

## ✓ CSE 256 FA24: NLP UCSD PA3:

### ✓ Retrieval-Augmented Generation (RAG) (40 points)

The goal of this assignment is to gain hands-on experience with aspects of **Retrieval-Augmented Generation (RAG)**, with a focus on retrieval.

You will use **LangChain**, a framework that simplifies integrating external knowledge into generation tasks by:

- Implementing various vector databases for efficient neural retrieval. You will use a vector database for storing our memories.
- Allowing seamless integration of pretrained text encoders, which you will access via HuggingFace models. You will use a text encoder to get text embeddings for storing in the vector database.

#### Data

You will build a retrieval system using the [QMSum Dataset](#), a human-annotated benchmark designed for question answering on long meeting transcripts. The dataset includes over 230 meetings across multiple domains.

**Release Date: November 6, 2024 | Due Date: November 18, 2024**

**IMPORTANT:** After copying this notebook to your Google Drive along with the two data files, paste a link to your copy below. To create a publicly accessible link:

1. Click the *Share* button in the top-right corner.
2. Select "Get shareable link" and copy the link.

**Link: paste your link here:**

<https://drive.google.com/drive/folders/1aKiRfJOakF5uTWbMyczwionPRCzPCyb?usp=sharing>

#### Notes:

Make sure to save the notebook as you go along.

Submission instructions are located at the bottom of the notebook.

```
1 # This mounts your Google Drive to the Colab VM.
2 from google.colab import drive
3 drive.mount('/content/drive')
4
5 # TODO: Enter the foldername in your Drive where you have saved this notebook
6 # e.g. 'CSE156/assignments/PA3/'
7 FOLDERNAME = "Zhecheng Li CSE256 PA3"
8 assert FOLDERNAME is not None, "[!] Enter the foldername."
9
10 # Now that we've mounted your Drive, this ensures that
11 # the Python interpreter of the Colab VM can load
12 # python files from within it.
13 import sys
14 sys.path.append('/content/drive/MyDrive/{}'.format(FOLDERNAME))
```

15

16 # This is later used to use the IMDB reviews

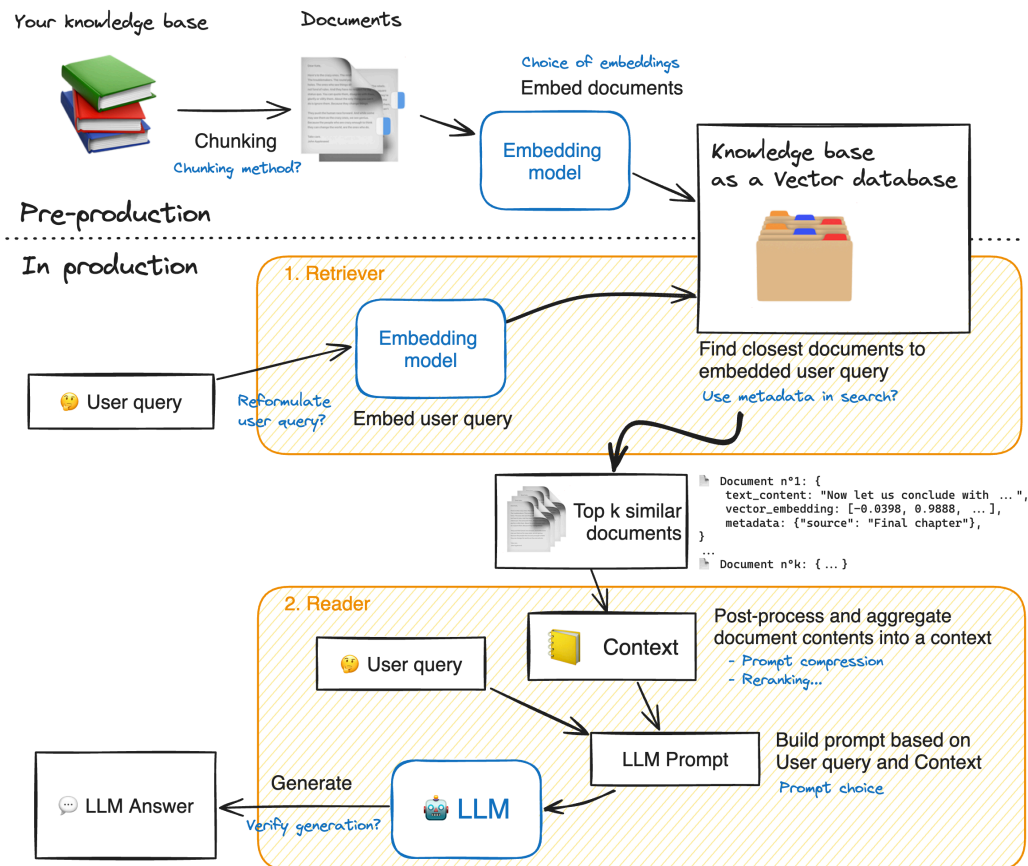
17 %cd /content/drive/MyDrive/\$FOLDERNAME/

Mounted at /content/drive  
/content/drive/MyDrive/Zhecheng Li CSE256 PA3

## RAG Workflow

Retrieval-Augmented Generation (RAG) systems involve several interconnected components. Below is a RAG workflow diagram from Hugging Face. Areas highlighted in blue indicate opportunities for system improvement.

In this assignment, we will focus on the **\*Retriever** so the PA does not cover any processes starting from "2. Reader" and below.



## ✓ First, install the required model dependancies.

```
1 pip install -q torch transformers langchain_chroma bitsandbytes langchain faiss-gpu langchain_huggingface langchain-community sentence-transformers pacmap tqdm matplotlib
```



```

647.5/647.5 kB 17.0 MB/s eta 0:00:00
Preparing metadata (setup.py) ... done
67.3/67.3 kB 5.3 MB/s eta 0:00:00
Installing build dependencies ... done
Getting requirements to build wheel ... done
Preparing metadata (pyproject.toml) ... done
122.4/122.4 MB 7.6 MB/s eta 0:00:00
85.5/85.5 MB 8.9 MB/s eta 0:00:00
2.4/2.4 MB 74.5 MB/s eta 0:00:00
1.0/1.0 MB 54.6 MB/s eta 0:00:00
615.5/615.5 kB 43.0 MB/s eta 0:00:00
2.4/2.4 MB 71.2 MB/s eta 0:00:00
94.7/94.7 kB 9.5 MB/s eta 0:00:00
408.7/408.7 kB 34.0 MB/s eta 0:00:00
3.1/3.1 MB 82.1 MB/s eta 0:00:00
273.8/273.8 kB 21.9 MB/s eta 0:00:00
1.9/1.9 MB 71.5 MB/s eta 0:00:00
49.5/49.5 kB 4.3 MB/s eta 0:00:00
93.2/93.2 kB 9.6 MB/s eta 0:00:00
13.3/13.3 MB 87.0 MB/s eta 0:00:00
55.8/55.8 kB 5.8 MB/s eta 0:00:00
159.2/159.2 kB 15.6 MB/s eta 0:00:00
64.3/64.3 kB 6.5 MB/s eta 0:00:00
118.7/118.7 kB 10.0 MB/s eta 0:00:00
54.4/54.4 kB 4.9 MB/s eta 0:00:00
73.3/73.3 kB 6.1 MB/s eta 0:00:00
63.7/63.7 kB 5.7 MB/s eta 0:00:00
442.1/442.1 kB 35.6 MB/s eta 0:00:00
316.6/316.6 kB 28.8 MB/s eta 0:00:00
3.8/3.8 MB 84.4 MB/s eta 0:00:00
425.7/425.7 kB 29.4 MB/s eta 0:00:00
167.5/167.5 kB 13.5 MB/s eta 0:00:00
46.0/46.0 kB 3.5 MB/s eta 0:00:00
86.8/86.8 kB 7.7 MB/s eta 0:00:00
Building wheel for annoy (setup.py) ... done
Building wheel for pypika (pyproject.toml) ... done

```

ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source of the following dependency conflicts.

google-cloud-datastore 2.19.0 requires protobuf!=3.20.0,!3.20.1,!4.21.0,!4.21.1,!4.21.2,!4.21.3,!4.21.4,!4.21.5,<5.0.0dev,>=3.19.5, but you have protobuf 5.28.3 which is incompatible.

google-cloud-firestore 2.16.1 requires protobuf!=3.20.0,!3.20.1,!4.21.0,!4.21.1,!4.21.2,!4.21.3,!4.21.4,!4.21.5,<5.0.0dev,>=3.19.5, but you have protobuf 5.28.3 which is incompatible.

tensorboard 2.17.0 requires protobuf!=4.24.0,<5.0.0,>=3.19.6, but you have protobuf 5.28.3 which is incompatible.

tensorflow 2.17.0 requires protobuf!=4.21.0,!4.21.1,!4.21.2,!4.21.3,!4.21.4,!4.21.5,<5.0.0dev,>=3.20.3, but you have protobuf 5.28.3 which is incompatible.

tensorflow-metadata 1.16.1 requires protobuf<4.21,>=3.20.3; python\_version < "3.11", but you have protobuf 5.28.3 which is incompatible.

```

1 from tqdm.notebook import tqdm
2 import pandas as pd
3 import os
4 import csv
5 import sys
6 import numpy as np
7 import time
8 import random
9 from typing import Optional, List, Tuple

```

```
10 import matplotlib.pyplot as plt
11 import textwrap
12 import torch
13
14 seed = 42
15 random.seed(seed)
16 np.random.seed(seed)
17 torch.manual_seed(seed)
18 torch.cuda.manual_seed_all(seed)
19 torch.backends.cudnn.deterministic = True
20 torch.backends.cudnn.benchmark = False
21
22 # Disable huggingface tokenizers parallelism
23 os.environ["TOKENIZERS_PARALLELISM"] = "false"
```

## ✓ Load the meetings dataset

```
1 from langchain.docstore.document import Document
2
3 def load_documents(doc_file):
4     """
5     Loads the document contents from the first file.
6
7     :param doc_file: Path to the document file (document ID <TAB> document contents).
8     :return: A dictionary {document_id: document_contents}.
9     """
10     max_size = sys.maxsize
11     csv.field_size_limit(max_size)
12
13
14     documents = {}
15     with open(doc_file, 'r', encoding='utf-8') as f:
16         reader = csv.reader(f, delimiter='\t')
17         for row in reader:
18             if len(row)==0: continue
19             doc_id, content = row
20             documents[doc_id] = content
21     return documents
22
23
24 docs = [] #
25 doc_file = '/content/drive/MyDrive/{}/meetings.tsv'.format(FOLDERNAME)
26 # doc_file = "./meetings.tsv"
27 documents = load_documents(doc_file)
28
29 for doc_id in documents:
30     doc = Document(page_content=documents[doc_id])
31     metadata = {'source': doc_id}
```

```
32 doc.metadata = metadata
33 docs.append(doc)
34
35 print(f"Total meetings (docs): {len(documents)}")
```

↗ Total meetings (docs): 230

## ✓ Retriever - Building the retriever 📁

The **retriever functions like a search engine**: given a user query, it returns relevant documents from the knowledge base.

These documents are then used by the Reader model to generate an answer. In this assignment, however, we are only focusing on the retriever, not the Reader model.

**Our goal**: Given a user question, find the most relevant documents from the knowledge base.

Key parameters:

- `top_k`: The number of documents to retrieve. Increasing `top_k` can improve the chances of retrieving relevant content.
- `chunk_size`: The length of each document. While this can vary, avoid overly long documents, as too many tokens can overwhelm most reader models.

Langchain **offers a huge variety of options for vector databases and allows us to keep document metadata throughout the processing.**

### ✓ 1. Specify an Embedding Model and Visualize Document Lengths

```
1 EMBEDDING_MODEL_NAME = "thenlper/gte-small"
2
3 from sentence_transformers import SentenceTransformer
4
5 print(f"Model's maximum sequence length: {SentenceTransformer(EMBEDDING_MODEL_NAME).max_seq_length}")
6
7 from transformers import AutoTokenizer
8
9 tokenizer = AutoTokenizer.from_pretrained(EMBEDDING_MODEL_NAME)
10 lengths = [len(tokenizer.encode(doc.page_content)) for doc in tqdm(docs)]
11
12 # Plot the distribution of document lengths, counted as the number of tokens
13 fig = pd.Series(lengths).hist()
14 plt.title("Distribution of document lengths in the knowledge base (in count of tokens)")
15 plt.show()
```

```

/usr/local/lib/python3.10/dist-packages/sentence_transformers/cross_encoder/CrossEncoder.py:13: TqdmExperimentalWarning: Using `tqdm.autonotebook.tqdm` in notebook mode. Use `tqdm.tqdm` instead to
from tqdm.autonotebook import tqdm, trange
/usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:89: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as secret in your Google Colab and restart your session.
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
warnings.warn(

```

```

modules.json: 100%                 385/385 [00:00<00:00, 6.56kB/s]

README.md: 100%                   68.1k/68.1k [00:00<00:00, 931kB/s]

sentence_bert_config.json: 100%    57.0/57.0 [00:00<00:00, 1.86kB/s]

config.json: 100%                 583/583 [00:00<00:00, 18.0kB/s]

model.safetensors: 100%           66.7M/66.7M [00:00<00:00, 135MB/s]

tokenizer_config.json: 100%        394/394 [00:00<00:00, 18.6kB/s]

vocab.txt: 100%                   232k/232k [00:00<00:00, 3.27MB/s]

tokenizer.json: 100%              712k/712k [00:00<00:00, 2.74MB/s]

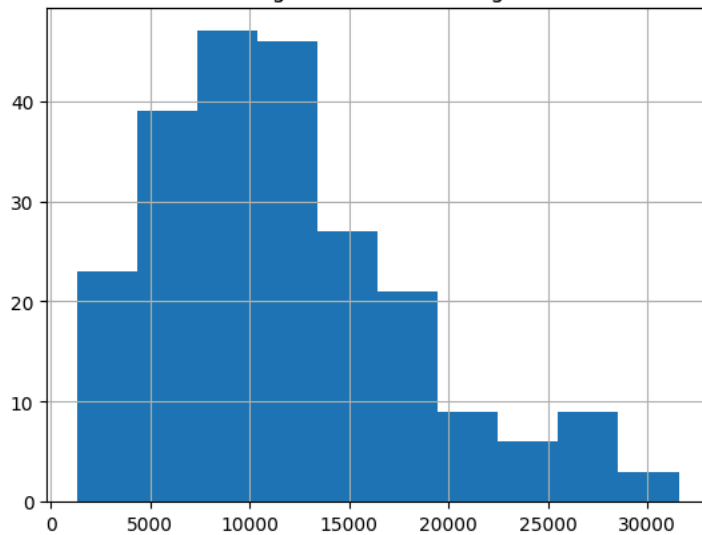
special_tokens_map.json: 100%      125/125 [00:00<00:00, 5.04kB/s]

1_Pooling/config.json: 100%        190/190 [00:00<00:00, 12.9kB/s]

Model's maximum sequence length: 512
100%                             230/230 [00:08<00:00, 30.60it/s]

```

Distribution of document lengths in the knowledge base (in count of tokens)



## 2. Split the Documents into Chunks

The documents (meeting transcripts) are very long—some up to 30,000 tokens! To make retrieval effective, we'll **split each document into smaller, semantically meaningful chunks**. These chunks will serve as the snippets the retriever compares to the query, returning the `top_k` most relevant ones.

**Objective:** Create Semantically Relevant Snippets

Chunks should be long enough to capture complete ideas but not so lengthy that they lose focus.

We will use Langchain's implementation of recursive chunking with `RecursiveCharacterTextSplitter`.

- Parameter `chunk_size` controls the length of individual chunks: this length is counted by default as the number of characters in the chunk.
- Parameter `chunk_overlap` lets adjacent chunks get a bit of overlap on each other. This reduces the probability that an idea could be cut in half by the split between two adjacent chunks.

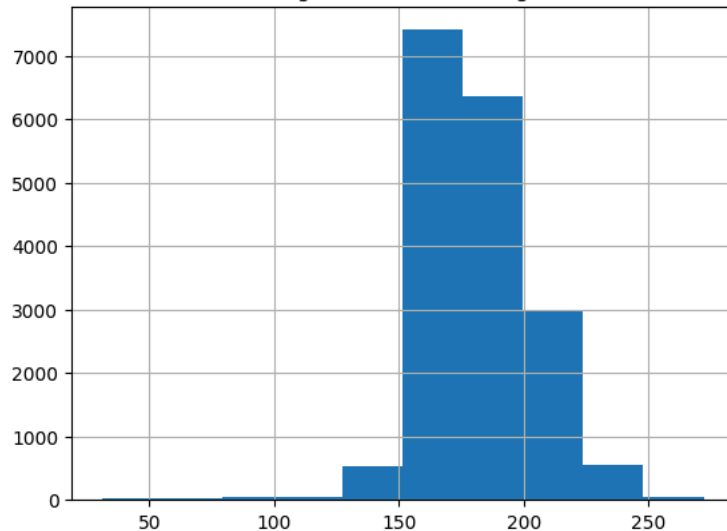
From the produced plot below, you can see that now the chunk length distribution looks better!

```
1 from langchain.text_splitter import RecursiveCharacterTextSplitter
2
3 text_splitter = RecursiveCharacterTextSplitter(chunk_size = 768, chunk_overlap = 128)
4
5 doc_snippets = text_splitter.split_documents(docs)
6 print(f"Total {len(doc_snippets)} snippets to be stored in our vector store.")
7
8 lengths = [len(tokenizer.encode(doc.page_content)) for doc in tqdm(doc_snippets)]
9
10 # Plot the distribution of document snippet lengths, counted as the number of tokens
11 fig = pd.Series(lengths).hist()
12 plt.title("Distribution of document lengths in the knowledge base (in count of tokens)")
13 plt.show()
```

↔ Total 18070 snippets to be stored in our vector store.

100% 18070/18070 [00:12<00:00, 1743.06it/s]

Distribution of document lengths in the knowledge base (in count of tokens)



### 3. Build the Vector Database

To enable retrieval, we need to compute embeddings for all chunks in our knowledge base. These embeddings will then be stored in a vector database.

#### How Retrieval Works

A query is embedded using an embedding model and a similarity search finds the closest matching chunks in the vector database.

The following cell builds the vector database consisting of all chunks in our knowledge base.

```
1 from langchain_huggingface import HuggingFaceEmbeddings
2 from langchain.vectorstores import FAISS
3 from langchain_community.vectorstores.utils import DistanceStrategy
4
5 # Automatically set the device to 'cuda' if available, otherwise use 'cpu'
6 device = "cuda" if torch.cuda.is_available() else "cpu"
7 print(f"Found device: {device}")
8
9 embedding_model = HuggingFaceEmbeddings(
10     model_name=EMBEDDING_MODEL_NAME,
11     multi_process=True,
12     model_kwargs={"device": device},
13     encode_kwargs={"normalize_embeddings": True} # Set `True` for cosine similarity
14 )
15
```



```

16 start_time = time.time()
17
18 KNOWLEDGE_VECTOR_DATABASE = FAISS.from_documents(
19     doc_snippets,
20     embedding_model,
21     distance_strategy=DistanceStrategy.COSINE
22 )
23
24 end_time = time.time()
25
26 elapsed_time = (end_time - start_time) / 60
27 print(f"Time taken: {elapsed_time} minutes")

```

Found device: cuda  
Time taken: 1.4476740638415018 minutes

#### ✓ 4. Querying the Vector Database

Using LangChain's vector database, the function `vector_database.similarity_search(query)` implements a Bi-Encoder (covered in class), independently encoding the query and each document into a single-vector representation, allowing document embeddings to be precomputed.

Let's define the Bi-Encoder ranking function and then use it on a sample query from the QMSum dataset.

```

1 ## The function for ranking documents given a query:
2 def rank_documents_biencoder(user_query, top_k=5):
3     """
4     Function for document ranking based on the query.
5
6     :param query: The query to retrieve documents for.
7     :return: A list of document IDs ranked based on the query (mocked).
8     """
9     retrieved_docs = KNOWLEDGE_VECTOR_DATABASE.similarity_search(query=user_query, k=top_k)
10    ranked_list = []
11    for i, doc in enumerate(retrieved_docs):
12        ranked_list.append(retrieved_docs[i].metadata['source'])
13
14    return ranked_list # ranked document IDs.
15
16
17 user_query = "what did kirsty williams am say about her plan for quality assurance ?"
18 retrieved_docs = rank_documents_biencoder(user_query)
19
20 print("\n=====Top-5 documents=====")
21 print("\n\nRetrieved documents:", retrieved_docs)
22 print("\n=====")

```

=====Top-5 documents=====

Retrieved documents: ['doc\_211', 'doc\_2', 'doc\_43', 'doc\_160', 'doc\_43']

## ▼ 5. TODO: Implementation of ColBERT as a Reranker for a Bi-Encoder (35 points)

The Bi-Encoder's ranking for the sample query is not optimal: the ground truth document is not ranked at position 1, instead the document ID, **doc\_211** is ranked at position 1. To determine the correct document ID for this query, refer to the `questions_answers.tsv` file.

In this task, you will implement the [ColBERT](#) approach by Khattab and Zaharia. We'll use a simplified version of ColBERT, focusing on the following key steps:

1. Retrieve the top (  $K = 15$  ) documents for query (  $q$  ) using the Bi-Encoder.
2. Re-rank these top (  $K = 15$  ) documents using ColBERT's fine-grained interaction scoring. This will involve:
  - Using frozen BERT embeddings from a HuggingFace BERT model (no training is required, thus our version is not expected to work as well as full-fledged ColBERT).
  - Calculating scores based on fine-grained token-level interactions between the query and each document.
3. Implement the method `rank_documents_finegrained_interactions()` to perform this re-ranking.
  - Test your method on the same query as in the cell from #4 above.
  - Print out the entire re-ranked document list of 5 document IDs, as done in #4 above (the code below does it for you)
4. Ensure that your ColBERT implementation ranks the correct document at position 1 for the sample query.

**Note:** Since the same document is divided into multiple chunks that retain the original document ID, you may see the same document ID appear multiple times in your top\_k results. However, each instance refers to a different chunk of the document's content.

```

1 import torch
2 import torch.nn.functional as F
3 from transformers import AutoTokenizer, AutoModel
4 from collections import defaultdict
5
6 # Load tokenizer and model BERT from HuggingFace
7 tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
8 model = AutoModel.from_pretrained("bert-base-uncased")
9
10
11 def rank_documents_finegrained_interactions(user_query, shortlist=15, top_k=5):
12
13     """
14     Rerank the top-K=15 retrieved documents from Bi-encoder using fine-grained token-level interactions
15     and return the top_k=5 most similar documents.
16
17     Args:
18     - user_query (str): The user query string.
19     - shortlist (list): Number of documents in the longer short list
20     - top_k (int): Number of top reranked documents to return.
21
22     Returns:

```

```

23 - ranked_list of document IDs.
24 """
25
26 retrieved_docs = KNOWLEDGE_VECTOR_DATABASE.similarity_search(query=user_query, k=shortlist)
27
28
29 # Tokenize the user query
30 query_inputs = tokenizer(user_query, return_tensors='pt', truncation=True, padding=True)
31
32 # Get query token embeddings from BERT
33 with torch.no_grad():
34     query_embeddings = model(**query_inputs).last_hidden_state # Shape: (1, seq_len_query, hidden_dim)
35
36 ranked_list = []
37
38 ### YOUR CODE HERE
39 doc_scores = []
40
41 for doc in retrieved_docs:
42     doc_inputs = tokenizer(doc.page_content, return_tensors='pt', truncation=True, padding=True)
43
44     with torch.no_grad():
45         doc_embeddings = model(**doc_inputs).last_hidden_state
46
47     query_embeddings_norm = F.normalize(query_embeddings, p=2, dim=2)
48     doc_embeddings_norm = F.normalize(doc_embeddings, p=2, dim=2)
49
50     similarity_matrix = torch.bmm(query_embeddings_norm, doc_embeddings_norm.transpose(1, 2))
51     max_similarities = torch.max(similarity_matrix, dim=2)[0]
52
53     attention_mask = query_inputs['attention_mask']
54     masked_similarities = max_similarities * attention_mask
55     doc_score = masked_similarities.sum() / attention_mask.sum()
56     doc_scores.append((doc.metadata['source'], doc_score.item()))
57
58 grouped_scores = defaultdict(list)
59 for doc_id, score in doc_scores:
60     grouped_scores[doc_id].append(score)
61
62 final_scores = [(doc_id, max(scores)) for doc_id, scores in grouped_scores.items()]
63
64 ranked_list = [doc_id for doc_id, _ in sorted(final_scores, key=lambda x: x[1], reverse=True)[:top_k]]
65
66 return ranked_list # ranked document IDs
67
68
69 user_query = "what did kirsty williams am say about her plan for quality assurance ?"
70 retrieved_docs = rank_documents_finegrained_interactions(user_query)
71
72 print("\n=====Top-5 documents=====")

```

```
73 print("\n\nRetrieved documents:", retrieved_docs)
74 print("=====\n")
```

```

tokenzier_config.json: 100%                               48.0/48.0 [00:00<00:00, 2.52kB/s]

config.json: 100%                                         570/570 [00:00<00:00, 42.2kB/s]

vocab.txt: 100%                                           232k/232k [00:00<00:00, 13.4MB/s]

tokenizer.json: 100%                                     466k/466k [00:00<00:00, 2.38MB/s]
/usr/local/lib/python3.10/dist-packages/transformers/tokenization_utils_base.py:1601: FutureWarning: `clean_up_tokenization_spaces` was not set. It will be set to `True` by default. This behavior
  warnings.warn(
model.safetensors: 100%                                  440M/440M [00:01<00:00, 206MB/s]

=====Top-5 documents=====

Retrieved documents: ['doc_2', 'doc_160', 'doc_43', 'doc_211', 'doc_176']
=====
```

## 6. TODO: ColBERT Max vs. Mean Pooling for Relevance Scoring of Documents: (5 points)

ColBERT uses a form of **max pooling**, where each query term's contribution to the relevance score of a document is determined by its maximum similarity to any document term. One alternative approach is **mean pooling**, where each query term's contribution is calculated as the average similarity across all document terms.

Discuss the merits and potential limitations of using mean pooling versus max pooling in ColBERT. In your answer, consider how each approach might affect retrieval accuracy, sensitivity to specific token matches, interpretability, and anything you deem relevant. You are welcome to use an analysis of ColBERT's performance on the provided sample query in your discussion, but this is not required.

Explain here (<= 5 sentences) :

- Max pooling in ColBERT excels at capturing the best semantic match for each query term, making it more robust to document length variation and better at detecting important single-token matches.
- Mean pooling may be better at capturing overall topical relevance by considering aggregate similarity across all document terms, but can dilute the impact of strong individual token matches.
- Max pooling provides better interpretability by clearly connecting query terms to their best document matches, while mean pooling makes it harder to explain model decisions.
- For general web search and information retrieval, the advantages of max pooling in terms of interpretability, precise matching, and computational efficiency likely make it the more suitable choice for ColBERT's architecture.

## ✓ Submission Instructions

1. Click the Save button at the top of the Jupyter Notebook.

2. Select Runtime -> Run All. This will run all the cells in order, and will take several minutes.
3. Once you've rerun everything, save a PDF version of your notebook. Make sure all your code and answers are displayed in the pdf, it's okay if the provided codes get cut off because lines are not wrapped in code cells).
4. Look at the PDF file and make sure all your code and answers are there, displayed correctly. The PDF is the only thing your graders will see!
5. Submit your PDF on Gradescope.

## ✓ 7. (Optional) Full evaluation pipeline for your own exploration.

For this assignment, we only ask you to explore one sample query. Running on many queries is super slow without the right compute. If you have compute/and/or time to wait, below is a more complete evaluation setup that works with all the queries in QMSum dataset, and reports the `precision@k=5` metric.

**Note:** you need to remove the comment markers from the code below.

ChatGPT

```

1 def load_questions_answers(qa_file):
2     """
3     Loads the questions and corresponding ground truth document IDs.
4
5     :param qa_file: Path to the question-answer file (document ID <TAB> question <TAB> answer).
6     :return: A list of tuples [(document_id, question, answer)].
7     """
8     qa_pairs = []
9     with open(qa_file, 'r', encoding='utf-8') as f:
10         reader = csv.reader(f, delimiter='\t')
11         for row in reader:
12             doc_id, question, answer = row
13             qa_pairs.append((doc_id, question, answer))
14
15     random.shuffle(qa_pairs)
16
17     return qa_pairs
18
19 def precision_at_k(ground_truth, retrieved_docs, k):
20     """
21     Computes Precision at k for a single query.
22
23     :param ground_truth: The name of the ground truth document.
24     :param retrieved_docs: The list of document names returned by the model in ranked order.
25     :param k: The cutoff for computing Precision.
26     :return: Precision at k.
27     """
28     return 1 if ground_truth in retrieved_docs[:k] else 0
29
30 def evaluate(doc_file, qa_file, ranking_fuction=None, k=5):
31     """
32     Evaluate the retrieval system based on the documents and question-answer pairs.
33
34     :param doc_file: Path to the document file

```

```
34 :param doc_file: Path to the document file.
35 :param qa_file: Path to the question-answer file.
36 :param k: The cutoff for Precision@k.
37 """
38 # Load the QA pairs
39 qa_pairs = load_questions_answers(qa_file)
40
41 precision_scores = []
42
43
44 for doc_id, question, _ in qa_pairs:
45     retrieved_docs = ranking_fuction(question)
46     precision_scores.append(precision_at_k(doc_id, retrieved_docs, k))
47
48     avg_precision_at_k = sum(precision_scores) / len(precision_scores)
49
50     if len(precision_scores) % 10 == 0:
51         print(f"After {len(precision_scores)} queries, Precision@{k}: {avg_precision_at_k}")
52
53
54 # Compute average Precision@k
55 avg_precision_at_k = sum(precision_scores) / len(precision_scores)
56
57 print(f"Precision@{k}: {avg_precision_at_k}")
58
59
60 qa_file = 'questions_answers.tsv' # document ID <TAB> question <TAB> answer
61
62 start_time = time.time()
63 evaluate(doc_file, qa_file, rank_documents_biencoder)
64 end_time = time.time()
65 elapsed_time = (end_time - start_time) / 60
66 print(f"Time taken: {elapsed_time} minutes")
```

```
After 10 queries, Precision@5: 0.3
After 20 queries, Precision@5: 0.4
After 30 queries, Precision@5: 0.4
After 40 queries, Precision@5: 0.375
After 50 queries, Precision@5: 0.44
After 60 queries, Precision@5: 0.4333333333333333
After 70 queries, Precision@5: 0.44285714285714284
After 80 queries, Precision@5: 0.4375
After 90 queries, Precision@5: 0.4333333333333333
After 100 queries, Precision@5: 0.45
After 110 queries, Precision@5: 0.42727272727272725
After 120 queries, Precision@5: 0.44166666666666665
After 130 queries, Precision@5: 0.45384615384615384
After 140 queries, Precision@5: 0.4714285714285714
After 150 queries, Precision@5: 0.4733333333333333
After 160 queries, Precision@5: 0.48125
After 170 queries, Precision@5: 0.4764705882352941
After 180 queries, Precision@5: 0.4777777777777778
After 190 queries, Precision@5: 0.46842105263157896
After 200 queries, Precision@5: 0.47
After 210 queries, Precision@5: 0.46190476190476193
After 220 queries, Precision@5: 0.45
After 230 queries, Precision@5: 0.45652173913043476
After 240 queries, Precision@5: 0.4583333333333333
After 250 queries, Precision@5: 0.456
After 260 queries, Precision@5: 0.45384615384615384
After 270 queries, Precision@5: 0.4666666666666667
After 280 queries, Precision@5: 0.4642857142857143
After 290 queries, Precision@5: 0.4586206896551724
After 300 queries, Precision@5: 0.4633333333333333
After 310 queries, Precision@5: 0.45161290322580644
After 320 queries, Precision@5: 0.453125
After 330 queries, Precision@5: 0.4575757575757576
After 340 queries, Precision@5: 0.4588235294117647
After 350 queries, Precision@5: 0.46
After 360 queries, Precision@5: 0.4583333333333333
After 370 queries, Precision@5: 0.4540540540540541
After 380 queries, Precision@5: 0.45789473684210524
After 390 queries, Precision@5: 0.4564102564102564
After 400 queries, Precision@5: 0.4575
After 410 queries, Precision@5: 0.4658536585365854
After 420 queries, Precision@5: 0.46904761904761905
After 430 queries, Precision@5: 0.46744186046511627
After 440 queries, Precision@5: 0.46136363636363636
After 450 queries, Precision@5: 0.4622222222222222
After 460 queries, Precision@5: 0.46956521739130436
After 470 queries, Precision@5: 0.46505744680851064
```