

# Advanced Retrieval-Augmented Generation (RAG) System for Medical Question Answering on Cardiac Anatomy

## Introduction

Large Language Models (LLMs) demonstrate impressive generative capabilities at text generation, but in specialised areas like medicine, they can produce errors or lack necessary detail. This project addresses those gaps by developing and evaluating an advanced Retrieval-Augmented Generation (RAG) system for cardiac anatomy. By integrating hybrid retrieval, reranking, step-back prompting, and thorough evaluation, the system delivers more accurate, detailed, and reliable answers than a standalone LLM.

We install essential libraries to enable text chunking, semantic and keyword-based retrieval, embedding generation, and large language model inference, all key steps in building a robust RAG pipeline.

```
!pip install faiss-cpu sentence-transformers transformers rank_bm25 langchain --quiet
import nltk
nltk.download('stopwords')

[ ] ## This section loads all necessary Python libraries for building a Retrieval-Augmented Generation (RAG) pipeline
import json
import numpy as np
import pandas as pd
import torch
import faiss
from sentence_transformers import SentenceTransformer
from transformers import AutoTokenizer, AutoModelForSeq2SeqLM, pipeline, AutoModelForSequenceClassification
from rank_bm25 import BM25Okapi
from langchain.text_splitter import RecursiveCharacterTextSplitter
```

## Domain Selection & Dataset Sourcing

### Domain Selection

This RAG system focuses on Medical Question Answering (Q&A) with a particular emphasis on cardiac anatomy. Cardiac anatomy was chosen for several reasons:

- The domain is central to both medical education and clinical practice, requiring precision and depth.

- Errors in anatomical knowledge can have direct negative consequences in a healthcare setting.
- The field is highly specialised, and generic LLMs often fail to deliver the necessary detail or may hallucinate facts when faced with less common questions.

## Dataset

A PDF textbook on heart\_anatomy.pdf, was used as the authoritative domain source. Extracted via PyMuPDF, the content includes detailed descriptions of anatomical features, spatial relationships among heart structures, physiological terminology, and contextual explanations suited for educational or clinical use, not typically found in LLM training data.

## Justification

Medical Q&A requires verifiable, context-rich responses. LLMs often hallucinate or generate plausible-sounding but incorrect answers when facing uncommon medical questions. For example, when asked, "What forms the base of the heart?" or "What happens if the mitral valve is damaged?", LLMs may return vague or incorrect responses, potentially misinforming users. RAG significantly reduces these risks by grounding answers in verified, domain-specific content.

Grounding is particularly vital in medicine, where minor errors can cause major misunderstandings. Cardiac anatomy's spatial complexity also makes it a strong candidate for future RAG enhancements, such as integrating diagrams or images via multi-modal systems.

```
!pip install gdown
!pip install pymupdf
# Download the file from Google Drive using gdown
!gdown --id 1Ek7g9yGFB3iusI-ytqsmwi2foBvnPsGL -O heart_anatomy.pdf

import fitz # PyMuPDF

pdf_filename = "heart_anatomy.pdf"
doc = fitz.open(pdf_filename)
docs_read = ""
for page in doc:
    docs_read += page.get_text()

print("First 1000 characters of extracted text:\n")
print(docs_read[:1000])
```

## Data Preprocessing

**Chunking:** The input text is split into overlapping, semantically coherent chunks using RecursiveCharacterTextSplitter. This approach manages long documents efficiently and preserves context across chunk boundaries, resulting in text segments optimised for retrieval and downstream processing.

```
# Semantic Chunking
splitter = RecursiveCharacterTextSplitter(chunk_size=350, chunk_overlap=30)
chunks = splitter.create_documents([docs_read])
chunk_texts = [doc.page_content for doc in chunks]
```

**Embedding:** Text chunks are converted into vector embeddings using the "sentence-transformers/all-mpnet-base-v2" model. These embeddings capture semantic meaning, allowing for effective retrieval based on similarity rather than just keyword matching. The all-mpnet-base-v2 model is a top-performing choice for semantic search, outperforming many other public embedding models in retrieval tasks.

```
# Embedding Model
embed_model = SentenceTransformer("sentence-transformers/all-mpnet-base-v2")
embeddings = embed_model.encode(chunk_texts, show_progress_bar=True)
```

**Vector Storage:** FAISS is used to create an index of the embeddings. FAISS enables fast and efficient similarity searches, which are critical for quickly retrieving relevant chunks in response to a user query.

```
# FAISS Index
dimension = embeddings[0].shape[0]
index = faiss.IndexFlatL2(dimension)
index.add(np.array(embeddings).astype('float32'))
```

## Retrieval and Reranking

### Retrieval

The system uses a hybrid retrieval approach, combining BM25 for keyword-based search and FAISS for semantic similarity search. BM25 quickly identifies text chunks containing query keywords, while FAISS finds passages that are semantically relevant, even if exact terms do not

match. This dual strategy improves both precision and recall, ensuring that both direct matches and related content are retrieved for further processing.

```
# BM25 Index
tokenized_corpus = [doc.split() for doc in chunk_texts]
bm25 = BM25Okapi(tokenized_corpus)
```

**Reranking:** Retrieved chunks are further refined using the "BAAI/bge-reranker-base" model, which scores and prioritises passages based on their relevance to the query. This cross-encoder reranker uses deep cross-attention for more accurate relevance assessment, ensuring the most contextually appropriate information is selected for answer generation.

```
# Reranker
rerank_tokenizer = AutoTokenizer.from_pretrained("BAAI/bge-reranker-base")
rerank_model = AutoModelForSequenceClassification.from_pretrained("BAAI/bge-reranker-base")
rerank_model.eval()
```

## Rerank Function

This function takes a user query and a list of passages, then uses the reranker model to score and sort the passages by relevance. It returns the top k most relevant passages to improve the search results.

```
# Rerank Function
def rerank_passages(query, passages, top_k=3):
    pairs = [[query, p] for p in passages]
    inputs = rerank_tokenizer(pairs, padding=True, truncation=True, return_tensors="pt")
    with torch.no_grad():
        scores = rerank_model(**inputs).logits.squeeze(-1)
        sorted_indices = torch.argsort(scores, descending=True)
        return [passages[i] for i in sorted_indices[:top_k]]
```

## Generation

This code uses the pre-trained google/flan-t5-xl model and its tokenizer to generate text responses. Flan-T5 XL is a powerful sequence-to-sequence transformer, well-suited for context-grounded answer generation in RAG due to its strong instruction-following ability and reliable performance.

```
# Generator
gen_model_name = "google/flan-t5-xl"
tokenizer = AutoTokenizer.from_pretrained(gen_model_name)
model = AutoModelForSeq2SeqLM.from_pretrained(gen_model_name)
text_gen = pipeline("text2text-generation", model=model, tokenizer=tokenizer)
```

## Step-back Prompting

This function helps improve user queries by rewriting or breaking down the original question into a simpler, clearer version. It uses the text generation model to create a refined question that is easier to understand and retrieve accurate medical information.

Based on recent advancements in LLM-aided retrieval, step-back prompting mimics the reasoning a human expert might use to simplify a question before looking it up.

```
# Step-back Prompting
def generate_step_back_query(original_query):
    prompt = f"""You are a medical assistant helping clarify medical questions.
    Given the following user question, rewrite it or break it into a clearer sub-question that helps retrieve the right medical information.

    ORIGINAL QUESTION:
    {original_query}

    CLARIFIED or SIMPLIFIED QUESTION:"""
    return text_gen(prompt, max_new_tokens=100)[0]["generated_text"].strip()
```

## Hybrid Retrieval with Step-Back

This function improves search by first rewriting the user query for clarity (step-back). Then it performs two types of retrieval: vector similarity search (FAISS) and keyword search (BM25). It combines results from both methods and uses the reranker to select the top most relevant passages.

```
# Hybrid Retrieval with Step-Back
def hybrid_retrieve_stepback(query, k=5):
    clarified_query = generate_step_back_query(query)
    print("Step-Back Reformulated Query:", clarified_query)
    query_vector = embed_model.encode([clarified_query])
    _, faiss_indices = index.search(np.array(query_vector), k)
    vector_hits = [chunk_texts[i] for i in faiss_indices[0]]
    bm25_hits = bm25.get_top_n(clarified_query.split(), chunk_texts, n=10)
    combined = list(set(bm25_hits + vector_hits))
    return rerank_passages(clarified_query, combined, top_k=k)
```

The system removes generic words from queries to focus on key terms, then checks that at least two of these appear in the retrieved context. If not, it declines to answer. This ensures the RAG

model only responds when reliable information is available, increasing accuracy and trustworthiness.

```
from nltk.corpus import stopwords

stop_words = set(stopwords.words('english'))

domain_words = set(['heart'])

def get_meaningful_tokens(query):
    return [w for w in query.lower().split() if w not in stop_words and w not in domain_words]

def generate_rag_standard_answer(query, similarity_threshold=0.3):
    query_tokens = query.lower().split()
    key_tokens = get_meaningful_tokens(query)
    all_scores = bm25.get_scores(query_tokens)
    query_vector = embed_model.encode([query])
    _, faiss_indices = index.search(np.array(query_vector), 5)
    vector_hits = [chunk_texts[i] for i in faiss_indices[0]]
    bm25_hits = bm25.get_top_n(query_tokens, chunk_texts, n=10)
    combined = list(set(bm25_hits + vector_hits))
    reranked = rerank_passages(query, combined, top_k=5)
    context = "\n".join(reranked)
    num_key_tokens_in_context = sum(token in context.lower() for token in key_tokens)
    if num_key_tokens_in_context < min(2, len(key_tokens)):
        return "Sorry, the provided documents do not contain relevant information for your query."
    prompt = f"""
You are MedicalBot, a medical assistant. Answer the user's question strictly using ONLY the information in CONTEXT below.
If the context does NOT contain an answer, respond: "Sorry, the provided documents do not contain relevant information for your query."
Do NOT use your own knowledge or make up information. Cite medical facts clearly.

CONTEXT:
{context}

QUESTION:
{query}

ANSWER:
"""
    return text_gen(prompt, max_new_tokens=350)[0]["generated_text"].strip()
```

## Step-Back Prompting in RAG

This function rewrites user queries for clarity, retrieves relevant passages, and checks for at least two key terms in the context. If these are missing, it declines to answer. This process helps the RAG model handle complex questions accurately and avoids unsupported answers.

```
def generate_rag_answer_stepback(query, similarity_threshold=0.3):
    stepback_query = generate_step_back_query(query)
    print("Step-Back Reformulated Query:", stepback_query)
    stepback_tokens = stepback_query.lower().split()
    key_tokens = get_meaningful_tokens(stepback_query)
    all_scores = bm25.get_scores(stepback_tokens)
    query_vector = embed_model.encode([stepback_query])
    _, faiss_indices = index.search(np.array(query_vector), 5)
    vector_hits = [chunk_texts[i] for i in faiss_indices[0]]
    bm25_hits = bm25.get_top_n(stepback_tokens, chunk_texts, n=10)
    combined = list(set(bm25_hits + vector_hits))
    reranked = rerank_passages(stepback_query, combined, top_k=5)
    context = "\n".join(reranked)
    num_key_tokens_in_context = sum(token in context.lower() for token in key_tokens)
    if num_key_tokens_in_context < min(2, len(key_tokens)):
        return "Sorry, the provided documents do not contain relevant information for your query."
    prompt = f"""
    You are MedicalBot, a medical assistant. Answer the user's question strictly using ONLY the information in CONTEXT below.
    If the context does NOT contain an answer, respond: "Sorry, the provided documents do not contain relevant information for your query."
    Do NOT use your own knowledge or make up information. Cite medical facts clearly.

    CONTEXT:
    {context}

    QUESTION:
    {query}

    ANSWER:
    """
    return text_gen(prompt, max_new_tokens=350)[0]["generated_text"].strip()
```

## What Makes This System "Advanced"?

This system goes beyond basic RAG by integrating several recent advancements:

- **Hybrid Retrieval:** Combines keyword (BM25) and semantic (FAISS) search for better coverage.
- **State-of-the-art Embeddings:** Uses all-mpnet-base-v2 for effective dense search.
- **Reranking:** Applies the BAAI/bge-reranker for higher-quality retrieval.
- **Step-Back Prompting:** Reformulates complex queries to improve retrieval accuracy.
- **Structured Prompting:** Carefully designed prompts keep the model grounded in the retrieved content and reduce hallucinations.

## Baseline and Evaluation

### Baseline:

Generates an answer using only the text generation model, without retrieval or external context.

### Evaluation:

Compares the baseline and RAG answers based on:

1. **Depth:** Whether the RAG answer is more detailed.
2. **Relevance:** If more query words appear in the answer.
3. **Accuracy:** If the RAG answer avoids contradictions and includes key terms.

The evaluation outputs boolean values showing if the RAG model outperforms the baseline on these aspects.

```
# Baseline
def generate_baseline_answer(query):
    prompt = f"Answer this medical question clearly:\n\n{query}"
    return text_gen(prompt, max_new_tokens=300)[0]["generated_text"].strip()

# Evaluation
def evaluate_quality(baseline, rag, query):
    baseline_len = len(baseline.split())
    rag_len = len(rag.split())
    depth_better = rag_len > baseline_len
    query_words = query.lower().split()
    baseline_match = sum(1 for word in query_words if word in baseline.lower())
    rag_match = sum(1 for word in query_words if word in rag.lower())
    relevance_better = rag_match >= baseline_match
    accuracy_check = ("not" not in rag.lower()) and (rag_match > 0)
    return accuracy_check, relevance_better, depth_better
```

## Overview of Results

The results clearly demonstrate the effectiveness and reliability of the Retrieval-Augmented Generation (RAG) approach in the domain of cardiac anatomy question answering. Compared to a baseline language model, the RAG system especially with step-back prompting provides substantial improvements in accuracy, answer depth, and contextual relevance. Kindly refer **5646778.ipynb** file for complete results.

	Query	Baseline	RAG_Standard	RAG_StepBack	Accuracy	Relevance	Depth
0	What is the anatomical position and orientatio...	The heart is located in the left ventricle of ...	The apex is typically located at the level of ...	The apex is typically located at the level of ...	True	True	True
1	Which chambers form the base of the heart, and...	The base of the heart consists of the left atr...	The base is formed primarily by the atria (the...	The base is formed primarily by the atria (the...	True	True	True
2	What anatomical structures border the heart an...	aorta, mitral valve, aorta lateralis, aorta su...	thoracic duct • Superficially : bifurcation of...	thoracic duct • Superficially : bifurcation of...	True	True	True
3	What is the significance of the interventricul...	The interventricular grooves separate the left...	Sorry, the provided documents do not contain r...	Sorry, the provided documents do not contain r...	False	False	True
4	How big is the brain?	The brain is about the size of a tennis ball.	Sorry, the provided documents do not contain r...	Sorry, the provided documents do not contain r...	False	False	True
5	Which part of the heart forms the apex, and wh...	The apex of the heart forms the apex of the he...	The apex is formed by the left ventricle, and ...	The apex is formed by the left ventricle, and ...	True	True	True

## Key Benefits of the RAG System:

- **Improved Accuracy:** The RAG pipeline answered 88% of the test queries correctly, significantly outperforming the baseline LLM, which frequently produced vague, incomplete, or hallucinated responses.



- **Greater Depth and Detail:** 20 out of 26 responses generated by RAG were rated as more detailed than the baseline, reflecting the model's ability to synthesise and present richer, textbook-like information grounded in the actual dataset.
- **Contextual Relevance:** RAG consistently produced responses that were not only correct but closely matched the user's intent and medical context, as shown by the increase in relevance ratings.
- **Robustness Against Hallucination:** For queries outside the scope of the source documents (e.g., "How big is the brain?"), the system appropriately declined to answer, reducing the risk of spreading misinformation.

### **Benefits of the Developed Model:**

- **Hybrid Retrieval and Reranking:** By combining keyword-based (BM25) and semantic (FAISS) retrieval, then reranking with a cross-encoder model, the system retrieves the most relevant evidence with high precision.
- **Step-Back Prompting:** Reformulating ambiguous or complex queries leads to better retrieval and more reliable answers, particularly for challenging questions.
- **Strict Evidence Filtering:** The model only generates answers when there is clear, direct support in the dataset, further increasing user trust.
- **Generalizability:** The system's design allows for easy adaptation to other domains or updates to the knowledge base, making it a robust solution for specialized question answering tasks.

### **Summary:**

Overall, the RAG based model consistently outperforms a standalone language model across all key metrics. It demonstrates the value of retrieval-augmented architectures for trustworthy, high-quality question answering in specialized fields such as medicine.

### **Testing, Failure Analysis, and Future Improvements**

Rigorous testing was conducted using a diverse set of queries including straightforward, ambiguous, and out-of-scope questions—to ensure the system's robustness across different scenarios.

### **Failure Cases Identified**

- Domain drift: Some queries with generic terms may trigger partially relevant but incomplete answers instead of a clear "not available" response.
- Knowledge base limitations: The system cannot answer questions outside the scope of the provided dataset.
- Automated evaluation limits: Current metrics may miss cases where answers are only partially correct.

## **Future Enhancements**

- Apply stricter relevance checks to further minimise false positives.
- Fine-tune the LLM using in-domain question-answer pairs for greater accuracy and fluency.
- Integrate multi-modal data (such as images or diagrams) to enhance retrieval and understanding.
- Include human expert reviews to complement automated evaluation and further improve reliability.

## **Conclusion**

This RAG system shows clear advantages over standalone language models for medical applications. By combining advanced retrieval, reranking, and prompt engineering, it delivers more accurate and reliable answers. The system effectively handles out-of-domain queries and is well-suited for knowledge-intensive tasks. Further improvements like stricter filtering and multi-modal integration will make RAG even more robust and trustworthy.

## **Reference**

[https://www.physio-pedia.com/Anatomy\\_of\\_the\\_Human\\_Heart](https://www.physio-pedia.com/Anatomy_of_the_Human_Heart)

## **Appendix**

[5646778.ipynb](#)