

**Zainab Hasan 27127 and Ikhlas Ahmed Khan 27096**

***Assignment 4***

***ITA RAG Assignment***

**20th April, 2025**

## 1. Platform Details

For this experimentation, Kaggle was used as the primary platform. All stages, including data processing, embedding generation, storage, and retrieval, were conducted within the Kaggle environment.

## 2. Data Details

The dataset consists of 221 pages from the book \*Harry Potter and the Philosopher's Stone\*. Each page was treated as a separate document. The book's PDF was processed, and the embedding of each document was generated and stored in a vector database to facilitate semantic retrieval.

## 3. Chunking Strategy

The experiment began with a chunk size of 250 tokens. During the process, both the chunk size and the overlap between chunks were gradually adjusted. This was done to preserve the context within the text and to ensure that no important information was lost when splitting the content into smaller parts. The overlap between chunks was especially important to maintain continuity and improve the quality of information being retrieved.

Through multiple trials, it was observed that increasing the chunk size led to improved performance in terms of both faithfulness (how true the answers were to the original text) and relevancy (how well the answers matched the questions). The performance kept improving until a chunk size of 1000 tokens with an overlap of 750 tokens. This specific combination gave the highest scores during testing.

However, increasing the chunk size or overlap beyond this point caused the scores to drop. This decline suggested that there is a global maximum—an optimal point for chunk size and overlap—after which adding more content in a single chunk may introduce too much information, causing the system to lose focus and reduce retrieval accuracy.

These observations helped identify the most effective chunking configuration for the given dataset and showed how important it is to balance chunk size and overlap for the best retrieval results.

## 4. Summarization Technique

The main summarization technique used in the experiment involved summarizing the response generated by the language model (LLM) by extracting the first 1–2 sentences of the generated answer. This aimed to simplify the content for quick evaluation and comparison while keeping it within the LLM's context limits.

However, the updated results using a different LLM revealed that this summarization step slightly decreased both the faithfulness and relevancy scores across all retrieval methods when compared to using the full response. For instance, the faithfulness score for Hybrid Search dropped from 0.9253 (without summarization) to 0.9059 (with summarization), and the relevancy score decreased from 0.8682 to 0.8694. Similar trends were observed for other retrieval strategies: for example, Reciprocal Rank Fusion scored 0.9206 faithfulness and 0.8721 relevancy with summarization, whereas without summarization, it improved to 0.9253 and 0.8671, respectively.

Overall, the non-summarized approach consistently produced better scores, indicating that providing the full, detailed response from the LLM—without truncation—results in more faithful and relevant answers. Based on these findings, it's evident that excluding summarization is the more effective strategy for this specific task.

## 5. Retrieval Techniques

In the retrieval system, five methods were employed to retrieve relevant documents for answering the query: Semantic Search, Keyword Search (BM25), Hybrid Search, MMR Search, and Reciprocal Rank Fusion. Each method operates differently and produces varying results based on how it retrieves information.

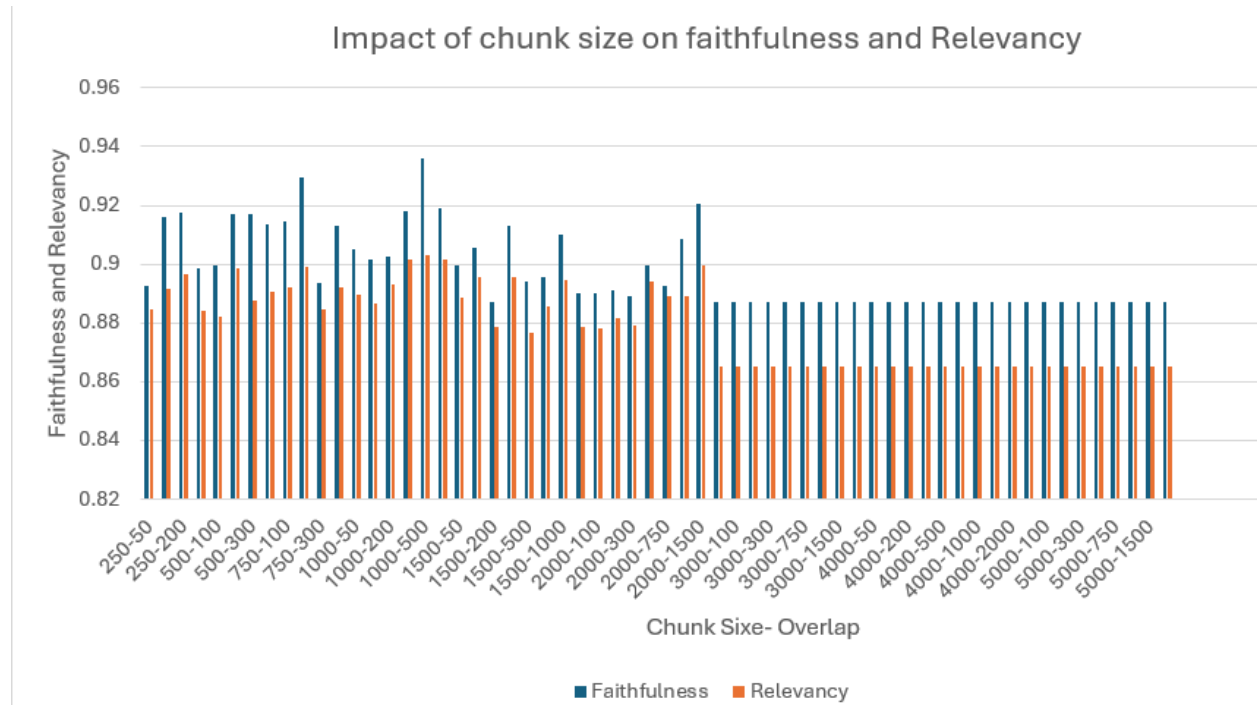


Figure 1: Semantic Search Performance Illustration

Semantic Search retrieves documents based on their semantic meaning, closely aligning the retrieved documents with the query and generally offering the highest performance. In contrast, Keyword Search (BM25) relies on exact keyword matching and, while effective, does not capture deeper semantic relationships between terms.

Hybrid Search combines both semantic and keyword-based retrievals using an ensemble approach, offering a good balance between the two retrieval techniques. MMR Search (Maximum Marginal Relevance) focuses on promoting diversity in the selection of documents, though this diversity can sometimes reduce the quality of retrieval.

Reciprocal Rank Fusion (RRF) merges results from both semantic and keyword-based searches, providing strong, balanced results by leveraging the strengths of each technique. Overall, Semantic Search remains the most effective standalone method, while Hybrid and RRF approaches offer robust alternatives by combining multiple retrieval strategies.

## 6. Model Evaluation

The evaluation shows a clear trend where models using powerful embedding architectures and instruction-tuned generative models consistently perform better in terms of both faithfulness and relevancy.

Among all the tested combinations, the best performing setup is the pair of `thenlper/gte-large` as the embedding model and `deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B` as the generative model. This combination achieves the highest faithfulness score of 0.9319 and relevancy score of 0.903, making it the most reliable and contextually accurate model overall.

The reason for this choice lies in the strengths of both components: `gte-large` is a state-of-the-art embedding model that captures nuanced semantic meaning effectively, and DeepSeek’s instruction-tuned model is designed to generate coherent, fact-based, and highly relevant responses—making them together the most robust choice for retrieval-augmented generation tasks.

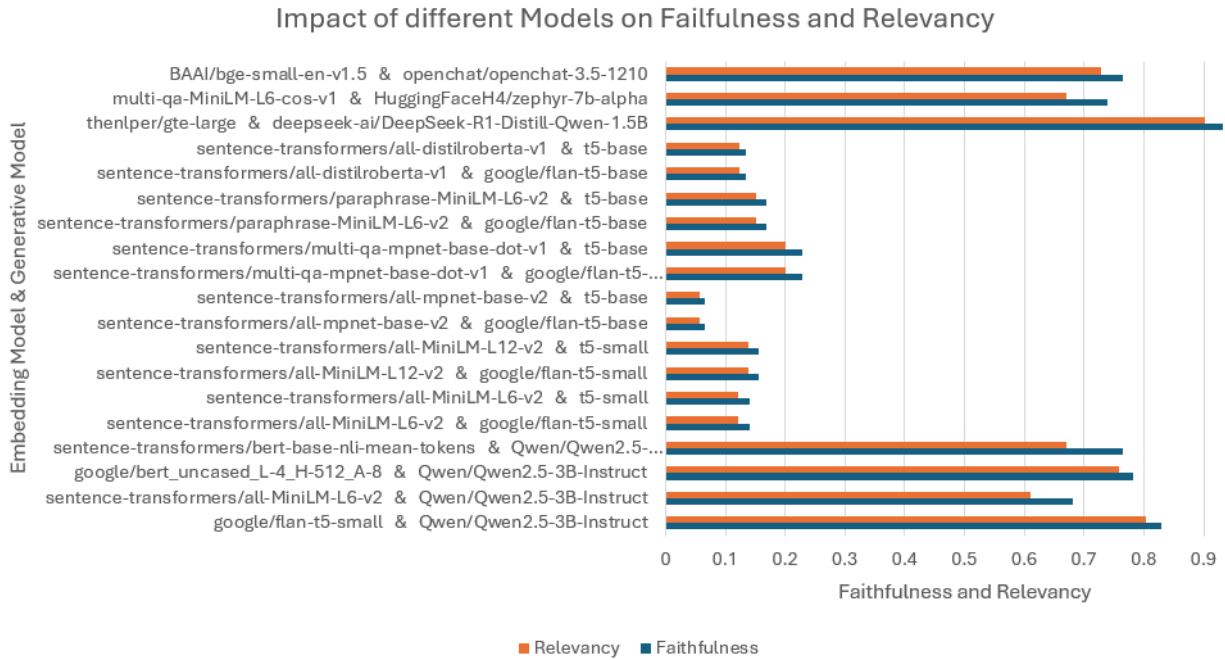


Figure 2: Impact of different models on Faithfulness and Relevance Scores

## 7. Performance

### Semantic Search

When comparing all the methods, Semantic Search stands out as the most effective. It gives the most accurate (faithful) and relevant answer to the question "Who is Hagrid?" with a faithfulness score of 0.9329 and a relevancy score of 0.8895, both of which are the highest among all methods. It's also fast, taking only 0.03 seconds to retrieve and 30.54 seconds to generate the answer. The final response was clear, context-aware, and avoided unnecessary repetition.

### Keyword Search (BM25)

On the other hand, Keyword Search (BM25) was the fastest overall with zero retrieval time, but it performed the worst in accuracy and relevancy, scoring the lowest on both metrics. It simply matched exact words and couldn't understand the deeper context of the question. The result was a bit generic and less meaningful.

### Hybrid Search

**Hybrid Search**, which mixes both keyword and semantic search, gave decent results. It had a good balance of accuracy and relevancy, but it suffered from repetitive phrasing, such as constantly describing Hagrid as "a very kind person." It was fast like semantic search but not as clear or relevant.

### MMR Search

MMR Search provided fairly accurate content with scores close to Hybrid, but it took the longest time to retrieve results—over 92 seconds, making it the slowest method overall with a total time of **122.08 seconds**. Although it introduced a bit of variety in the answer, it repeated itself and didn't add enough value to justify the wait.

## Reciprocal Rank Fusion

Reciprocal Rank Fusion, which combines the results from different ranking methods, also performed well in terms of accuracy and relevancy, but again, it included unnecessary repetition and wasn't as clear or concise as Semantic Search. It was fast but didn't offer any unique advantages over the other methods.

## Conclusion

In conclusion, **Semantic Search is the best method** among all. It offers the most faithful, most relevant, and quickest well-rounded answer with no confusing repetition. It understands the meaning behind the words and provides a cleaner, smarter response, making it the top choice for answering questions like "Who is Hagrid?"

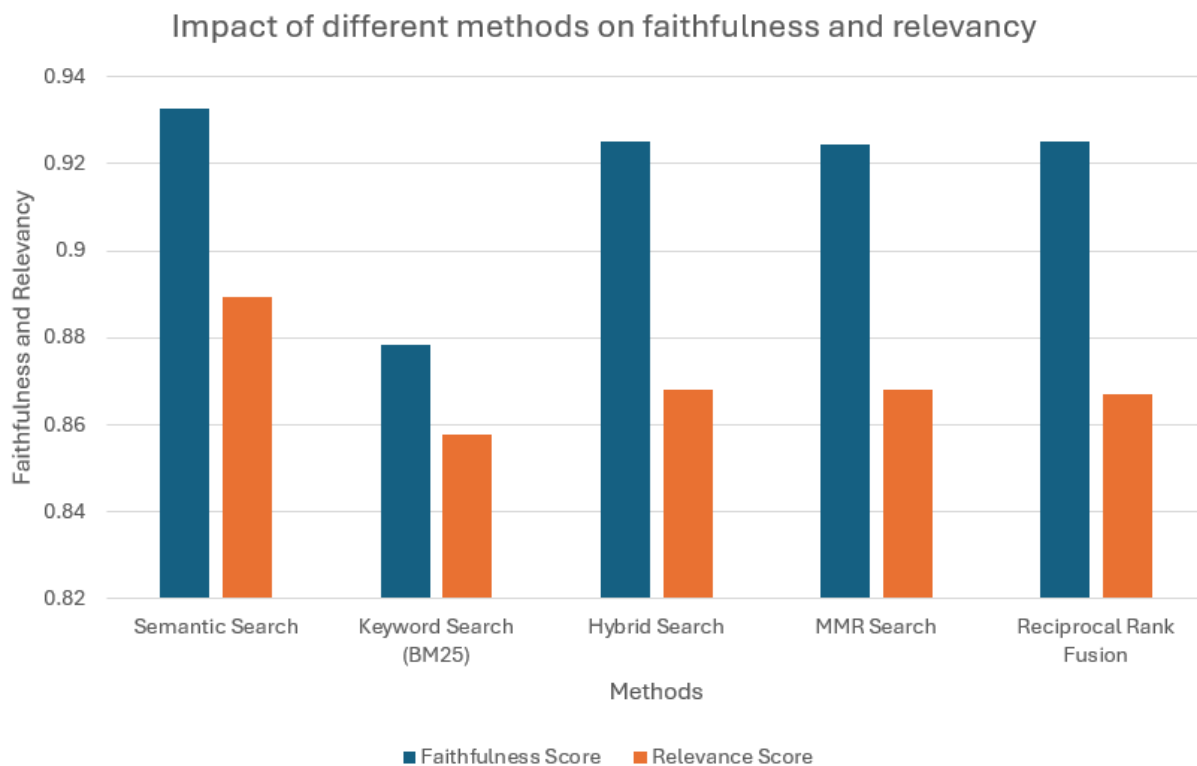


Figure 3: Comparison of Faithfulness and Relevance Scores for Different Search Methods

## Variation in Value of K

Among all values of  $k$ ,  $k = 7$  performed the best with the highest faithfulness score (0.9426) and a strong relevancy score (0.9079). This shows that retrieving 7 documents gave the most accurate and relevant answer overall, making it the optimal choice.

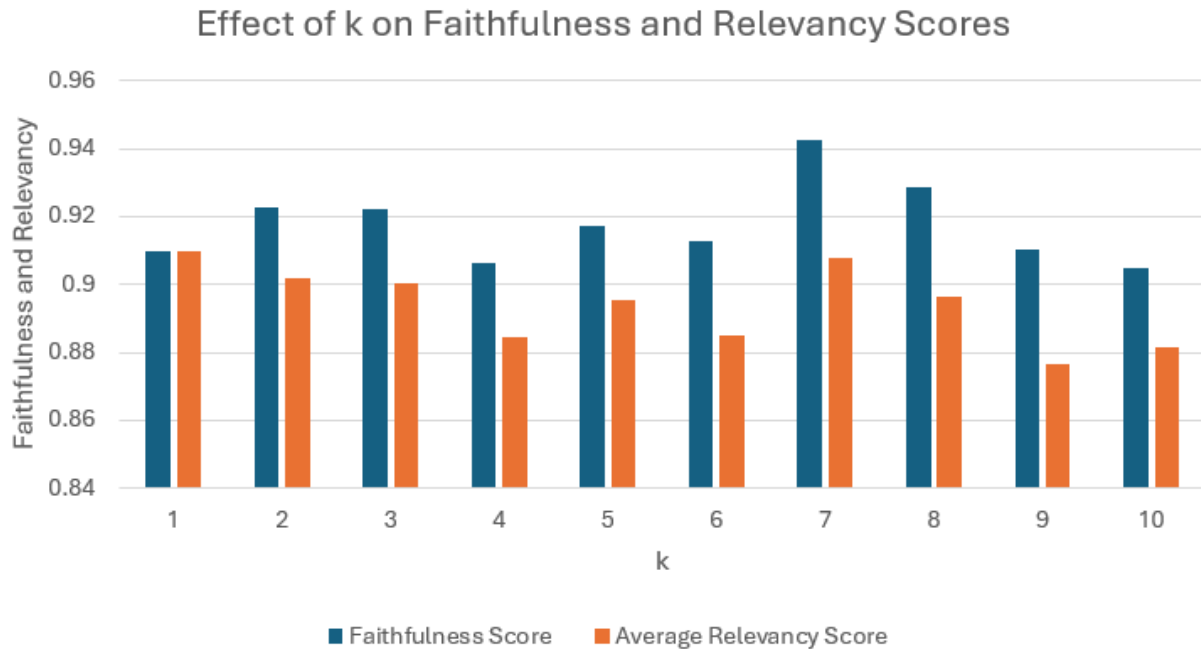


Figure 4: Effect of  $k$  on Faithfulness and Relevance Scores

## 8. Best Model Selection

Our best RAG setup used semantic search with  $k = 7$ , combining the embedding model **thenlper/gte-large** and the generative model **deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B**. This combination gave the most accurate and relevant answers by retrieving meaningful context and generating coherent responses. The **semantic search** technique helped match the question with the most related content, making this setup the most effective. The code for the best model is given below.

## 9. Best Model CODE

```

1 import os
2 from langchain_community.document_loaders import PyMuPDFLoader
3 from langchain.embeddings import HuggingFaceEmbeddings
4 from langchain.vectorstores import FAISS
5 from langchain.text_splitter import RecursiveCharacterTextSplitter
6 from langchain.schema.document import Document
7
8 from transformers import AutoModelForCausalLM, AutoTokenizer
9 from transformers import pipeline
10 import warnings
11 warnings.filterwarnings("ignore")
12 import textwrap
13
14 from ragas.metrics import faithfulness, answer_relevancy

```

```

15 from ragas.evaluation import evaluate
16
17 # Load and split PDF document
18 loader = PyMuPDFLoader("/kaggle/input/boooooook/Harry Potter Book.pdf")
19 pages = loader.load_and_split()
20 print(len(pages))
21
22 # Initialize embedding model
23 embedding_model = HuggingFaceEmbeddings(model_name="thenlper/gte-large")
24
25 text_splitter = RecursiveCharacterTextSplitter(
26     chunk_size=1000, # Adjust chunk size as needed
27     chunk_overlap=750 # Ensures context continuity
28 )
29
30 # Apply chunking to pages
31 chunks = text_splitter.split_documents(pages)
32 len(chunks)
33
34 # Create FAISS vector database
35 #vectordb = FAISS.from_documents(pages, embedding_model)
36 vectordb = FAISS.from_documents(chunks, embedding_model)
37
38
39 # Save FAISS index to disk for later use
40 vectordb.save_local("faiss_index")
41
42 # Check the number of stored documents
43 print(f"Number of documents in the vector store: {vectordb.index.ntotal}")
44
45 # Query processing
46 question = "Who was Hagrid?"
47 retriever = vectordb.as_retriever(search_kwargs={"k": 8})
48 docs = retriever.get_relevant_documents(question)
49
50 # Print results
51 for i, doc in enumerate(docs, 1):
52     page_number = doc.metadata.get('page', 'Unknown')
53     print(f"Document {i} - Page {page_number} - Score: {doc.metadata.get('score', 'N/A')}")
54     print(doc.page_content[:500]) # Print first 500 characters of each result
55     print("-" * 80)
56
57 model_name = "deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B"
58 model = AutoModelForCausalLM.from_pretrained(
59     model_name,
60     device_map="auto", #device_map='cuda'
61     torch_dtype="auto",
62     trust_remote_code=True,
63 )
64 tokenizer = AutoTokenizer.from_pretrained(model_name)
65
66 print(model.dtype)
67 total_params = sum(p.numel() for p in model.parameters())
68 print(f"Total Parameters: {total_params / 1e6} million")
69 memory_footprint = total_params * 2 / (1024 ** 2) # Convert to MB
70 print(f"Estimated Memory Footprint: {memory_footprint:.2f} MB")
71
72 # Create a pipeline
73 generator = pipeline(
74     "text-generation",
75     model=model,
76     tokenizer=tokenizer,
77     return_full_text=False,
78     max_new_tokens=5000,
79     do_sample=False
80 )
81
82 # Extract page content and metadata properly

```

```

83 def format_response(doc):
84     return f"Page {doc.metadata.get('page', 'Unknown')}: {doc.page_content.strip()}"
85
86 # Handle cases where fewer than 3 results are returned
87 retrieved_responses = [format_response(doc) for doc in docs[:7]]
88 while len(retrieved_responses) < 7:
89     retrieved_responses.append("No relevant information.") # Fill missing slots
90
91 # Construct the RAG prompt
92 prompt = f"""
93 You are an AI assistant tasked with answering questions based on retrieved knowledge.
94
95 ### **Retrieved Information**:
96 {retrieved_responses[0]}
97 {retrieved_responses[1]}
98 {retrieved_responses[2]}
99 {retrieved_responses[3]}
100 {retrieved_responses[4]}
101 {retrieved_responses[5]}
102 {retrieved_responses[6]}
103
104 ### **Question**:
105 {question}
106
107 ### **Instructions**:
108 - Integrate the key points from all retrieved responses into a **cohesive, well-structured
109   answer**.
110 - If the responses are **contradictory**, mention the different perspectives.
111 - If none of the retrieved responses contain relevant information, reply:
112   **"I couldn't find a good response to your query in the database."**
113
114
115 # Use Qwen2.5 3B with the correct message format
116 messages = [
117     {"role": "user", "content": prompt}
118 ]
119
120 # Generate output using the model
121 output = generator(messages)
122
123 # Print formatted response
124 print(textwrap.fill(output[0]["generated_text"], width=80))
125 import numpy as np
126 from sklearn.metrics.pairwise import cosine_similarity
127
128 # Function to get embeddings for a list of texts (documents)
129 def get_embeddings_for_documents(documents, model):
130     # Use embed_documents to get embeddings for a list of documents
131     return model.embed_documents(documents)
132
133 # Function to get embedding for a single text (generated text)
134 def get_embedding_for_query(query, model):
135     # Use embed_query to get embedding for a single query
136     return model.embed_query(query)
137
138 # Get embeddings for the retrieved documents and the generated text
139 retrieved_documents = [doc.page_content for doc in docs[:8]]
140 retrieved_embeddings = get_embeddings_for_documents(retrieved_documents, embedding_model)
141 generated_text_embedding = get_embedding_for_query(output[0]["generated_text"],
142     embedding_model)
143
144 # Calculate cosine similarity between each retrieved document and the generated text
145 similarities = []
146 for doc_embedding in retrieved_embeddings:
147     similarity = cosine_similarity([doc_embedding], [generated_text_embedding])
148     similarities.append(similarity[0][0])

```



```
149 # Print out cosine similarity scores for faithfulness and relevancy
150 print("Cosine Similarities (Relevancy Scores):")
151 for i, similarity in enumerate(similarities, 1):
152     print(f"Document {i}: Similarity = {similarity:.4f}")
153
154 # Calculate the Faithfulness Score
155 # Faithfulness Score: This is the maximum cosine similarity between the generated response
    and any document
156 faithfulness_score = max(similarities)
157 print(f"\nFaithfulness Score: {faithfulness_score:.4f}")
158
159 # Calculate the average relevancy score (mean cosine similarity)
160 average_relevancy_score = np.mean(similarities)
161 print(f"Average Relevancy Score: {average_relevancy_score:.4f}")
162
163 print(f"Average Faithfulness Score")
```

Listing 1: Best Model Code