune the retrieval model, giving higher weights to documents marked as relevant.
3. **Re-rank Documents:**
 - Adjust the ranking of retrieved documents dynamically based on updated relevance scores.

Benefits:

1. **Contextual Query Expansion:**
 - Improves the relevance of retrieved documents by considering broader context.
 - Reduces the likelihood of missing relevant documents due to narrowly defined queries.
2. **Dynamic Retrieval Ranking with Feedback Loop:**
 - Continuously improves the retrieval process based on actual user interactions.
 - Adapts to changing information needs and preferences over time.