

Below, I'll address your four requests based on the provided code: a line-by-line explanation, an explanation of related concepts, a textual architecture diagram, and points about potential missing elements or improvements.

1. Line-by-Line Explanation of the Code

Setup and Installation

- `!pip install langchain_community`
 - Installs the `langchain_community` package, providing community-driven tools for LangChain.
- `!pip install langchain_openai`
 - Installs `langchain_openai`, enabling OpenAI-specific integrations with LangChain.
- `!pip install faiss-cpu`
 - Installs the CPU version of FAISS (Facebook AI Similarity Search) for efficient vector similarity search.

Document Loading and Splitting

- `from langchain_community.document_loaders import TextLoader`
 - Imports `TextLoader` to load text files as LangChain documents.
- `from langchain_community.vectorstores import FAISS`
 - Imports FAISS for vector storage and retrieval.
- `from langchain_openai import OpenAIEmbeddings`
 - Imports OpenAI's embedding model to convert text into vectors.
- `from langchain_text_splitters import CharacterTextSplitter`
 - Imports a text splitter that divides text based on character count.
- `documents = TextLoader("/content/state_of_the_union.txt").load()`
 - Loads the `state_of_the_union.txt` file into a LangChain document object.
- `text_splitter = CharacterTextSplitter(chunk_size=500, chunk_overlap=100)`
 - Initializes a splitter with 500-character chunks and 100-character overlap for context continuity.
- `texts = text_splitter.split_documents(documents)`
 - Splits the loaded document into smaller chunks stored in `texts`.

- texts
 - Displays the list of chunks (though this line doesn't print unless explicitly called).

API Key Configuration

- from google.colab import userdata
 - Imports Colab's utility to access user-stored secrets.
- OPENAI_API_KEY = userdata.get('OPENAI_API_KEY')
 - Retrieves the OpenAI API key from Colab's user data.
- import os
 - Imports the os module for environment variable management.
- os.environ["OPENAI_API_KEY"] = OPENAI_API_KEY
 - Sets the API key as an environment variable for LangChain to use.

Initial Retrieval Setup

- retriever = FAISS.from_documents(texts, OpenAIEmbeddings()).as_retriever()
 - Creates a FAISS vector store from the text chunks using OpenAI embeddings, then converts it to a retriever (default k=4 documents returned).
- docs = retriever.invoke("What did the president say about Ketanji Brown Jackson")
 - Retrieves documents relevant to the query about Ketanji Brown Jackson.

Helper Function

- def pretty_print_docs(docs):
 - Defines a function to format and print documents.
- print(f"\n{'-' * 100}\n".join([f"Document {i+1}:\n\n" + d.page_content for i, d in enumerate(docs)]))
 - Joins document contents with 100-dash separators, printing each with an index (no metadata here, unlike the previous code).

Basic Retrieval Examples

- pretty_print_docs(docs)
 - Prints the retrieved documents for the Ketanji Brown Jackson query.

- docs2 = retriever.invoke("What were the top three priorities outlined in the most recent State of the Union address?")
 - Retrieves documents for a query about top priorities.
- pretty_print_docs(docs2)
 - Prints those documents.
- docs3 = retriever.invoke("How did the President propose to tackle the issue of climate change?")
 - Retrieves documents for a climate change query.
- pretty_print_docs(docs3)
 - Prints those documents.

Retrieval-Augmented QA

- from langchain_openai import OpenAI
 - Imports the OpenAI LLM class (non-chat version).
- llm = OpenAI(temperature=0)
 - Initializes the OpenAI LLM with zero temperature for deterministic output.
- from langchain.chains import RetrievalQA
 - Imports RetrievalQA for question answering over retrieved documents.
- chain = RetrievalQA.from_chain_type(llm=llm, retriever=retriever)
 - Creates a QA chain with the LLM and basic FAISS retriever.
- query = "What were the top three priorities outlined in the most recent State of the Union address?"
 - Defines a query for the QA chain.
- chain.invoke(query)
 - Runs the QA chain, returning a dictionary with the answer in result.
- print(chain.invoke(query)['result'])
 - Prints just the answer text.

Contextual Compression with LLMChainExtractor

- from langchain.retrievers import ContextualCompressionRetriever
 - Imports a retriever that compresses results for relevance.
- from langchain.retrievers.document_compressors import LLMChainExtractor
 - Imports a compressor that extracts relevant parts of documents using an LLM.
- compressor = LLMChainExtractor.from_llm(llm)

- Initializes the extractor with the OpenAI LLM.
- `compression_retriever = ContextualCompressionRetriever(base_compressor=compressor, base_retriever=retriever)`
 - Combines the FAISS retriever with the LLM extractor.
- `compressed_docs = compression_retriever.invoke("What did the president say about Ketanji Jackson Brown")`
 - Retrieves and compresses documents for a slightly different Ketanji query.
- `compressed_docs`
 - Displays the compressed documents (not printed unless formatted).
- `compressed_docs = compression_retriever.invoke("What were the top three priorities outlined in the most recent State of the Union address?")`
 - Compresses documents for the priorities query.
- `pretty_print_docs(compressed_docs)`
 - Prints the compressed documents.
- `compressed_docs2 = compression_retriever.invoke("How did the President propose to tackle the issue of climate change?")`
 - Compresses documents for the climate change query.
- `pretty_print_docs(compressed_docs2)`
 - Prints those compressed documents.

Contextual Compression with LLMChainFilter

- `from langchain.retrievers.document_compressors import LLMChainFilter`
 - Imports a compressor that filters out irrelevant documents entirely.
- `filter = LLMChainFilter.from_llm(llm)`
 - Initializes the filter with the OpenAI LLM.
- `compression_retriever2 = ContextualCompressionRetriever(base_compressor=filter, base_retriever=retriever)`
 - Combines the FAISS retriever with the filter.
- `compressed_docs3 = compression_retriever2.invoke("What were the top three priorities outlined in the most recent State of the Union address?")`
 - Retrieves and filters documents for the priorities query.
- `pretty_print_docs(compressed_docs3)`
 - Prints the filtered documents.

Compression Metrics

- `original_contexts_len = len("\n\n".join([d.page_content for i, d in enumerate(docs2)]))`
 - Calculates the total character length of the original retrieved documents for the priorities query.
- `compressed_contexts_len = len("\n\n".join([d.page_content for i, d in enumerate(compressed_docs)]))`
 - Calculates the length of the compressed documents (from LLMChainExtractor).
- `print("Original context length:", original_contexts_len)`
 - Prints the original length.
- `print("Compressed context length:", compressed_contexts_len)`
 - Prints the compressed length.
- `print("Compressed Ratio:", f"{original_contexts_len/(compressed_contexts_len + 1e-5):.2f}x")`
 - Prints the compression ratio, adding a small epsilon (1e-5) to avoid division by zero.

Contextual Compression with EmbeddingsFilter

- `from langchain.retrievers.document_compressors import EmbeddingsFilter`
 - Imports a compressor that filters documents based on embedding similarity.
- `embeddings = OpenAIEmbeddings()`
 - Initializes OpenAI embeddings (default text-embedding-ada-002).
- `embeddings_filter = EmbeddingsFilter(embeddings=embeddings)`
 - Creates a filter using the embeddings (default similarity threshold).
- `compression_retriever3 = ContextualCompressionRetriever(base_compressor=embeddings_filter, base_retriever=retriever)`
 - Combines the FAISS retriever with the embeddings filter.
- `compressed_docs4 = compression_retriever3.invoke("What were the top three priorities outlined in the most recent State of the Union address?")`
 - Retrieves and filters documents using embeddings.
- `pretty_print_docs(compressed_docs4)`
 - Prints the filtered documents.
- `print("Original context length:", original_contexts_len)`

- Reprints the original length.
- `compressed_contexts_len = len("\n\n".join([d.page_content for i, d in enumerate(compressed_docs)]))`
 - Recalculates the compressed length (still using `compressed_docs` from `LLMChainExtractor`).
- `print("Compressed context length:", compressed_contexts_len)`
 - Prints the compressed length.
- `print("Compressed Ratio:", f"{original_contexts_len/(compressed_contexts_len + 1e-5):.2f}x")`
 - Prints the compression ratio (note: this uses the wrong `compressed_docs`; should use `compressed_docs4`).

Pipeline Compression

- `from langchain.retrievers.document_compressors import DocumentCompressorPipeline`
 - Imports a pipeline for chaining multiple compressors.
- `from langchain_community.document_transformers import EmbeddingsRedundantFilter`
 - Imports a filter to remove redundant documents based on embeddings.
- `splitter = CharacterTextSplitter(chunk_size=300, chunk_overlap=0, separator=". ")`
 - Initializes a splitter with 300-character chunks, no overlap, splitting on periods.
- `redundant_filter = EmbeddingsRedundantFilter(embeddings=embeddings)`
 - Creates a filter to remove redundant documents using embeddings.
- `relevant_filter = EmbeddingsFilter(embeddings=embeddings, similarity_threshold=0.76)`
 - Creates a filter to keep only documents with similarity above 0.76.
- `pipeline_compressor = DocumentCompressorPipeline(transformers=[splitter, redundant_filter, relevant_filter])`
 - Chains the splitter, redundant filter, and relevance filter into a pipeline.
- `compression_retriever = ContextualCompressionRetriever(base_compressor=pipeline_compressor, base_retriever=retriever)`
 - Combines the FAISS retriever with the pipeline.
- `compressed_docs = compression_retriever.invoke("What were the top three priorities outlined in the most recent State of the Union address?")`

- Retrieves and compresses documents using the pipeline.
- `pretty_print_docs(compressed_docs)`
 - Prints the pipeline-compressed documents.

Final QA with ChatOpenAI

- `from langchain_openai import ChatOpenAI`
 - Imports the chat version of OpenAI's LLM.
- `llm = ChatOpenAI(temperature=0)`
 - Initializes the chat LLM with zero temperature.
- `chain = RetrievalQA.from_chain_type(llm=llm, retriever=compression_retriever)`
 - Creates a QA chain with the chat LLM and pipeline retriever.
- `query = "What were the top three priorities outlined in the most recent State of the Union address?"`
 - Defines the query again.
- `chain.invoke(query)`
 - Runs the QA chain.
- `print(chain.invoke(query)['result'])`
 - Prints the generated answer.

Hardcoded Output

- The top three priorities outlined in the most recent State of the Union address were: ...
 - A hardcoded example output (not executed from the code but provided for illustration).

2. Explanation of Concepts Related to the Code

- 1. Document Loading and Splitting:**
 - a. Loading text files (`TextLoader`) and splitting them into chunks (`CharacterTextSplitter`) prepares large documents for retrieval by fitting them within model context limits.
- 2. Embeddings:**
 - a. OpenAI embeddings (`OpenAIEmbeddings`) convert text into dense vectors capturing semantic meaning, enabling similarity-based retrieval.
- 3. Vector Search (FAISS):**

- a. FAISS stores document embeddings and retrieves the top k most similar to a query embedding, forming the basis of the initial retrieval.
- 4. **Retrieval:**
 - a. The retriever fetches relevant documents based on vector similarity, a foundational step in retrieval-augmented generation (RAG).
- 5. **Retrieval-Augmented Generation (RAG):**
 - a. RetrievalQA combines retrieval (FAISS) with generation (OpenAI LLM) to answer queries using external context, reducing hallucination.
- 6. **Contextual Compression:**
 - a. Enhances retrieval by refining initial results:
 - i. **LLMChainExtractor:** Extracts relevant snippets from documents using an LLM.
 - ii. **LLMChainFilter:** Discards irrelevant documents entirely.
 - iii. **EmbeddingsFilter:** Filters based on embedding similarity to the query.
 - iv. **DocumentCompressorPipeline:** Chains multiple compressors (splitting, redundancy removal, relevance filtering).
- 7. **LLM (Language Model):**
 - a. OpenAI's OpenAI (non-chat) and ChatOpenAI (chat) models generate answers or assist in compression. temperature=0 ensures deterministic output.
- 8. **Compression Metrics:**
 - a. Measures like context length and compression ratio quantify how much irrelevant content is removed, improving efficiency and focus.
- 9. **Pipeline Compression:**
 - a. Combines splitting, redundancy removal, and relevance filtering to iteratively refine documents, offering a modular approach to compression.

3. Architecture Diagram

Here's a textual representation of the architecture:

text

CollapseWrapCopy

```
[Query] --> [FAISS Retriever] | v [Initial Documents] --> [Contextual
Compression] | | | v | [Compressor Options] | - LLMChainExtractor | -
LLMChainFilter | - EmbeddingsFilter | - Pipeline (Splitter ->
```



```
RedundantFilter -> RelevantFilter) v | [Original Docs] [Compressed Docs] | | v v [RetrievalQA Chain] <--- [LLM (OpenAI/ChatOpenAI)] | v [Generated Answer]
```

- **Flow:**

- Query enters the FAISS retriever.
- Initial documents are retrieved.
- Compression refines them using one or more methods.
- The QA chain processes the compressed documents with an LLM to generate an answer.

4. Missing Points or Suggestions

1. Error Handling:

- a. No try-except blocks for file loading (state_of_the_union.txt) or API key retrieval, risking runtime errors.

2. Metadata Usage:

- a. The pretty_print_docs function omits metadata (unlike the previous code), which could be useful for debugging or context.

3. Compression Metrics Bug:

- a. The second compression ratio calculation uses compressed_docs (from LLMChainExtractor) instead of compressed_docs4 (from EmbeddingsFilter), leading to incorrect reporting.

4. Default Retriever Parameters:

- a. The FAISS retriever uses a default k=4, which isn't explicitly set or justified. Adjusting k could affect results.

5. Query Consistency:

- a. Queries vary slightly (e.g., "Ketanji Brown Jackson" vs. "Ketanji Jackson Brown"), which might affect consistency.

6. Redundant Calculations:

- a. compressed_contexts_len is recalculated multiple times; it could be stored or computed once per compressor.

7. Pipeline Threshold Tuning:

- a. The similarity_threshold=0.76 in the pipeline is arbitrary; tuning or justifying it could improve results.

8. Performance Evaluation:

- a. No evaluation of answer quality (e.g., accuracy, relevance) beyond compression ratios.

9. File Dependency:

- a. Assumes `state_of_the_union.txt` exists without a fallback or documentation.

10. Hardcoded Output:

- a. The final output is hardcoded rather than dynamically generated from the code, suggesting a disconnect between execution and presentation.

Below, I'll address your requests: explaining the concepts related to the provided code, providing detailed notes on how Contextual Compression Retrieval works in generative AI, and outlining the theoretical basics of its operation.

1. Concepts Related to the Provided Code

The code implements a retrieval-augmented generation (RAG) pipeline with a focus on enhancing retrieval through contextual compression. Here are the key concepts involved:

- **Document Loading and Splitting:**
 - Loading a text file (`TextLoader`) and splitting it into manageable chunks (`CharacterTextSplitter`) ensures large documents can be processed for retrieval and fit within model context limits.
- **Embeddings:**
 - OpenAI embeddings (`OpenAIEmbeddings`) convert text into dense vectors that encode semantic meaning, enabling similarity-based retrieval.
- **Vector Search (FAISS):**
 - FAISS (Facebook AI Similarity Search) stores document embeddings and retrieves the top k documents most similar to a query embedding, providing the initial retrieval step.
- **Retrieval:**

- The process of fetching potentially relevant documents from a corpus based on a query, serving as the foundation for RAG.
- **Retrieval-Augmented Generation (RAG):**
 - Combines retrieval (via FAISS) with generation (via OpenAI LLMs like OpenAI or ChatOpenAI) to answer queries using external context, improving accuracy over pure generation.
- **Contextual Compression:**
 - Refines retrieved documents to focus on the most relevant content, reducing noise and improving efficiency. Methods include:
 - **LLMChainExtractor:** Extracts key snippets using an LLM.
 - **LLMChainFilter:** Filters out irrelevant documents using an LLM.
 - **EmbeddingsFilter:** Filters based on embedding similarity.
 - **DocumentCompressorPipeline:** Chains splitting, redundancy removal, and relevance filtering.
- **Large Language Models (LLMs):**
 - OpenAI's OpenAI (non-chat) and ChatOpenAI (chat) models are used for generation and compression tasks, with temperature=0 ensuring deterministic output.
- **Compression Metrics:**
 - Metrics like context length and compression ratio quantify the reduction in irrelevant content, enhancing focus and reducing computational load.

These concepts work together to retrieve documents, refine them via contextual compression, and generate answers, optimizing for relevance and efficiency in generative AI.

2. Detailed Notes on How Contextual Compression Retrieval Works in Generative AI

Contextual Compression Retrieval enhances traditional retrieval by refining the initial set of retrieved documents to provide more relevant and concise context to generative models. Here's a detailed explanation of its role and operation in generative AI:

Overview of Contextual Compression Retrieval

- **Purpose:** To improve the quality of retrieved documents by filtering or extracting the most relevant content, ensuring that generative models (e.g., LLMs) receive high-signal context for accurate and concise responses.
- **Core Idea:** Instead of passing all retrieved documents directly to the LLM, contextual compression reranks, filters, or extracts content to reduce noise and focus on what matters most to the query.
- **Implementation in the Code:** The `ContextualCompressionRetriever` wraps a base retriever (FAISS) with various compressors (`LLMChainExtractor`, `LLMChainFilter`, `EmbeddingsFilter`, `DocumentCompressorPipeline`) to refine results.

How It Works in the Code

1. Base Retrieval:

- a. **FAISS Retriever:** Retrieves an initial set of documents (e.g., top 4 by default) based on embedding similarity to the query.
- b. Example: For "What did the president say about Ketanji Brown Jackson," FAISS returns docs with chunks from `state_of_the_union.txt`.

2. Compression Step:

- a. **ContextualCompressionRetriever:** Takes the base retriever's output and applies a compressor:

i. **LLMChainExtractor:**

1. Uses the LLM to extract relevant snippets from each document.
2. Example: From a 500-character chunk, it might extract "The President praised Ketanji Brown Jackson for her judicial record."
3. Output: `compressed_docs` with reduced but highly relevant content.

ii. **LLMChainFilter:**

1. Uses the LLM to classify each document as relevant or not, discarding irrelevant ones.
2. Example: Keeps only chunks explicitly mentioning "Ketanji Brown Jackson."
3. Output: `compressed_docs3` with fewer documents.

iii. **EmbeddingsFilter:**

1. Recomputes embedding similarity between the query and documents, filtering out those below a threshold.
2. Example: Keeps documents with cosine similarity > default threshold (e.g., 0.8).
3. Output: compressed_docs4 with a refined subset.

iv. **DocumentCompressorPipeline:**

1. Chains multiple transformers:
 - a. **Splitter:** Breaks documents into smaller units (e.g., sentences).
 - b. **RedundantFilter:** Removes near-duplicates based on embedding similarity.
 - c. **RelevantFilter:** Keeps units above a similarity threshold (e.g., 0.76).
2. Example: Splits a chunk into sentences, removes redundant ones, and keeps only those relevant to "top three priorities."
3. Output: compressed_docs with concise, unique, relevant content.

3. **Integration with Generative AI:**

- a. **RetrievalQA Chain:** Feeds the compressed documents to the LLM (OpenAI or ChatOpenAI) to generate an answer.
- b. Example: For "What were the top three priorities...", the pipeline-compressed documents are passed to ChatOpenAI, yielding a structured response like:

text

CollapseWrapCopy

1. Beating the opioid epidemic... 2. Strengthening infrastructure...
3. Promoting domestic production...

c. **Process:**

- i. Base retriever fetches raw context.
- ii. Compressor refines it.
- iii. LLM generates a response grounded in the refined context.

Role in Generative AI

- **Context Quality:** Ensures the LLM receives only the most relevant information, reducing the risk of hallucination or irrelevant tangents.
- **Efficiency:** Shrinks the context size (e.g., from 2000 to 500 characters), lowering token usage and speeding up generation.

- **Flexibility:** Offers multiple compression strategies to balance precision (LLM-based) and speed (embedding-based), adaptable to different use cases.

Benefits

- **Improved Relevance:** By reranking or extracting content, it prioritizes what's most pertinent to the query.
- **Reduced Noise:** Filters out filler or off-topic content common in initial retrievals.
- **Scalability:** Works with large corpora by keeping only what's needed, making RAG practical for real-world applications.

3. Theoretical Basics of Contextual Compression Retrieval Working

Let's explore the theoretical underpinnings of Contextual Compression Retrieval, focusing on how it refines retrieval for generative AI.

Theoretical Foundation

- **Base Retrieval:**
 - Relies on a vector space model where documents and queries are embedded into a high-dimensional space (e.g., 1536 dimensions with OpenAIEmbeddings).
 - Similarity metric (e.g., cosine similarity) ranks documents:

$$S_i = \cos(\vec{Q}, \vec{D}_i) = \frac{\vec{Q} \cdot \vec{D}_i}{\|\vec{Q}\| \|\vec{D}_i\|} = \frac{\|\vec{Q}\| \|\vec{D}_i\| \cos(\angle(\vec{Q}, \vec{D}_i))}{\|\vec{Q}\| \|\vec{D}_i\|} = \cos(\angle(\vec{Q}, \vec{D}_i))$$
 - Output: Top k documents $D = \{D_1, D_2, \dots, D_k\}$, ranked by S_i .
- **Compression Goal:** Refine D into a smaller, more relevant set $D' = \{D'_1, D'_2, \dots, D'_m\}$ where $m \leq k$, using semantic or similarity-based methods.

Mechanisms by Compressor Type

1. LLMChainExtractor:

- Input:** Query Q and document D_i .
- Process:**

- i. LLM (e.g., OpenAI) is prompted: "Extract the part of $D_i D_i$ most relevant to Q ."
 - ii. Internally:
 1. Tokenizes Q and $D_i D_i$.
 2. Uses attention mechanisms to weigh token interactions (e.g., "Ketanji" in Q with "Jackson" in $D_i D_i$).
 3. Outputs a substring $D'_i D'_i$ via generation or span selection.
 - iii. Example: $D_i = D_i = D_i =$ "The President discussed many topics, including praising Ketanji Brown Jackson" $\rightarrow D'_i = D'_i = D'_i =$ "The President praised Ketanji Brown Jackson."
 - c. **Theory:** Leverages the LLM's contextual understanding, trained on vast data to identify relevance.
 - d. **Output:** $D' D'$ with reduced length but preserved meaning.
2. **LLMChainFilter:**
- a. **Input:** Q and $D_i D_i$.
 - b. **Process:**
 - i. LLM is prompted: "Is $D_i D_i$ relevant to Q ?" (yes/no).
 - ii. Internally:
 1. Encodes Q and $D_i D_i$ into a joint representation.
 2. Classifies relevance using a softmax over a binary output (relevant/irrelevant).
 - iii. Example: $D_i = D_i = D_i =$ "The weather was sunny" \rightarrow discarded; $D_i = D_i = D_i =$ "Ketanji Brown Jackson was nominated" \rightarrow kept.
 - c. **Theory:** Acts as a binary classifier, using the LLM's discriminative power.
 - d. **Output:** $D' \subseteq D$ | $subseq D' \subseteq D$, subset of relevant documents.
3. **EmbeddingsFilter:**
- a. **Input:** Q and $\{D_1, D_2, \dots, D_k\}$ | $\{vec\{D_1\}, vec\{D_2\}, \dots, vec\{D_k\}\}$ | $\{D_1, D_2, \dots, D_k\}$.
 - b. **Process:**
 - i. Recompute similarity $S_i = \cos(Q, D_i)$ | $S_i = \cos(vec\{Q\}, vec\{D_i\})$ | $S_i = \cos(Q, D_i)$.
 - ii. Filter: Keep D_i if $S_i > \theta$ | $S_i > \theta$ (threshold, e.g., 0.8).
 - iii. Example: If $S_1 = 0.85$, $S_2 = 0.6$, and $\theta = 0.7$, keep D_1 , discard D_2 .
 - c. **Theory:** Refines ranking by enforcing a stricter similarity cutoff, relying on pre-trained embedding semantics.
 - d. **Output:** $D' \subseteq D$ | $subseq D' \subseteq D$, subset with high similarity.

4. DocumentCompressorPipeline:

- a. **Input:** Q and $D = \{D_1, \dots, D_k\}$
- b. **Process:**
 - i. **Splitter:** Breaks D_i into smaller units: $D_i \rightarrow \{D_{i1}, D_{i2}, \dots, D_{im}\}$
 1. Example: Splits on "." to get sentences.
 - ii. **RedundantFilter:** Removes near-duplicates:
 1. Compute $\cos(\vec{Q}, \vec{D_{ij}})$ and $\cos(\vec{D_{ij}}, \vec{D_{ik}})$; if > 0.95 , discard one.
 - iii. **RelevantFilter:** Keeps units with $\cos(\vec{Q}, \vec{D_{ij}}) > \theta$ (e.g., 0.76).
 - iv. Example: From "The President spoke. Ketanji was praised. The weather was nice," keeps "Ketanji was praised."
- c. **Theory:** Multi-stage refinement using syntactic splitting and semantic filtering.
- d. **Output:** D' , a concise, unique, relevant set.

Theoretical Workflow

- **Step 1: Base Retrieval:**
 - $D = \text{FAISS}(Q, \text{corpus}, k)$
- **Step 2: Compression:**
 - $D' = \text{Compressor}(Q, D)$
 - Compressor applies scoring/extraction/filtering, effectively reranking D into D' .
- **Step 3: Generation:**
 - $\text{Answer} = \text{LLM}(Q, D')$

Complexity

- **LLM-Based:** $O(k \cdot T_{LLM})$, where T_{LLM} is inference time per document.
- **Embeddings-Based:** $O(k \cdot d)$, where d is embedding dimension.
- **Pipeline:** Sum of each stage's complexity, typically dominated by LLM or embedding computation.

Why It Works in Generative AI

- **Precision:** Filters out noise, ensuring the LLM focuses on high-relevance content.
- **Efficiency:** Reduces context size, lowering token costs and speeding up generation.
- **Adaptability:** Combines LLM's semantic power with embeddings' speed, fitting various generative tasks.

Summary

- **Concepts:** The code integrates retrieval, compression, and generation for RAG, with contextual compression as a key enhancer.
- **Contextual Compression Retrieval:** Refines retrieved documents using LLM or embedding-based methods, improving context quality for generative AI.
- **Theory:** Operates by reranking or filtering via semantic extraction (LLM), binary classification (LLM), similarity thresholds (embeddings), or multi-stage pipelines.