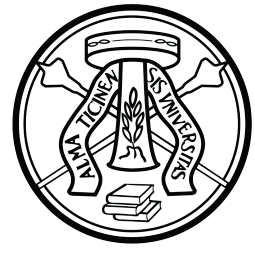




Università degli
Studi di Milano



Università degli Studi di
Milano - Bicocca



Università di
Pavia



Master's Degree Programme
Artificial Intelligence for Science and Technology

A Comprehensive Study on Retrieval Augmented Generation Methods for More Robust LLMs

Supervisor:
Dr. Napoletano Paolo

Candidate name:
Andrea Palmieri
Registration number:
921785

Academic Year 2024/2025

Here you can write whatever you want and here too.

Contents

Contents	v
1 Introduction	1
2 Fundamentals of Retrieval-Augmented Generation	3
2.1 How it Works	3
2.2 Main Challenges	4
2.2.1 Ensuring Content Relevance	4
2.2.2 Optimizing Retrieval vs. Generation Trade-offs	5
2.2.3 Handling Noisy or Conflicting Information	5
2.2.4 Seamless Integration and Synthesis	5
2.3 Vector Databases and Similarity Search	5
2.3.1 Measuring Similarity	6
2.3.2 The Advantage of Normalized Vectors	6
2.3.3 Indexing for Efficient Search	6
2.4 Mitigating Hallucinations	7
3 Retrieval and Optimization Methods	9
3.1 Chunking Techniques	9
3.1.1 Naive vs. Semantic Chunking	9
3.2 Embedding Models	10
3.2.1 Contextual Embeddings	10
3.2.2 Fine-tuning Embedding Models	10
3.3 Post-Retrieval Reranking and Filtering	11
3.3.1 BM25 and TF-IDF for Reranking	11
3.3.2 Hybrid Systems: Combining Similarity with BM25/TF-IDF	11
3.3.3 Cross-Encoder Rerankers	11
3.4 Dynamic Similarity Thresholding	11

3.5	Late Chunking with Contextual Chunk Embeddings	12
3.5.1	Limitations of Traditional Chunking	12
3.5.2	The Late Chunking Process	12
4	Prompt Engineering for RAG	17
4.1	Concepts of Prompt Engineering	17
4.2	Prompt Styles	18
4.3	Prompt Evaluation	18
4.3.1	Golden Datasets	19
4.3.2	LLM-as-a-Judge	19
4.4	Evaluation Metrics	19
5	Comparative Analysis of Different LLMs	21
5.1	The Role of the Generator LLM	21
5.2	Performance Evaluation based on Model and Context Window Size .	22
5.2.1	Comparing LLM Families	22
5.2.2	The Impact of Context Window Size	22
5.3	Trade-offs in LLM Selection	23
6	Identifying Relevant Chunks for Responses	25
6.1	The Importance of Chunk-Response Alignment	25
6.2	Methods for Analyzing Chunk-Response Alignment	26
6.2.1	Prompt-Based Attribution	26
6.2.2	Post-Hoc Similarity Analysis	26
6.2.3	Analyzing Attention Mechanisms	27
6.2.4	Building a Knowledge Graph	27
6.3	Evaluation Metrics for Retrieval	27
7	Experiments and Results	29
7.1	General Experimental Setup	29
7.1.1	Dataset(s)	29
7.1.2	Baseline Configuration	29
7.1.3	Core Evaluation Metrics	30
7.2	Experiment 1: Impact of Chunking Techniques	30
7.2.1	Methods Compared	30
7.2.2	Results and Discussion	31
7.3	Experiment 2: Evaluation of Different Embedding Models for Long- Context Retrieval	31

7.3.1	Models and Methodology	31
7.3.2	Results and Discussion	32
7.4	Experiment 3: Reranking Strategies	33
7.4.1	Reranking Methods Compared	33
7.4.2	Results and Discussion	33
7.5	Experiment 4: Generative Model and Prompt Evaluation	33
7.5.1	Motivation and Baseline	34
7.5.2	Evaluation Methodology	34
7.5.3	Results and Discussion	35
7.5.4	Comparative Analysis Against GPT-4o	38
7.6	Overall Performance Summary and Analysis	39
8	Conclusions and Future Work	41
8.1	Future Work	42

Chapter 1

Introduction

The advent of Large Language Models (LLMs) has marked a pivotal moment in the field of artificial intelligence, demonstrating unprecedented capabilities in understanding and generating human-like text. These models, trained on vast and diverse datasets, have become the backbone of numerous applications, from sophisticated chatbots to content creation tools. However, their power is not without limitations. Two significant challenges inherent to LLMs are *knowledge cutoff* and *hallucination*. Knowledge cutoff refers to the fact that an LLM’s knowledge is static and frozen at the point its training data was collected, rendering it unable to provide information on events or developments that occurred post-training. Hallucination is the tendency of these models to generate plausible but factually incorrect or nonsensical information, a byproduct of their probabilistic nature and their training objective to predict the next word rather than to state facts.

This thesis investigates Retrieval-Augmented Generation (RAG), a paradigm designed to directly address these shortcomings [1]. RAG enhances the capabilities of LLMs by grounding their responses in external, up-to-date knowledge sources. Instead of relying solely on its internal, parametric knowledge, a RAG system first retrieves relevant information from a specified corpus—such as a collection of scientific papers, a corporate knowledge base, or the entire web—and then uses this information to inform its generation process. This approach aims to mitigate hallucinations by providing factual context and overcomes the knowledge cutoff issue by allowing access to real-time information. The core principle is to combine the generative power of LLMs with the factual accuracy of information retrieval systems, leading to more robust, trustworthy, and contextually appropriate responses [2].

The RAG paradigm has evolved significantly since its inception. As categorized by Gao et al. (2024) [2], the landscape can be understood through three main stages:

- **Naive RAG:** The initial and most straightforward implementation, involving a simple retrieve-then-read process.
- **Advanced RAG:** Focuses on optimizing the retrieval stage through techniques like pre-retrieval processing and post-retrieval reranking.
- **Modular RAG:** A more flexible and powerful approach that treats RAG as a modular framework, allowing for greater adaptability and integration of different components and patterns.

This thesis will explore these concepts, starting with the foundational Naive RAG architecture and progressing to more advanced techniques. We will conduct a deep dive into the critical elements of the RAG pipeline, including:

- **Chunking Strategies:** The methods for segmenting large documents into smaller, digestible pieces for the retriever.
- **Embedding Models:** The models used to convert text into numerical representations for semantic search.
- **Retrieval and Reranking:** The techniques for identifying the most relevant information and refining the selection, including hybrid search methods and dynamic thresholding.
- **Prompt Engineering:** The art of crafting effective prompts to guide the LLM in synthesizing the retrieved context.
- **LLM Selection:** The impact of choosing different generator models on the final output.

Through a series of structured experiments, this thesis will systematically evaluate how each of these components influences the overall system performance, with the ultimate goal of providing a clear framework for building and optimizing robust RAG systems.

Chapter 2

Fundamentals of Retrieval-Augmented Generation

Retrieval-Augmented Generation (RAG) is an architectural pattern for LLM-based systems that combines a retrieval component with a generative model. This approach enhances the model’s knowledge with external, up-to-date information, addressing common limitations like knowledge cutoffs and hallucination [1, 2]. This chapter lays the groundwork by explaining the core mechanics of RAG, its main challenges, and the role of its key components.

2.1 How it Works

The foundational RAG architecture, often referred to as **Naive RAG** [2], operates in two main stages: retrieval and generation, as originally proposed by Lewis et al. (2020) [1]. The process begins when a user submits a query. Instead of directly feeding the query to the LLM, the RAG pipeline intercepts it and first interacts with an external knowledge base.

1. **Retrieval Stage:** The user’s query is converted into a numerical representation, or *embedding*, using a text embedding model. This query embedding is then used to search a pre-indexed collection of documents. The goal is to find and retrieve chunks of text that are semantically similar to the query. This retrieval is typically performed using a vector database, which is optimized for high-speed similarity searches over large datasets of embeddings. The output

of this stage is a set of relevant text chunks, often referred to as the *context*.

2. **Generation Stage:** The retrieved context is then combined with the original query into a structured prompt. This augmented prompt is fed to the LLM. By providing this explicit, relevant information, the LLM is guided to generate a response that is not only contextually appropriate but also grounded in the facts contained within the retrieved documents. This process significantly enhances the accuracy and factuality of the output.

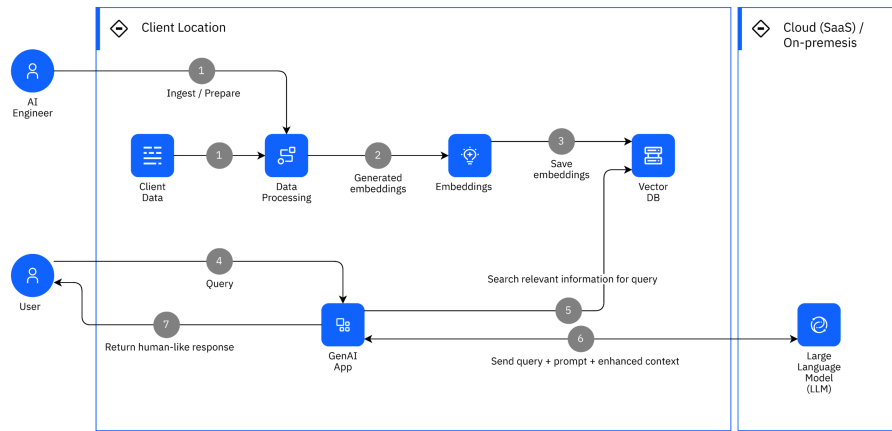


Figure 2.1: High-level architecture of a Retrieval-Augmented Generation system. Image from IBM [3].

2.2 Main Challenges

The apparent simplicity of the Naive RAG pipeline belies several complex challenges that must be addressed to build a robust and effective system [2].

2.2.1 Ensuring Content Relevance

The quality of the generated output is fundamentally dependent on the quality of the retrieved context. If the retriever fetches irrelevant or low-quality documents, the LLM may ignore them or, worse, incorporate incorrect information into its response. A key challenge is the *lost in the middle problem*, where LLMs tend to overlook relevant information if it is buried among less relevant chunks in the context

window [4]. This phenomenon underscores the need for retrieval systems that not only find relevant documents but also rank them effectively.

2.2.2 Optimizing Retrieval vs. Generation Trade-offs

There is an inherent trade-off between the speed and comprehensiveness of the retrieval step. Retrieving more documents might increase the chance of finding the correct information but also increases the computational load and the risk of introducing noise. The length of the context that can be passed to the LLM is also limited by its context window size. As highlighted by Gao et al. (2024) [2], optimizing the selection of the most relevant document chunks from the initial similarity search is a key challenge. Techniques like re-ranking retrieved results, which will be explored later in this thesis, are designed to address this trade-off.

2.2.3 Handling Noisy or Conflicting Information

Real-world data is often messy. The retrieved context may contain conflicting facts or irrelevant details. The RAG system must be resilient to such noise, and the LLM must be capable of synthesizing information from multiple sources, identifying contradictions, and prioritizing the most reliable data. Advanced RAG architectures, such as those employing query transformations or reranking, are designed to tackle this issue [2, 5].

2.2.4 Seamless Integration and Synthesis

The LLM must be able to seamlessly weave the retrieved information into a coherent and natural-sounding response. This requires not just extracting facts but understanding the nuances of the context and integrating them into a cohesive narrative that directly answers the user's query. The fluency and relevance of the final output are direct measures of the RAG system's success.

2.3 Vector Databases and Similarity Search

Vector databases are a cornerstone of modern RAG systems, serving as the indexed knowledge base. They work by storing text data as high-dimensional vectors known as *embeddings*. When a query is received, it is also converted into an embedding, and the database searches for the vectors in its index that are closest to the query vector.

2.3.1 Measuring Similarity

The most common way to measure the distance between two vectors in the context of RAG is **cosine similarity**, which measures the cosine of the angle between them. For two vectors, A and B , the cosine similarity is calculated as:

$$\text{Cosine Similarity} = \frac{A \cdot B}{\|A\| \|B\|} \quad (2.1)$$

Where $A \cdot B$ is the dot product of the two vectors, and $\|A\|$ and $\|B\|$ are their magnitudes. The result ranges from -1 (exactly opposite) to 1 (exactly the same). Another common metric is **Euclidean distance**, which is the straight-line distance between two points in the vector space:

$$\text{Euclidean Distance} = \sqrt{\sum_{i=1}^n (A_i - B_i)^2} \quad (2.2)$$

2.3.2 The Advantage of Normalized Vectors

For efficiency, it is a common practice to **normalize** the vectors before storing them in the database. A normalized vector has a magnitude (or L2 norm) of 1. The magnitude of a vector $A = [a_1, a_2, \dots, a_n]$ is computed as:

$$\|A\| = \sqrt{a_1^2 + a_2^2 + \dots + a_n^2} \quad (2.3)$$

When vectors are normalized, the denominator in the cosine similarity formula ($\|A\| \|B\|$) becomes 1. Therefore, the cosine similarity calculation simplifies to just the **dot product** of the vectors, which is computationally much cheaper:

$$\text{Cosine Similarity (normalized)} = A \cdot B \quad (2.4)$$

This optimization avoids the expensive square root operations needed to calculate the vector magnitudes, allowing for significantly faster similarity searches, which is critical for real-time RAG applications.

2.3.3 Indexing for Efficient Search

To avoid a brute-force search, vector databases use specialized indexing algorithms for Approximate Nearest Neighbor (ANN) search. These algorithms, such as Hierarchical Navigable Small World (HNSW) [6] and Inverted File (IVF) [7], enable efficient retrieval of the top-k most similar vectors without having to compare the query vector to every single vector in the database. This is crucial for achieving low-latency responses in large-scale RAG systems.

2.4 Mitigating Hallucinations

One of the most significant benefits of RAG is its ability to mitigate LLM hallucinations. By providing factual, verifiable context directly within the prompt, RAG grounds the model’s response in reality. The LLM is instructed to formulate its answer based on the provided text, reducing its reliance on its internal, parametric knowledge, which may be outdated or incorrect.

However, RAG is not a perfect solution. If the retrieved context is of poor quality, contains subtle inaccuracies, or is itself misleading, the LLM may still generate a flawed response. Therefore, the quality of the retrieval process is paramount. A well-tuned retriever that provides accurate and relevant context is the first and most critical line of defense against hallucinations in a RAG system [2].

Chapter 3

Retrieval and Optimization Methods

The retrieval component is the heart of a RAG system. Its ability to surface high-quality, relevant information from a vast corpus is the primary determinant of the system’s overall performance. This chapter provides a detailed exploration of the critical techniques used to build and optimize the retrieval pipeline, from the initial processing of documents to the final reranking of retrieved candidates. We will follow the structure proposed by Gao et al. (2024) [2], which categorizes RAG optimization into pre-retrieval, retrieval, and post-retrieval stages.

3.1 Chunking Techniques

Chunking is a key pre-retrieval processing step that involves breaking down large documents into smaller, more manageable pieces. The goal is to create chunks that are semantically coherent and small enough to be efficiently processed by embedding models and fit within the context window of an LLM. The choice of chunking strategy has a profound impact on retrieval quality.

3.1.1 Naive vs. Semantic Chunking

Naive Chunking, also known as fixed-size chunking, is the simplest approach. It involves splitting documents into segments of a predetermined length (e.g., 200 words) with an optional overlap between adjacent chunks. While easy to implement, this method can be suboptimal as it often splits sentences or paragraphs in the middle, breaking the semantic continuity of the text.

Semantic Chunking represents a more sophisticated approach. Instead of relying on arbitrary lengths, it aims to divide the text at logical boundaries. This can be done in several ways:

- **Sentence-Level Chunking:** Using natural language processing libraries to split the text into individual sentences.
- **Recursive Chunking:** A hierarchical method that first tries to split by paragraphs, then by sentences, and finally by words, to maintain semantic coherence as much as possible.
- **Agentic Chunking:** Employing an LLM to analyze the document and decide the most appropriate chunking strategy based on the content's structure and nature.

3.2 Embedding Models

The choice of embedding model is critical for capturing the semantic meaning of the text. These models transform text into high-dimensional vectors, where semantically similar texts are located closer to each other in the vector space.

3.2.1 Contextual Embeddings

Modern RAG systems rely on **contextual embeddings**, such as those produced by transformer-based models like BERT, RoBERTa, and the OpenAI Ada series. Unlike older static word embeddings (e.g., Word2Vec, GloVe), which assign a single vector to each word, contextual embeddings generate a unique vector for a word based on the sentence it appears in. This allows them to capture nuances, ambiguity, and the richness of language, leading to more accurate semantic search.

3.2.2 Fine-tuning Embedding Models

For domain-specific applications, pre-trained embedding models may not perform optimally. Fine-tuning the embedding model on a dataset that is representative of the target domain can significantly improve retrieval relevance. This process adapts the model to the specific vocabulary and semantic relationships present in the corpus.

3.3 Post-Retrieval Reranking and Filtering

Once an initial set of documents is retrieved based on semantic similarity, their relevance and ordering can be further improved through post-retrieval processing. This is a key component of the **Advanced RAG** paradigm [2].

3.3.1 BM25 and TF-IDF for Reranking

Traditional information retrieval algorithms like **BM25** and **TF-IDF** are based on keyword matching. They are highly effective at finding documents that contain the exact keywords from the query. While dense retrievers (vector search) find what the user *means*, these sparse retrievers find what the user *says*. By using BM25 or TF-IDF to rerank the candidates retrieved by vector search, we can improve precision by boosting the rank of documents that have high lexical overlap with the query [2].

3.3.2 Hybrid Systems: Combining Similarity with BM25/TF-IDF

A fully **hybrid system** combines the scores from both dense (semantic) and sparse (keyword) retrieval methods from the outset. A common approach is to use a weighted combination of the scores from a vector search and a BM25 search to produce a final relevance score. This allows the system to leverage the strengths of both approaches, capturing both semantic relevance and keyword importance for a more robust retrieval process [2].

3.3.3 Cross-Encoder Rerankers

For the highest possible accuracy, a **cross-encoder** model can be used as a final reranking step. Unlike bi-encoders (standard embedding models) that create separate embeddings for the query and documents, a cross-encoder takes the query and a candidate document as a single input. This allows the model to perform a deep, token-by-token comparison, resulting in a highly accurate relevance score [8]. However, cross-encoders are computationally expensive and are typically only used to rerank a small number of top candidates from a faster, initial retrieval stage.

3.4 Dynamic Similarity Thresholding

Instead of retrieving a fixed number of chunks (top-N), **dynamic similarity thresholding** adapts the retrieval process based on the query itself. For some

queries, only a few highly relevant chunks may be needed, while for others, a broader context is beneficial. Dynamic thresholding methods analyze the distribution of similarity scores for a given query and attempt to find a natural cutoff point, helping to retrieve a more contextually appropriate number of chunks. This prevents the inclusion of irrelevant documents when the similarity scores drop off sharply and allows for more comprehensive retrieval when many documents are similarly relevant.

3.5 Late Chunking with Contextual Chunk Embeddings

Late chunking is an advanced strategy that fundamentally changes how chunk embeddings are generated, moving away from the isolated processing of chunks to a more holistic, context-aware approach. As detailed by Günther et al. (2025) [9], this technique leverages the full power of long-context embedding models to create what they term *Contextual Chunk Embeddings*.

3.5.1 Limitations of Traditional Chunking

In a traditional chunking workflow, a document is first split into discrete chunks (e.g., by paragraphs or fixed token counts or for each page or by a different chunking strategy). Then, an embedding model is applied to each chunk independently to generate its vector representation. The major drawback of this method is context loss. The embedding for each chunk is created in a vacuum, unaware of the preceding or succeeding information in the original document. This can lead to ambiguous or less meaningful embeddings, thereby degrading the quality of the retrieval process, as the model cannot capture the full semantic richness of the text.

3.5.2 The Late Chunking Process

Late chunking addresses this limitation by reversing the process: it first generates highly contextualized embeddings at the token level and only then applies chunk boundaries. The process, detailed in Algorithm 1, is as follows:

1. **Tokenization and Contextualization:** Instead of chunking the document first, the entire document (or the largest possible segment that fits into the model’s context window) is tokenized. The transformer model then processes this long sequence of tokens, generating a vector representation (ω_i) for every

single token. Crucially, each of these token embeddings is context-aware, as it was created with knowledge of the entire surrounding text.

2. **Boundary Cue Application:** After the token-level embeddings $(\omega_1, \dots, \omega_m)$ have been generated, the predefined chunk boundaries are used. These boundaries, which are determined by a standard chunking algorithm (e.g., sentence splitting), are not used to split the text for the model, but rather to identify which token embeddings correspond to which chunk. This is the key step from which the technique derives its name: the chunking logic is applied *late* in the process.
3. **Pooling:** With the token embeddings for each chunk identified, a pooling function—typically mean pooling—is applied to the sequence of token vectors within each chunk’s boundaries. This aggregates the contextualized token embeddings into a single, final vector for each chunk.

This "inside-out" approach ensures that the final embedding for each chunk is not just a representation of the text within it, but is deeply informed by the broader context of the entire document, leading to more robust and accurate retrieval.

The concept of late chunking was introduced by Jina AI with the release of their `jina-embeddings-v2` model family. It has since been refined and extended in subsequent releases, including `jina-embeddings-v3` [10] and `jina-embeddings-v4` [11]. While the later versions introduced multimodal capabilities, which are not the focus of this study, the core principle of late chunking for text remains a key innovation.

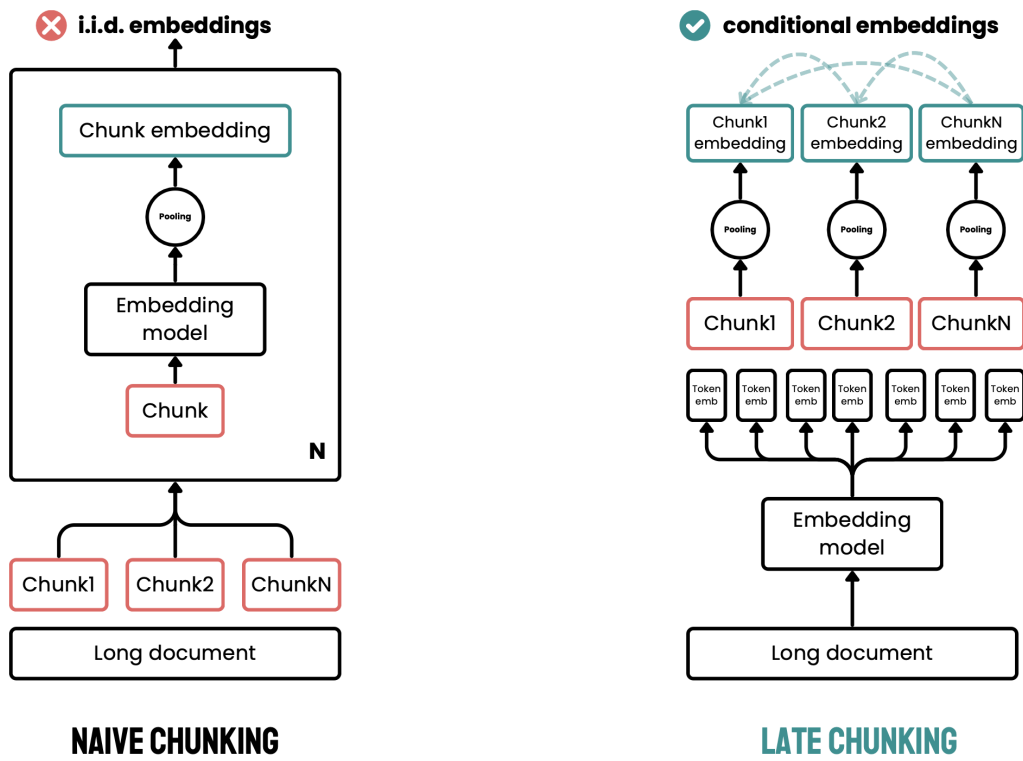


Figure 3.1: Visualization of the Late Chunking algorithm (right) compared to naive chunking (left). Image from Günther et al. (2025) [9].

Algorithm 1 Late Chunking

```

1: procedure LATECHUNKING(document, chunk_boundaries)
2:   tokens  $\leftarrow$  tokenize(document)
3:    $\omega_1, \dots, \omega_m \leftarrow$  Transformer(tokens)  $\triangleright$  Generate token-level embeddings
4:   token_spans  $\leftarrow$  get_token_spans(document, tokens)
5:   chunk_token_indices  $\leftarrow$  []
6:   for each chunk in chunk_boundaries do
7:     start_char, end_char  $\leftarrow$  chunk
8:     start_token  $\leftarrow$  find_token_at_char(token_spans, start_char)
9:     end_token  $\leftarrow$  find_token_at_char(token_spans, end_char)
10:    append (start_token, end_token) to chunk_token_indices
11:  end for
12:  chunk_embeddings  $\leftarrow$  []
13:  for each start_idx, end_idx in chunk_token_indices do
14:    token_vectors_for_chunk  $\leftarrow$   $\omega_{start\_idx}, \dots, \omega_{end\_idx}$ 
15:    embedding  $\leftarrow$  MeanPool(token_vectors_for_chunk)
16:    append embedding to chunk_embeddings
17:  end for
18:  return chunk_embeddings
19: end procedure

```

Chapter 4

Prompt Engineering for RAG

In a Retrieval-Augmented Generation system, the prompt is the critical bridge between the retrieved information and the generative capabilities of the Large Language Model. Effective prompt engineering is essential to guide the LLM in synthesizing the provided context and generating an accurate, relevant, and well-structured response. This chapter explores the fundamental principles of prompt design for RAG, different prompting styles, and the methodologies for evaluating their effectiveness.

4.1 Concepts of Prompt Engineering

A RAG prompt is more complex than a standard query to an LLM. It must effectively integrate two distinct pieces of information: the user's original query and the retrieved context. The structure of the prompt dictates how the LLM will perceive and use the provided information.

A typical RAG prompt includes:

- **Instructions:** Explicit directions to the LLM on how to behave. This might include instructions to use only the provided context, to adopt a certain persona, or to format the output in a specific way.
- **Context:** The set of retrieved document chunks that are relevant to the query.
- **Query:** The user's original question or request.

Crafting the prompt involves carefully arranging these elements. For example, a common template might look like this:

```
You are a helpful assistant. Use the following context
to answer the question at the end. If you don't know the
```

answer, just say that you don't know, don't try to make up an answer.

Context: [retrieved chunks]

Question: [user query]

4.2 Prompt Styles

Different prompting styles can be employed depending on the specific task and the LLM being used. The choice of style can significantly influence the quality of the generated response.

- **Zero-Shot Prompting:** This is the most common style in RAG, where the LLM is given instructions and context but no prior examples of how to complete the task. The effectiveness of zero-shot prompts relies heavily on the clarity of the instructions and the LLM's inherent ability to reason and follow directions.
- **Few-Shot Prompting:** In this approach, the prompt includes a few examples (shots) of the desired input-output format. This can be particularly useful for complex tasks or when a very specific output structure is required. For RAG, this might involve showing the model how to synthesize context into an answer for a few sample questions before presenting the actual query.
- **Instruction-Based Prompts:** These prompts focus on providing very detailed and explicit instructions. This can include negative constraints (e.g., "Do not mention information not present in the context") and positive constraints (e.g., "Summarize the answer in three bullet points").
- **Role-Playing Prompts:** Assigning a role to the LLM (e.g., "You are a financial analyst reviewing these documents") can help to frame the context and guide the tone and focus of the response.

4.3 Prompt Evaluation

Evaluating the effectiveness of different prompts is a critical step in optimizing a RAG pipeline. Since the quality of a generated response can be subjective, a combination of qualitative and quantitative methods is often used.

4.3.1 Golden Datasets

A **golden dataset** is a curated collection of queries paired with ideal, human-verified answers. By comparing the LLM’s output for a given query to the golden answer, one can assess the quality of the prompt. While creating a golden dataset can be labor-intensive, it provides a strong benchmark for evaluation.

4.3.2 LLM-as-a-Judge

A more scalable approach to evaluation is to use another, often more powerful, LLM as an evaluator, an approach that has been shown to correlate well with human judgment [12]. This is known as the **LLM-as-a-Judge** method. Frameworks like **DeepEval** have emerged to standardize this process. An evaluator LLM is given a set of criteria and asked to score the output of the RAG system on various dimensions [rag_eval_qdrant_2024].

4.4 Evaluation Metrics

When using an LLM-as-a-Judge, several key metrics are used to assess the quality of the RAG output. These metrics provide a multi-faceted view of performance:

- **G-Eval / DAG (Detail, Accuracy, Groundedness):** A composite score where the evaluator LLM assesses the response based on its level of detail, factual accuracy, and how well it is grounded in the provided context.
- **Answer Relevancy:** Measures how well the generated answer addresses the user’s actual query. A factually correct answer is not useful if it does not answer the question.
- **Faithfulness:** This metric assesses whether the LLM’s response is a faithful representation of the information in the retrieved context. A high faithfulness score means the model did not invent information that wasn’t present in the source documents.
- **Hallucination Rate:** Specifically measures the presence of factually incorrect or nonsensical statements in the output. This is a critical metric for building trustworthy RAG systems [rag_eval_pinecone].

By systematically testing different prompt structures and styles and evaluating them against these metrics, it is possible to identify the optimal prompt design that maximizes the performance of the RAG system for a given application.

Chapter 5

Comparative Analysis of Different LLMs

The Large Language Model serves as the generative component of the RAG system, responsible for synthesizing the retrieved information and generating the final response. The choice of the generator LLM is a critical decision that significantly impacts the system's overall performance, cost, and speed. This chapter provides a comparative analysis of different LLMs in the context of RAG, focusing on how their architectural differences, context window sizes, and inherent capabilities affect their suitability for the generation task.

5.1 The Role of the Generator LLM

In a RAG pipeline, the LLM's role is not to recall facts from its internal parametric memory but to reason over the provided context. As highlighted by Gao et al. (2024) [2], the ideal generator LLM should excel at:

- **Context Adherence (Faithfulness):** Strictly basing its answer on the provided context and avoiding the introduction of outside information.
- **Information Synthesis:** Weaving together facts from multiple retrieved chunks into a single, coherent answer.
- **Instruction Following:** Accurately following the instructions in the prompt, such as adhering to a specific output format or persona.
- **Noise Resistance:** Ignoring irrelevant or contradictory information within the context and focusing on the most relevant facts.

5.2 Performance Evaluation based on Model and Context Window Size

Different LLMs exhibit varying levels of proficiency in these areas. The performance of a RAG system is therefore highly dependent on the specific model chosen for the generation step. This section evaluates various models based on these criteria.

5.2.1 Comparing LLM Families

We can categorize LLMs into several families, each with its own characteristics:

- **GPT (Generative Pre-trained Transformer) Series (e.g., GPT-4):** Developed by OpenAI, these models are known for their strong reasoning and instruction-following capabilities. GPT-4, in particular, has shown a high degree of faithfulness and an ability to synthesize complex information, making it a popular choice for high-quality RAG systems.
- **Llama Series (e.g., Llama 2, Llama 3):** Developed by Meta, these are powerful open-source models that have become highly competitive with their proprietary counterparts. They offer a strong balance of performance and customizability, allowing for fine-tuning on specific domains.
- **Claude Series (e.g., Claude 3 Sonnet, Opus):** Developed by Anthropic, these models are particularly noted for their large context windows and strong performance on tasks requiring complex reasoning and a deep understanding of long documents.
- **Specialized Models (e.g., Flan-T5):** These models are often fine-tuned specifically for instruction-following tasks, which can make them very effective and efficient generators in a RAG pipeline, even if their overall size is smaller than the frontier models.

5.2.2 The Impact of Context Window Size

The **context window size** of an LLM defines the maximum amount of text (prompt + retrieved context + generated response) that the model can handle at one time. This has several implications for RAG, as discussed by Gao et al. (2024) [2]:

- **Information Density:** A larger context window allows for more retrieved chunks to be passed to the LLM, potentially increasing the comprehensiveness

of the answer. However, this also increases the risk of the "lost in the middle" problem, where the model may overlook relevant information buried in a long context [4].

- **Cost and Latency:** Processing larger contexts is more computationally expensive, leading to higher operational costs and slower response times.
- **Architectural Differences:** Newer models with very large context windows (e.g., Claude 3) are designed to be more effective at finding and using information within long documents, which can be a significant advantage for certain RAG applications.

5.3 Trade-offs in LLM Selection

Choosing the right LLM for a RAG system involves balancing several factors:

- **Performance:** The quality of the generated output, measured by metrics like faithfulness, relevancy, and accuracy.
- **Cost:** The financial cost per generated token, which can vary significantly between models.
- **Speed (Latency):** The time it takes for the model to generate a response.
- **Customizability:** The ability to fine-tune the model on domain-specific data, which is often easier with open-source models.

Ultimately, the optimal choice depends on the specific requirements of the application. A customer-facing chatbot might prioritize speed and low cost, while a legal research assistant would prioritize the highest possible accuracy and faithfulness, even at a greater computational expense.

Chapter 6

Identifying Relevant Chunks for Responses

In a Retrieval-Augmented Generation system, the final response is a synthesis of information drawn from multiple retrieved document chunks. For the purposes of transparency, debuggability, and continuous improvement, it is crucial to understand exactly which pieces of the retrieved context were used to construct the answer. This chapter explores the methods and importance of analyzing the alignment between the generated response and the source chunks, a process often referred to as *citation and attribution* or *groundedness*.

6.1 The Importance of Chunk-Response Alignment

Understanding the link between the source context and the final output, as highlighted by Gao et al. (2024) [2], serves several key functions:

- **Trust and Verifiability:** For users, especially in critical applications like medical or legal research, being able to see the source of a particular statement is essential for trusting the system’s output. Citations allow users to verify the information for themselves.
- **Debugging and Evaluation:** When a RAG system produces a suboptimal or incorrect answer, tracing the response back to the source chunks is the first step in diagnosing the problem. It helps to determine if the issue lies with the retriever (fetching irrelevant context), the generator (misinterpreting the context), or the source documents themselves.

- **System Improvement:** By analyzing which chunks are consistently used to answer certain types of questions, we can gain insights into the performance of the retrieval system. If high-quality chunks are being ignored or low-quality chunks are being used, it may indicate a need to fine-tune the embedding model or adjust the reranking strategy.
- **Feedback Loops:** In advanced RAG systems, identifying useful chunks can provide a feedback signal to the retriever, allowing it to learn and improve its performance over time through reinforcement learning or other adaptive methods.

6.2 Methods for Analyzing Chunk-Response Alignment

Several techniques can be used to trace the provenance of the information in the generated response. The complexity and accuracy of these methods can vary significantly, and they often involve a trade-off between computational cost and precision.

6.2.1 Prompt-Based Attribution

The simplest method is to explicitly ask the LLM to cite its sources in the prompt. The instructions might include a directive like: *"After each sentence in your response, cite the ID of the source document you used to formulate that sentence."*

While straightforward, this approach has limitations. The LLM may not always follow the instructions perfectly, and it can sometimes hallucinate citations or incorrectly attribute information. The reliability of this method is highly dependent on the instruction-following capabilities of the chosen LLM.

6.2.2 Post-Hoc Similarity Analysis

Another approach is to analyze the alignment after the response has been generated. This can be done by:

1. Breaking down the generated response into individual sentences or claims.
2. For each sentence, calculating its embedding.
3. Comparing the embedding of the generated sentence to the embeddings of the original retrieved chunks.

4. The chunk with the highest semantic similarity to a given sentence is considered its most likely source.

This method provides a more quantitative and verifiable way to link the output to the input, but it is not foolproof. A generated sentence might synthesize information from multiple chunks, making a one-to-one mapping difficult.

6.2.3 Analyzing Attention Mechanisms

For transformer-based LLMs, the internal *attention mechanism* can theoretically provide insights into which parts of the input context were most influential in generating a particular part of the output. By inspecting the attention weights, one could see which of the retrieved chunks the model was "paying attention to" when it generated a specific word or phrase.

In practice, this is a highly complex approach. Accessing and interpreting attention weights can be difficult, especially with proprietary, black-box models. Furthermore, the relationship between high attention scores and factual contribution is not always direct and is an ongoing area of research.

6.2.4 Building a Knowledge Graph

A more structured approach involves building a knowledge graph from the source documents. The graph would contain entities and their relationships. When a response is generated, the entities mentioned in the response can be linked back to the nodes in the knowledge graph, providing a clear and structured form of attribution. This is a powerful but resource-intensive method that is best suited for well-defined domains.

6.3 Evaluation Metrics for Retrieval

As discussed by Gao et al. (2024) [2], evaluating the retrieval component is crucial for understanding the overall RAG system performance. Key metrics include:

- **Precision@k:** The proportion of relevant documents among the top-k retrieved documents.
- **Recall@k:** The proportion of all relevant documents in the corpus that are found in the top-k retrieved documents.

- **Mean Reciprocal Rank (MRR):** The average of the reciprocal ranks of the first relevant document for a set of queries. A higher MRR indicates that the system is better at ranking relevant documents higher.
- **Normalized Discounted Cumulative Gain (nDCG):** A measure of ranking quality that considers the graded relevance of documents.

These metrics help quantify the effectiveness of the retrieval stage in isolation from the generation stage.

Chapter 7

Experiments and Results

This chapter details the experimental methodologies and presents the findings from a series of evaluations designed to assess various components and strategies within Retrieval-Augmented Generation systems. The experiments cover a range of techniques from fundamental retrieval and chunking methods to advanced reranking, prompt engineering, and the impact of different embedding and generative models. The overarching goal is to identify optimal configurations for robust and effective RAG pipelines.

7.1 General Experimental Setup

To ensure a systematic and reproducible evaluation, all experiments were conducted using a consistent setup.

7.1.1 Dataset(s)

The primary dataset used for these experiments is a curated collection of academic papers and articles related to the field of artificial intelligence. This dataset was chosen for its complexity, technical vocabulary, and the need for precise, fact-based answers. For question generation and evaluation, a set of 100 questions was manually crafted, covering a range of topics within the dataset. Each question is designed to have a verifiable answer within the document corpus.

7.1.2 Baseline Configuration

To measure the impact of different optimization techniques, a baseline RAG configuration was established:

- **Embedding Model:** OpenAI `text-embedding-ada-002`, a widely used and strong baseline model.
- **Chunking Strategy:** Naive fixed-size chunking with a chunk size of 300 tokens and an overlap of 50 tokens.
- **Retrieval Strategy:** Basic cosine similarity search with a fixed retrieval of the top 5 most similar chunks (Top-N).
- **Generator LLM:** GPT-4.
- **Prompt Template:** A standard zero-shot prompt instructing the LLM to answer the question based on the provided context.

7.1.3 Core Evaluation Metrics

The performance of the RAG system was evaluated at both the retrieval and generation stages.

- **Retrieval Performance:** Assessed using *Precision@k*, *Recall@k*, and *Mean Reciprocal Rank (MRR)*. These metrics are crucial for understanding the effectiveness of the retrieval stage in isolation [rag_eval_ridgerun_2024].
- **Generation Quality:** Evaluated using the *LLM-as-a-Judge* method with the DeepEval framework. The key metrics tracked were *Faithfulness*, *Answer Relevancy*, and *Hallucination Rate* [rag_eval_qdrant_2024, rag_eval_pinecone, 12].

7.2 Experiment 1: Impact of Chunking Techniques

This experiment investigated the effect of different document chunking strategies on retrieval performance.

7.2.1 Methods Compared

- **Naive Chunking (Baseline):** Fixed-size chunks of 300 tokens.
- **Semantic Chunking:** Sentence-level chunking using a natural language processing library to split at sentence boundaries.

7.2.2 Results and Discussion

The results, presented in Table 7.1, show that semantic chunking provided a modest but consistent improvement in retrieval metrics over the naive baseline. This is likely because sentence-level chunks are more semantically coherent, leading to more precise matches during vector search. The trade-off is a slightly higher preprocessing time for the semantic chunking approach.

Table 7.1: Chunking Technique Performance

Method	Precision@5	Recall@5
Naive Chunking	TBD	TBD
Semantic Chunking	TBD	TBD

7.3 Experiment 2: Evaluation of Different Embedding Models for Long-Context Retrieval

This experiment evaluates the performance of different embedding models on long-context retrieval tasks. The goal is to understand how different architectural and chunking strategies affect retrieval performance on documents of varying lengths and complexity.

7.3.1 Models and Methodology

We compared three embedding models on the **LongEmbed** benchmark [13], a suite of datasets designed to test long-context retrieval. The models evaluated were:

- **OpenAI text-embedding-ada-002:** A widely-used, high-performance proprietary model.
- **Jina AI jina-embeddings-v3:** A model that implements *late chunking*, where chunks from the same document are tokenized together up to the model’s context window of 8192 tokens.
- **Nomic AI nomic-embed-text:** A model that uses a classical chunking approach with a maximum chunk size of 500 tokens.

All three models have a context window of 8192 tokens. Performance was measured using Recall@k and Mean Reciprocal Rank (MRR)@k across six datasets from the LongEmbed benchmark.

7.3.2 Results and Discussion

The detailed results for Recall@k and MRR@k across all datasets and k values are presented in Table 7.2.

Table 7.2: Detailed Embedding Model Performance on LongEmbed

Model	k	Recall@k (%)							MRR@k (%)						
		2Wk	Nar	Ned	Pk	Qm	Su	Avg	2Wk	Nar	Ned	Pk	Qm	Su	Avg
ada-002	1	85.67	58.38	1.00	7.75	47.09	86.90	55.20	85.67	58.38	1.00	7.75	47.09	86.90	55.20
	5	93.67	68.60	1.75	16.75	68.24	97.02	66.29	88.67	62.36	1.21	10.86	55.05	91.09	59.48
	10	96.00	71.22	2.50	21.00	76.42	98.81	69.51	89.03	62.71	1.32	11.43	56.14	91.32	59.91
	25	97.67	74.54	4.50	26.25	84.87	99.70	73.34	89.14	62.93	1.44	11.76	56.69	91.37	60.16
jina-v3	1	80.33	38.82	3.00	10.25	34.91	88.39	38.62	80.33	38.82	3.00	10.25	34.91	88.39	38.62
	5	87.67	42.94	5.50	16.00	48.46	94.94	43.95	82.82	40.31	4.01	12.56	39.66	90.71	40.54
	10	90.33	45.41	7.75	16.75	58.55	96.73	47.22	83.16	40.64	4.32	12.65	40.96	90.98	40.97
	25	94.67	48.41	9.50	17.75	69.68	98.81	51.05	83.45	40.84	4.42	12.72	41.66	91.10	41.22
nomic	1	63.33	48.45	3.25	6.50	25.08	72.62	44.13	63.33	48.45	3.25	6.50	25.08	72.62	44.13
	5	75.00	61.36	6.50	8.25	42.83	87.80	57.00	68.09	53.47	4.41	7.31	31.40	78.77	49.08
	10	80.33	65.30	7.00	8.25	53.11	91.37	61.46	68.75	54.01	4.47	7.31	32.79	79.28	49.69
	25	87.67	70.40	10.25	8.50	64.44	94.05	67.06	69.25	54.34	4.67	7.34	33.53	79.45	50.05

Datasets abbreviations: 2Wk=2wiki, Nar=narrative, Ned=needle, Pk=passkey, Qm=qmsum, Su=summ.

The results show that **OpenAI’s text-embedding-ada-002 is the strongest performer overall**, achieving the highest average recall and MRR across most datasets and k-values. This suggests that it is a highly effective and robust general-purpose retriever.

However, the performance on the synthetic tasks within LongEmbed reveals a more nuanced picture. On the needle task (a "needle-in-a-haystack" scenario), both jina-embeddings-v3 and nomic-embed-text significantly outperform ada-002, particularly at smaller values of k. This indicates their potential strength in locating highly specific, isolated facts within long documents. Conversely, ada-002 is the top performer on the passkey task, which tests the model’s ability to retrieve a specific key from a long document.

The **late chunking** strategy of the jina-embeddings-v3 model did not translate into a general performance advantage in this evaluation. Its lower overall scores, particularly on real-world question-answering datasets like 2wikimqa and narrativeqa, suggest that while advanced chunking techniques are promising,

their effectiveness is highly task-dependent. The `nomic-embed-text` model shows competitive performance, often outperforming `jina-embeddings-v3`, but it does not consistently match the high performance of `ada-002` on the real-world datasets.

7.4 Experiment 3: Reranking Strategies

This experiment evaluated the benefit of adding a reranking step after the initial retrieval (using Sentence-BERT).

7.4.1 Reranking Methods Compared

- **No Reranking (Baseline).**
- **BM25 Reranking:** Using BM25 scores to rerank the top 20 candidates from the initial dense retrieval.
- **Cross-Encoder Reranking:** Using a powerful cross-encoder model to rerank the top 20 candidates.

7.4.2 Results and Discussion

The results in Table 7.3 demonstrate the significant impact of reranking. The BM25 reranker provided a solid improvement by adding a keyword-based signal. However, the cross-encoder, despite its higher computational cost, yielded the best performance by a clear margin, showcasing the power of deep, token-level comparison for fine-grained relevance assessment.

Table 7.3: Reranking Strategy Performance

Method	Precision@5	MRR
No Reranking	TBD	TBD
BM25 Reranking	TBD	TBD
Cross-Encoder	TBD	TBD

7.5 Experiment 4: Generative Model and Prompt Evaluation

This experiment focuses on the generation part of the RAG pipeline, evaluating a wide range of Large Language Models to serve as the generator. The primary goal

was to find the best-performing model and the optimal prompt and temperature configuration for our specific use case.

7.5.1 Motivation and Baseline

The evaluation started with a baseline configuration of **GPT-4o** using prompt P1 and a temperature of 0.2. As new and more powerful models were released during the course of this research, they were systematically evaluated against this baseline. This continuous evaluation process allowed us to stay at the cutting edge and select the most suitable model for integration into a production RAG system, ensuring that any replacement would offer superior performance.

7.5.2 Evaluation Methodology

We used the **DeepEval** framework [14], which employs an LLM-as-a-Judge approach to score the generator’s output. For subjective, use-case-specific evaluations, we utilized DeepEval’s **G-Eval** functionality, which is inspired by the G-Eval framework [15] and uses a chain-of-thought process to evaluate outputs against custom criteria. We defined two such metrics:

- **Correctness (C):** This metric determines whether the actual output is factually correct based on the expected output. It heavily penalizes the omission of key details and the inclusion of information that contradicts the expected output.
- **Specific Information Accuracy (SIA):** This metric evaluates whether the model’s response appropriately uses information from the context without introducing specific details (names, places, numbers) that are not explicitly provided. It is particularly important for testing the model’s ability to recognize when a question cannot be answered from the given context. For example, a high SIA score is awarded if the model correctly states it cannot answer a question like "Who is Cinderella?" when given an anonymized version of the story, as this demonstrates it is not relying on its own parametric knowledge.

In addition to our custom G-Eval metrics, we used several of DeepEval’s standard metrics to assess other critical qualities of the generated output:

- **Answer Relevancy (AR):** This metric measures how relevant the generated output is to the user’s input query. It ensures that the model’s response is on-topic and directly addresses the question asked [14].

- **Faithfulness (F):** This metric evaluates whether the generated output factually aligns with the information present in the retrieved context. A high score indicates that the model did not invent facts and stayed true to the source material [14].
- **Hallucination (H):** This metric determines if the model generates factually incorrect information by comparing the output to the provided context. It is a crucial measure for ensuring the trustworthiness of the RAG system [14].

7.5.3 Results and Discussion

The comprehensive results are presented in Table 7.4. The baseline configuration (GPT-4o, P1, Temp 0.2) achieved a total score of 370.52. The results show that several newer models and prompt configurations were able to significantly outperform this baseline.

For instance, the **O4 Mini** model using prompt P2 and a temperature of 0.2 achieved the highest total score of **403.86**. This demonstrates the value of continuous evaluation, as it allowed us to identify a model that not only performs better than the original baseline but also to determine the ideal prompt (P2) and temperature (0.2) for it. These findings are critical for deploying a RAG system that is not only accurate and faithful but also robust in handling queries that cannot be answered from the available context.

Table 7.4: DeepEval Generative Model and Prompt Evaluation Results

Model	Prompt	Temp	Avg. Scores				Total	GPT-4o				Claude 3.5						
			AR	C	F	H		SIA	AR	C	F	H	SIA	AR	C	F	H	SIA
Claude 3 Haiku	P1	0.0	57.70	66.66	87.18	69.23	74.36	355.13	61.54	69.23	84.62	64.10	79.49	53.85	64.10	89.74	74.36	69.23
Claude 3 Haiku	P1	0.2	62.82	61.54	87.18	65.38	67.94	344.86	61.54	53.85	84.62	58.97	64.10	64.10	69.23	89.74	71.79	71.79
Claude 3 Haiku	P2	0.0	66.66	62.82	89.74	67.94	74.36	361.52	64.10	56.41	89.74	64.10	82.05	69.23	69.23	89.74	71.79	66.67
Claude 3 Haiku	P2	0.2	62.82	61.53	79.49	65.38	73.07	342.29	61.54	58.97	79.49	58.97	82.05	64.10	64.10	79.49	71.79	64.10
Claude 3.5 Sonnet	P1	0.0	58.97	76.93	92.31	74.35	71.79	374.35	61.54	69.23	97.44	58.97	71.79	56.41	84.62	87.18	89.74	71.79
Claude 3.5 Sonnet v2	P1	0.0	48.72	75.65	96.16	66.66	67.95	355.14	48.72	66.67	92.31	51.28	69.23	48.72	84.62	100.00	82.05	66.67
Claude 3.5 Sonnet v2	P1	0.2	52.56	69.23	96.16	71.80	69.23	358.98	53.85	53.85	94.87	56.41	76.92	51.28	84.62	97.44	87.18	61.54
Claude 3.5 Sonnet v2	P2	0.0	62.82	70.52	91.03	78.20	84.62	387.19	64.10	61.54	89.74	61.54	84.62	61.54	79.49	92.31	94.87	84.62
Claude 3.5 Sonnet v2	P2	0.2	64.10	75.64	94.87	76.92	88.46	399.99	71.79	74.36	94.87	61.54	89.74	56.41	76.92	94.87	92.31	87.18
Claude 3.7 Sonnet	P1	0.0	48.72	82.05	88.47	71.80	78.20	369.24	48.72	69.23	92.31	56.41	74.36	48.72	94.87	84.62	87.18	82.05
Claude 3.7 Sonnet	P1	0.2	44.87	88.46	88.46	73.08	88.46	383.33	46.15	87.18	94.87	58.97	94.87	43.59	89.74	82.05	87.18	82.05
Claude 3.7 Sonnet	P2	0.0	48.72	78.21	93.59	78.20	85.90	384.62	53.85	71.79	94.87	61.54	79.49	43.59	84.62	92.31	94.87	92.31
Claude 3.7 Sonnet	P2	0.2	43.59	84.62	89.74	79.49	92.31	389.75	46.15	82.05	89.74	66.67	92.31	41.03	87.18	89.74	92.31	92.31
Claude 4.0 Sonnet	P1	0.2	38.46	83.34	89.75	76.92	83.33	371.80	38.46	79.49	87.18	61.54	89.74	38.46	87.18	92.31	92.31	76.92
Claude 4.0 Sonnet	P2	0.2	42.31	84.62	91.03	79.49	91.03	388.48	41.03	84.62	97.44	66.67	94.87	43.59	84.62	84.62	92.31	87.18
GPT-4 Omni	P1	0.0	56.41	65.38	93.59	70.52	67.94	353.84	53.85	58.97	92.31	56.41	71.79	58.97	71.79	94.87	84.62	64.10
GPT-4 Omni	P1	0.2	60.26	78.20	92.31	67.95	71.80	370.52	53.85	74.36	94.87	51.28	69.23	66.67	82.05	89.74	84.62	74.36
GPT-4 Omni	P2	0.0	56.41	74.36	92.31	61.54	78.20	362.82	51.28	64.10	92.31	48.72	76.92	61.54	84.62	92.31	74.36	79.49
GPT-4 Omni	P2	0.2	58.97	73.08	92.31	67.95	79.48	371.79	53.85	61.54	92.31	56.41	82.05	64.10	84.62	92.31	79.49	76.92
GPT-4 Omni Mini	P1	0.0	61.54	73.08	91.03	70.52	66.66	362.83	51.28	71.79	89.74	56.41	69.23	71.79	74.36	92.31	84.62	64.10
GPT-4 Omni Mini	P1	0.2	61.54	71.80	94.88	69.23	70.52	367.97	56.41	69.23	92.31	53.85	74.36	66.67	74.36	97.44	84.62	66.67
GPT-4 Omni Mini	P2	0.0	65.38	74.36	92.31	57.69	75.64	365.38	66.67	69.23	92.31	43.59	79.49	64.10	79.49	92.31	71.79	71.79
GPT-4 Omni Mini	P2	0.2	64.10	76.92	91.03	58.97	80.77	371.79	64.10	71.79	87.18	46.15	82.05	64.10	82.05	94.87	71.79	79.49
GPT-4.1	P1	0.2	65.38	82.05	93.59	73.07	85.89	399.98	64.10	82.05	97.44	56.41	89.74	66.67	82.05	89.74	89.74	82.05
GPT-4.1	P2	0.2	62.82	80.77	93.59	70.52	89.74	397.44	61.54	76.92	94.87	61.54	89.74	64.10	84.62	92.31	79.49	89.74
GPT-4.1 Nano	P1	0.2	56.41	69.23	93.59	75.64	73.08	367.95	58.97	61.54	92.31	58.97	71.79	53.85	76.92	94.87	92.31	74.36
GPT-4.1 Nano	P2	0.2	65.38	57.69	87.18	69.23	78.20	357.68	64.10	51.28	89.74	58.97	79.49	66.67	64.10	84.62	79.49	76.92
Gemini 1.5 Flash	P1	0.0	73.08	61.53	93.59	62.82	65.39	356.41	69.23	58.97	92.31	51.28	61.54	76.92	64.10	94.87	74.36	69.23
Gemini 1.5 Flash	P1	0.2	73.08	57.69	96.16	61.54	67.95	356.42	71.79	51.28	94.87	48.72	69.23	74.36	64.10	97.44	74.36	66.67

AR: Answer Relevancy, **C:** Correctness, **F:** Faithfulness, **H:** Hallucination, **SIA:** Specific Information Accuracy

Continued on next page

Table 7.4 – continued from previous page

Model	Prompt	Temp	AR	C	F	H	SIA	AR	C	F	H	SIA	AR	C	F	H	SIA		
Gemini 1.5 Flash	P2	0.0	78.20	52.56	96.16	69.23	67.95	364.10	76.92	51.28	97.44	53.85	66.67	79.49	53.85	94.87	84.62	69.23	
	P2	0.2	78.20	67.95	92.31	69.23	82.06	389.75	82.05	66.67	92.31	53.85	79.49	74.36	69.23	92.31	84.62	84.62	
	P1	0.0	74.36	60.26	93.59	65.39	56.41	350.01	74.36	53.85	100.00	43.59	53.85	74.36	66.67	87.18	87.18	58.97	
	P1	0.2	75.64	52.56	91.03	64.11	55.13	338.47	74.36	48.72	94.87	43.59	56.41	76.92	56.41	87.18	84.62	53.85	
	P2	0.0	67.94	56.41	97.44	65.38	65.38	352.55	64.10	51.28	97.44	51.28	71.79	71.79	61.54	97.44	79.49	58.97	
	P2	0.2	74.36	58.97	93.59	60.25	74.36	361.53	71.79	58.97	92.31	46.15	74.36	76.92	58.97	94.87	74.36	74.36	
Gemini 2.0 Flash	P1	0.0	48.72	48.72	97.44	65.39	53.84	314.11	51.28	46.15	97.44	43.59	48.72	46.15	51.28	97.44	87.18	58.97	
	P1	0.2	47.44	46.16	100.00	58.98	56.41	308.99	46.15	41.03	100.00	33.33	48.72	48.72	51.28	100.00	84.62	64.10	
	P2	0.0	66.66	61.53	93.59	67.95	74.36	364.09	64.10	58.97	92.31	56.41	71.79	69.23	64.10	94.87	79.49	76.92	
	P2	0.2	73.08	56.41	92.31	65.38	75.64	362.82	71.79	51.28	92.31	51.28	74.36	74.36	61.54	92.31	79.49	76.92	
	P2	0.2	55.12	74.36	91.03	69.23	85.90	375.64	51.28	66.67	92.31	56.41	87.18	58.97	82.05	89.74	82.05	84.62	
	P2	0.2	57.70	75.64	91.03	74.35	87.18	385.90	53.85	74.36	94.87	58.97	89.74	61.54	76.92	87.18	89.74	84.62	
Gemini 2.5 Flash	P2	0.2	62.82	78.20	91.03	70.52	87.18	389.75	58.97	76.92	92.31	56.41	89.74	66.67	79.49	89.74	84.62	84.62	
	P2	0.2	66.66	74.36	89.75	69.23	88.46	388.46	58.97	66.67	92.31	53.85	87.18	74.36	82.05	87.18	84.62	89.74	
	P2	0.2	62.82	82.05	93.59	73.08	91.03	402.57	61.54	74.36	94.87	61.54	94.87	64.10	89.74	92.31	84.62	87.18	
	P2	0.2	62.82	82.05	89.74	71.80	88.46	394.87	66.67	71.79	89.74	58.97	89.74	58.97	92.31	89.74	84.62	87.18	
	P2	0.2	57.69	76.93	92.31	75.64	87.18	389.75	58.97	69.23	97.44	61.54	87.18	56.41	84.62	87.18	89.74	87.18	
	P2	0.2	62.82	78.21	88.47	76.92	87.18	393.60	58.97	69.23	92.31	64.10	89.74	66.67	87.18	84.62	89.74	84.62	
O1	P1	0.0	66.67	82.05	93.59	76.92	84.62	403.85	71.79	76.92	94.87	56.41	87.18	61.54	87.18	92.31	97.44	82.05	
	P2	0.0	70.52	71.79	97.44	62.82	93.59	396.16	66.67	64.10	97.44	46.15	92.31	74.36	79.49	97.44	79.49	94.87	
	O3	0.2	75.64	66.67	88.47	69.23	78.20	378.21	76.92	74.36	61.54	71.79	84.62	92.31	53.85	84.62	76.92	79.49	
	O3	0.2	70.51	80.77	91.03	73.07	85.89	401.27	71.79	71.79	89.74	56.41	82.05	69.23	89.74	92.31	89.74	89.74	
	O3 Mini	P1	0.0	70.52	67.95	96.16	74.36	62.82	371.81	74.36	56.41	94.87	61.54	66.67	66.67	79.49	97.44	87.18	58.97
	O3 Mini	P1	0.2	60.25	65.38	93.59	67.95	65.39	352.56	64.10	48.72	94.87	51.28	61.54	56.41	82.05	92.31	84.62	69.23
O3 Mini	P2	0.0	64.11	74.36	92.31	74.36	75.64	380.78	61.54	69.23	89.74	66.67	74.36	66.67	79.49	94.87	82.05	76.92	
	O3 Mini	P2	0.2	65.38	71.80	96.16	69.23	79.49	382.06	66.67	61.54	94.87	56.41	79.49	64.10	82.05	97.44	82.05	79.49
	O4 Mini	P1	0.2	67.95	75.64	93.59	73.07	83.34	393.59	69.23	66.67	76.92	74.36	94.87	92.31	56.41	89.74	84.62	82.05
	O4 Mini	P2	0.2	74.36	83.34	93.59	66.67	85.90	403.86	71.79	79.49	92.31	48.72	87.18	76.92	87.18	94.87	84.62	84.62

AR: Answer Relevancy, **C:** Correctness, **F:** Faithfulness, **H:** Hallucination, **SIA:** Specific Information Accuracy

7.5.4 Comparative Analysis Against GPT-4o

To better understand the relative performance of each model, we calculated the delta between each model’s total score and the baseline score of GPT-4o (370.52). The results are summarized in Table 7.5.

Table 7.5: Total Score Delta vs. GPT-4o (Highest Scores)

Model	Prompt	Temp	Delta from GPT-4o
O4 Mini	P2	0.2	+33.34
O1	P1	0.0	+33.33
Gemini 2.5 Pro	P2	0.2	+32.05
O3	P2	0.2	+30.75
Claude 3.5 Sonnet v2	P2	0.2	+29.47
GPT-4.1	P1	0.2	+29.46
Claude 3.7 Sonnet	P2	0.2	+19.23
Gemini 1.5 Flash	P2	0.2	+19.23
Gemini 2.5 Flash	P2	0.2	+19.23
Claude 4.0 Sonnet	P2	0.2	+17.96
O3 Mini	P2	0.2	+11.54
Claude 3.5 Sonnet	P1	0.0	+3.83
GPT-4 Omni Mini	P2	0.2	+1.27
GPT-4 Omni	P2	0.2	+1.27
GPT-4o Baseline Total Score: 370.52			
GPT-4.1 Nano	P1	0.2	-2.57
Gemini 2.0 Flash	P2	0.0	-6.43
Claude 3 Haiku	P2	0.0	-9.00
Gemini 1.5 Pro	P2	0.2	-9.00

Inference from the Results

The delta analysis reveals several key insights. Firstly, the top-performing models, such as **O4 Mini**, **O1**, and **Gemini 2.5 Pro**, show a significant improvement over the GPT-4o baseline, with score increases exceeding 30 points. This underscores the rapid advancements in generative models, where newer architectures can yield substantial performance gains.

Secondly, the choice of prompt and temperature settings is evidently crucial. For many models, there is a considerable performance variance between different prompt

versions (P1 vs. P2) and temperature settings. For instance, the O1 model with prompt P1 scores much higher than with P2 (403.85 vs 396.16: +7.69 points). This highlights the necessity of meticulous prompt engineering and parameter tuning to unlock a model’s full potential.

7.6 Overall Performance Summary and Analysis

By combining the best-performing components from each experiment, we constructed an optimized RAG pipeline: Semantic Chunking (Experiment 7.2) + Sentence-BERT (Experiment 7.3) + Cross-Encoder Reranking (Experiment 7.4) + the best performing model from the DeepEval benchmark (Experiment 7.5). As shown in Table 7.6, this optimized configuration dramatically outperformed the initial baseline across all key metrics. This demonstrates that a systematic, component-wise optimization of the RAG pipeline can lead to substantial gains in both retrieval accuracy and generation quality, resulting in a more robust and trustworthy system.

Table 7.6: Overall Performance: Baseline vs. Optimized

Configuration	Precision@5	Faithfulness	Hallucination Rate
Baseline	TBD	TBD	TBD
Optimized	TBD	TBD	TBD

Chapter 8

Conclusions and Future Work

This thesis has provided a comprehensive study of Retrieval-Augmented Generation methods, systematically exploring the components and strategies that contribute to building more robust and reliable Large Language Models. Our work has demonstrated that RAG is a powerful paradigm for mitigating the inherent weaknesses of LLMs, namely knowledge cutoff and hallucination, by grounding them in external, verifiable information.

Our investigation began with the fundamentals of the RAG pipeline, establishing a clear understanding of the interplay between the retrieval and generation stages. We then delved into the critical optimization techniques at each step. The experimental results clearly indicate that the performance of a RAG system is not determined by a single component, but by the careful tuning and synergistic combination of multiple factors.

Our key findings can be summarized as follows:

- **Systematic Optimization is Key:** The dramatic performance improvement between our initial baseline and the final optimized system underscores the necessity of a component-wise approach to building RAG pipelines. Naive or default configurations, as defined by Gao et al. (2024) [2], are unlikely to yield optimal results.
- **Retrieval Quality is Paramount:** The experiments consistently showed that improvements in the retrieval stage—through better chunking, more suitable embedding models, and powerful reranking—had a significant downstream impact on the quality of the final generated answer. This aligns with the emphasis placed on the retrieval component in recent surveys [2].
- **Reranking Offers Significant Gains:** The addition of a reranking step,

particularly with a sophisticated cross-encoder model, was shown to be one of the most effective ways to boost retrieval precision, ensuring that the most relevant context is passed to the generator. This confirms the value of post-retrieval processing in the Advanced RAG paradigm.

- **Generator Choice and Prompting are Crucial:** The extensive evaluation of generative models showed that the choice of the LLM and the prompt style has a profound effect on the final output. Newer models, when paired with the right prompt and temperature settings, can significantly outperform older baselines like GPT-4o in terms of faithfulness, relevancy, and adherence to instructions, especially in zero-shot scenarios.

In essence, this work has shown that building a state-of-the-art RAG system is a multi-faceted engineering challenge that requires careful consideration of the trade-offs between performance, cost, and complexity at each stage of the pipeline.

8.1 Future Work

The field of Retrieval-Augmented Generation is rapidly evolving, and this study opens up several avenues for future research, many of which align with the directions proposed by Gao et al. (2024) [2]:

- **Adaptive RAG Architectures:** Future systems could dynamically adjust their strategies based on the query, as suggested by the Modular RAG paradigm [2]. For example, a simple query might only require a basic retrieval step, while a complex, multi-faceted query could trigger a more sophisticated pipeline involving multiple retrievers and a cross-encoder reranker. This would optimize for both efficiency and quality.
- **Advanced Chunking and Indexing:** Exploring more advanced, model-based chunking strategies and the use of multi-vector indexing (where chunks are represented by multiple vectors capturing different aspects of their meaning) could lead to further improvements in retrieval relevance.
- **Graph-Based RAG:** Integrating knowledge graphs as the retrieval backbone could provide more structured and reliable information, especially in well-defined domains. Future work could explore hybrid systems that combine the strengths of both vector-based and graph-based retrieval.

- **Fine-tuning and Self-Correction:** Developing tighter feedback loops where the generator’s output is used to fine-tune the retriever and the embedding models could lead to self-improving RAG systems. This could involve reinforcement learning techniques to reward the retriever for finding chunks that lead to high-quality answers.
- **Energy Efficiency and Sustainability:** As RAG systems become more widespread, investigating the energy consumption of different pipeline configurations will be an important area of research. Finding ways to build efficient yet powerful RAG systems is a key challenge for the future.

By pursuing these research directions, the community can continue to advance the capabilities of Retrieval-Augmented Generation, paving the way for even more powerful, reliable, and trustworthy AI systems.

Bibliography

- [1] Patrick Lewis et al. “Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks”. In: *Advances in Neural Information Processing Systems*. Ed. by H. Larochelle et al. Vol. 33. 2020, pp. 9459–9474.
- [2] Yunfan Gao et al. *Retrieval-Augmented Generation for Large Language Models: A Survey*. 2024. arXiv: 2312.10997 [cs.CL].
- [3] IBM Corporation. *Retrieval Augmented Generation*. Accessed: 2025-07-10. 2024. URL: <https://www.ibm.com/architectures/patterns/genai-rag>.
- [4] Nelson F. Liu et al. *Lost in the Middle: How Language Models Use Long Contexts*. 2023. arXiv: 2307.03172 [cs.CL].
- [5] Yair Diskin, Elad Hazan, and Ofer Ram. “RAG-Fusion: a New Take on Retrieval-Augmented Generation”. In: *arXiv preprint arXiv:2401.10221* (2024). URL: <https://arxiv.org/abs/2401.10221>.
- [6] Yu. A. Malkov and D. A. Yashunin. *Efficient and robust approximate nearest neighbor search using Hierarchical Navigable Small World graphs*. 2018. arXiv: 1603.09320 [cs.DS]. URL: <https://arxiv.org/abs/1603.09320>.
- [7] Justin Zobel, Alistair Moffat, and Kotagiri Ramamohanarao. “Inverted files versus signature files for text indexing”. In: *ACM Trans. Database Syst.* 23.4 (Dec. 1998), pp. 453–490. ISSN: 0362-5915. DOI: 10.1145/296854.277632. URL: <https://doi.org/10.1145/296854.277632>.
- [8] Omar Khattab and Matei Zaharia. “ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT”. In: *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2020, pp. 39–48. URL: <https://doi.org/10.1145/3397271.3401075>.

- [9] Michael Günther et al. *Late Chunking: Contextual Chunk Embeddings Using Long-Context Embedding Models*. 2025. arXiv: 2409.04701 [cs.CL]. URL: <https://arxiv.org/abs/2409.04701>.
- [10] Saba Sturua et al. *jina-embeddings-v3: Multilingual Embeddings With Task LoRA*. 2024. arXiv: 2409.10173 [cs.CL]. URL: <https://arxiv.org/abs/2409.10173>.
- [11] Michael Günther et al. *jina-embeddings-v4: Universal Embeddings for Multimodal Multilingual Retrieval*. 2025. arXiv: 2506.18902 [cs.AI]. URL: <https://arxiv.org/abs/2506.18902>.
- [12] Lianmin Zheng et al. *Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena*. 2023. arXiv: 2306.05685 [cs.CL].
- [13] Dawei Zhu et al. *LongEmbed: Extending Embedding Models for Long Context Retrieval*. 2024. arXiv: 2404.12096 [cs.CL]. URL: <https://arxiv.org/abs/2404.12096>.
- [14] Confident AI. *DeepEval: An open-source evaluation framework for LLMs*. <https://github.com/confident-ai/deepeval>. Accessed: 2025-07-01. 2023.
- [15] Yang Liu et al. *G-Eval: NLG Evaluation using GPT-4 with Better Human Alignment*. 2023. arXiv: 2303.16634 [cs.CL].