# AutoRAG: Automated Framework for optimization of Retrieval Augmented Generation Pipeline

Dongkyu Kim*
Markr
jeffrey@markr.ai

Byoungwook Kim*
Markr
bwook@markr.ai

Donggeon Han*
Markr
eastsidegunn@markr.ai

Matouš Eibich
Predli
matous.eibich@datera.cz

October 29, 2024

## Abstract

Using LLMs (Large Language Models) in conjunction with external documents has made RAG (Retrieval-Augmented Generation) an essential technology. Numerous techniques and modules for RAG are being researched, but their performance can vary across different datasets. Finding RAG modules that perform well on specific datasets is challenging.

In this paper, we propose the AutoRAG framework, which automatically identifies suitable RAG modules for a given dataset. AutoRAG explores and approximates the optimal combination of RAG modules for the dataset.

Additionally, we share the results of optimizing a dataset using AutoRAG. All experimental results and data are publicly available and can be accessed through our GitHub repository.

## 1  Introduction

Large Language Models (LLMs) have significantly advanced the field of natural language processing (NLP), enabling applications from text generation to question answering. However, optimizing the integration of dynamic, external information remains challenging. Retrieval Augmented Generation (RAG) techniques address this by incorporating external knowledge sources into the generation process, enhancing the contextual relevance and accuracy of LLM outputs.

While RAG has proven successful, the process of selecting individual RAG techniques is often not automated or optimized, limiting the potential and scalability of this technology. This lack of systematic automation leads to inefficiencies and prevents the comprehensive exploration of RAG configurations, resulting in suboptimal performance.

AutoRAG aims to bridge this gap by introducing an automated framework that systematically evaluates numerous RAG setups across different stages of the pipeline. AutoRAG optimizes the selection of RAG techniques through extensive experimentation, similar to AutoML practices in traditional machine learning. This approach streamlines the evaluation process and improves the performance and scalability of RAG systems, enabling more efficient and effective integration of external knowledge into LLM outputs.

---

*These authors contributed equally to this work.
Special thanks to Shivay Nagpal and Alexander Fred-Ojala for their helpful feedback and thoughtful suggestions.
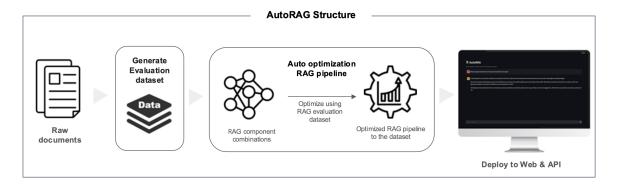
Figure 1: Structural diagram showing the overall structure of AutoRAG.

## 2 RAG Techniques

This section explores various RAG techniques evaluated in our study. We examine strategies for query expansion, retrieval, passage augmentation, passage reranking, and prompt creation. Each technique is aimed at optimizing the integration of external knowledge sources into the generation process to enhance the relevance and accuracy of LLM outputs. See figure 2 to check out all RAG techniques used in this paper.
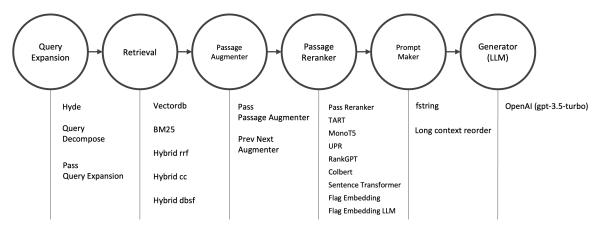


Figure 2: All RAG techniques used in this paper

### 2.1 Query Expansion

It is common to use the user's query directly as a search query in the retrieval system. However, augmenting the query can help enhance retrieval performance. A query expansion module modifies the user's query to create a better search query, making it easier to find the right passage.

#### 2.1.1 Query Decompose

The query decomposition process is used to break down a multi-hop question into single-hop questions using a LLM. This process leverages a default decomposition prompt inspired by the StrategyQA few-shot prompt from (Pereira et al. (2022)).

For instance, consider the multi-hop question: "What is the capital of the country where the inventor of the telephone was born?"

The query decomposition module would break this down into the following single-hop questions:
1. "Who invented the telephone?"
2. "Where was the inventor of the telephone born?"
3. "What is the capital of that country?"
By decomposing the original question into simpler queries, the retrieval system facilitates more accurate passage retrieval.

### 2.1.2   HyDE

The Hypothetical Document Embedding (HyDE) (Gao et al. (2022)) technique improves document retrieval by utilizing LLMs to generate a hypothetical passage to a query. This approach increases the embedding vector similarity that the hypothetical passage will be semantically similar to the actual relevant passage more than the user's query. For instance, consider the following example:

Question: What is Ars-HDGunn structure?

HyDE expanded query: The Ars-HDGunn structure is a specialized architectural design that integrates advanced materials and innovative engineering techniques to create a highly efficient and sustainable building. This structure is characterized by its unique combination of aesthetic appeal and functional performance, often incorporating features such as green roofs, solar panels, and energy-efficient systems.
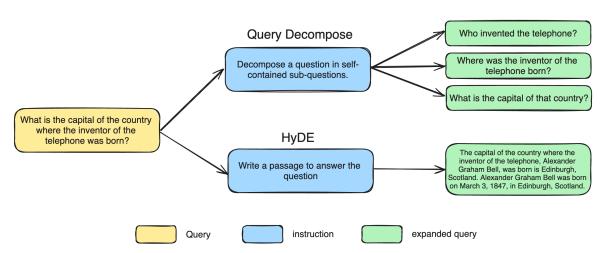


Figure 3: An illustration of Query Decompose and HyDE query expansion modules.

## 2.2   Retrieval

### 2.2.1   Vectordb

For retrieval with vector DB, passage embedding vectors are generated using a pre-trained embedding model. (Karpukhin et al. (2020)) These vectors represent the passages in a high-dimensional space. Subsequently, a query embedding vector is created using the same embedding model. The semantic similarity between the query embedding and each passage embedding is then computed.

The final step involves identifying the passage embedding that has the highest similarity score with the query embedding. This approach, known as cosine similarity search, efficiently retrieves the most relevant passages by leveraging dense vector representations and similarity computations.

### 2.2.2 BM25

BM25 (Best Matching 25) (Robertson and Zaragoza (2009)) is an information retrieval module based on the probabilistic relevance framework. It extends the classic TF-IDF (Term Frequency-Inverse Document Frequency) model by introducing term frequency saturation and document length normalization. BM25 scores are calculated using a set of formulas that account for the term frequency, inverse document frequency, and the length of documents.

In a nutshell, it uses words in queries for searching documents. The more words in a query appear in the document, the more relevance scores it gets. While searching, it uses only infrequent words in the document, ignoring frequent words. Because frequent terms like 'you', 'and', or 'like' can be irrelevant to the query and document meaning.

In this paper, we used 'rank-bm25' library(Brown et al. (2022)) for BM25 calculation.

### 2.2.3 Hybrid Retrievals

Hybrid retrieval is the fusion of sparse retrieval methods like BM25 and dense retrieval methods like vector databases that use embedding models. In this paper, we used four hybrid retrieval variants. Using hybrid retrieval can enhance performance because it leverages both the lexical similarity of sparse retrieval and the semantic similarity of dense retrieval.

**Hybrid RRF** (Cormack et al. (2009)) uses the reciprocal rank fusion algorithm to fuse results. RRF can be represented by the following formula:

$$f_{RRF}(q,d) = \frac{1}{\eta + \pi_{LEX}(q,d)} + \frac{1}{\eta + \pi_{SEM}(q,d)}$$

The $\pi(q,d)$ is the rank function, which gets the rank of the document $d$ to the given query $q$ in the set of documents. In other words, the document $d$ is a subset of the retrieved document set to the query $q$. The $\pi_{LEX}(q,d)$ is the rank of the passage in the lexical retrieval, which is BM25 in this paper. The $\pi_{SEM}(q,d)$ is the rank of the passage in the semantic retrieval, which is 'vectordb' in this paper. The $\eta$ is a free parameter.

**Hybrid CC** (Bruch et al. (2023)) and **Hybrid DBSF** use convex combination to fuse sparse and dense retrieval. Both methods use a convex combination for fusing, but the difference lies in the normalization method. The Hybrid CC retrieval uses min-max normalization, whereas the Hybrid DBSF uses the 3-sigma value as the min and max values in normalization. Both retrieval methods can be represented by the following formula:

$$f_{Convex}(q,d) = \alpha\phi_{LEX}(f_{LEX}(q,d)) + (1-\alpha)\phi_{SEM}(f_{SEM}(q,d))$$

Here, $f_{LEX}(d)$ is the lexical retrieval relevance score of the document, which is the BM25 relevance score in this paper. $f_{SEM}(d)$ is the semantic retrieval relevance score of the document, which is the vectordb relevance score in this paper. The $\phi_{LEX}$ and the $\phi_{SEM}$ are the normalized functions of each retrieval method. The range of the lexical retrieval score and the semantic retrieval score is different, so the normalization process is crucial. The *alpha* is the weight parameter and must be in the range of 0 to 1.

The min-max normalization process for hybrid cc can be represented by the following formula:

$$\phi_{MM}(f_o(q,d)) = \frac{f_o(q,d) - m_o(q)}{M_o(q) - m_o(q)}$$

Here, $m_o(q)$ is the minimum relevance score value of the given retrieval method, and $M_o(q)$ is the maximum relevance score value of the given retrieval method. In other words, the semantic and lexical normalization uses different min and max values in this paper.

The 3-sigma normalization used in hybrid DBSF retrieval can be represented by the following formula:

$$\phi_{ts}(f_o(q,d)) = \frac{f_o(q,d) - (\mu_o - 3\sigma_o)}{(\mu_o + 3\sigma_o) - (\mu_o - 3\sigma_o)}$$

In this formula, the $\mu_o$ is the mean value of the given retrieval method. And the $\sigma_o$ is the standard deviation value of the given retrieval method.

## 2.3 Passage Augmenter

The passage augmenter is a technique designed to enhance the performance of retrieval by obtaining additional relevant passages. This method expands the initial set of retrieved passages, thereby increasing the comprehensiveness and relevance of the information retrieved.

The process begins with an initial retrieval phase where a set of passages is obtained based on a query. Following this, the passage augmenter utilizes the metadata associated with these passages—such as author information, publication date, and keywords to perform a secondary search. In this paper, we only use metadata about neighboring passages. The secondary search aims to find more passages that are contextually related to the initially retrieved set.

### 2.3.1 prev next augmenter

The Prev-Next Passage Augmenter designed to enhance the retrieval performance by incorporating neighboring passages.

During the chunking process, each passage can be annotated with its preceding and succeeding passages. The Prev-Next Passage Augmenter leverages this structural information to retrieve not only the initially relevant passages but also their neighbors. This approach is based on the hypothesis that adjacent passages may contain additional context or relevant information that can improve the overall retrieval performance.

## 2.4 Passage Reranker

The passage reranker is a component in information retrieval systems, tasked with re-ranking passages after the initial retrieval phase. this method, while computationally expensive, has been suggested to enhance accuracy of retrieval modules(Lin (2019)).

In the context of RAG, the passage reranker plays a vital role. RAG are often constrained by the token limit and the high computational cost of LLM. By employing a passage reranker, the system can achieve higher accuracy in identifying the most relevant passages to the query, ensuring efficient use of prompt tokens.

| Passage Reranker | Type |
|---|---|
| MonoT5 | LM-based |
| Sentence Transformer Reranker | |
| Flag Reranker | |
| Flag LLM Reranker | |
| TART | |
| RankGPT | LLM-based |
| Colbert | Embedding-based |
| UPR | Log prob-based |

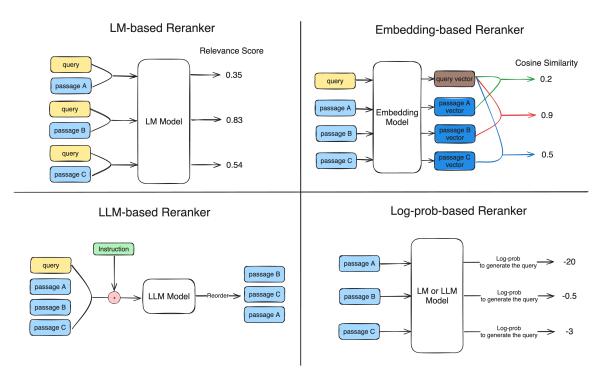Table 1: The passage rerankers and its type we used in this paper

Figure 4: An illustration of each passage reranker module.

### 2.4.1 LM-based Reranker

LM-based rerankers utilize fine-tuned language models to score the relevance of query-passage pairs. The training dataset comprises query-passage pairs annotated with relevance labels. The training process involves:

For relevant query-passage pairs, the model is trained to output the token 'True'. For non-relevant pairs, the model is trained to output the token 'False'. During inference, the model calculates the probability of outputting the 'True' token. This probability is used as the relevance score for the passage.

**MonoT5 Reranker** (Nogueira et al. (2020)) uses the T5 model (Raffel et al. (2023)) and is trained with query-passage pairs labeled for relevance. In this paper, the model is based on the T5-3B variant and is fine-tuned using the MS MARCO dataset (Bajaj et al. (2018)) over 10,000 steps (equivalent to 1 epoch).

The **Sentence Transformer Reranker** employs a cross-encoder model to rerank retrieved passages. In this paper, the ms-marco-MiniLM-L-2-v2 model is utilized for this purpose. This model is fine-tuned to classify query-passage relevance using the MSMARCO dataset (Bajaj et al. (2018)), which is based on the BERT model.

**Flag reranker** uses a BGE model to rerank passages. This model is based on xlm-roberta-base model(Conneau et al. (2019)) and fine-tuned using multilingual datasets like Mr.Tydi dataset (Zhang et al. (2021)).

**Flag llm reranker** is based on a LLM that has many weights compared to other rerankers. In this paper, it uses gemma-2b model (Team and et al. (2024)), fine-tuned for reranker usage.

The **TART reranker** (Asai et al. (2022)) is an LM-based reranker but uses task-specific instructions in the reranking process. TART lies in its ability to include instructions during the reranking phase, thereby capturing the user's intent beyond the explicit query. Other rerankers rely solely on the user's query, which may not fully encapsulate the user's underlying intent. TART addresses this limitation by allowing the inclusion of additional instructions that specify the user's intent.

The major distinction between TART and other rerankers lies in the training methodology. Other LM-based rerankers are trained on query-passage pairs with relevance labels. In contrast, TART trains the reranker model using an instruction-query-passage set. This approach allows TART to understand and incorporate the user's intent, as specified by the instructions, into the reranking process.

### 2.4.2 LLM-based Reranker

LLM-based rerankers leverage LLMs to reorder passages based on their relevance to a given query. Unlike other rerankers that require fine-tuning the language model, LLM-based rerankers utilize prompt engineering to achieve passage reranking.

**RankGPT** (Sun et al. (2023)) inputs the query and the passages to be reranked into the LLM, which then generates permutations of the passages. These generated permutations are the reranked order of the passages.

### 2.4.3 Embedding-based Reranker

Embedding-based rerankers leverage dense vector representations to capture semantic similarities between queries and documents. For reranking passages, they generate embedding vectors for both queries and passages and then calculate the similarity of each vector. Embedding-based rerankers can be an ensemble of different embedding models. For better results, it is common to choose an embedding model different from the one used in the retrieval phase.

**ColBERT reranker** (Khattab and Zaharia (2020)) independently encoding queries and passages with BERT. Then it calculates senmantic simliarity using encoded vectors. ColBERTv2 (Santhanam et al. (2022)) uses lightweight token representations, optimizing computation cost and maintaining rerank effectiveness. We employ ColBERTv2 in this paper.

### 2.4.4 Log prob-based Reranker

Log-probability based rerankers leverage the log probability of generating a query from a given passage to assess relevance. If the log probability of generating a query is higher, the passage is more relevant.

The **UPR** (Unsupervised Passage Reranker) (Sachan et al. (2023)) exemplifies this approach by utilizing a pre-trained language model to compute the log probability of generating a query at the given passage. In this paper, we use T5-large model(Raffel et al. (2020)) for UPR reranker.

## 2.5 Prompt Maker

For in-context learning, the retrieved passages must be included in the prompt to the LLM. Prompt maker is the module that concatenates retrieved passage contents, user queries, and instructions.

### 2.5.1 f-string

This module concatenates the user's query, retrieved passages, and instructions. The higher relevance passage will be the first, and the lowest will be the last. In this paper, the top-k was set as five, so it uses the top five relevant passages.

### 2.5.2 long context reorder

Long context reorder addresses the 'Lost in the Middle' phenomenon (Liu et al. (2023a)), where large language models (LLMs) tend to prioritize the beginning and end of input prompts, often neglecting the middle content.

To mitigate this, the long context reorder module appends the most relevant passage at the end of the input prompt, ensuring it appears both at the beginning and the end. This redundancy helps LLMs maintain focus on critical information, thereby enhancing the performance of generated responses.

# 3  Experiment

## 3.1  AutoRAG

Despite the development of numerous RAG techniques and metrics through research, these advancements have been scattered, making it challenging to identify the appropriate RAG pipeline for real-world applications. To address this issue, we propose AutoRAG, an open-source framework designed for RAG experimentation and optimization. AutoRAG leverages a greedy algorithm to efficiently search for the optimal initial pipeline. By organizing the model into modular nodes, each performing distinct tasks, AutoRAG dynamically selects the most promising node at each step using a strategy defined by metrics available for each node. This approach enables AutoRAG to construct near-optimal pipelines without requiring exhaustive search methods, offering both scalability and computational efficiency across diverse machine learning tasks.

### 3.1.1  Node

A 'node' corresponds to a specific step within RAG and operates as a high-level concept that acts as a container for modules. The modules that can be placed within the same node must have the same input and output formats as the node. Except for the initial node input (User query) and the final node output (Answer), the output of one node is used as the input for the subsequent node. The concepts of nodes and modules used in this experiment are illustrated in Figure 2.

### 3.1.2  Strategy

In AutoRAG, the term "strategy" refers to how we choose which module to use within a node. This involves selecting and arranging different RAG techniques(modules). Each node can use performance metrics and the time a module takes to complete its task. We can also use statistical measures like the average of all these metrics to help make decisions. By defining what makes a RAG module "good" using these performance metrics, we can compare different modules. In AutoRAG, the function that serves as the criteria for selecting and optimizing modules is called the "strategy."

### 3.1.3  Optimization

Evaluating nodes like 'query_expansion' or 'prompt_maker' based solely on their outputs is challenging. In such cases, we utilize the evaluation of the subsequent node. For instance, the output of the 'query_expansion' node is one or more queries for retrieval, which makes it difficult to evaluate the modules. The 'query_expansion' node is positioned before the 'retrieval' node which is relatively easier to evaluate. Therefore, we fix the modules of the 'retrieval' node and only change the modules of the 'query_expansion' node for experiments. Similarly, the 'prompt_maker' node is evaluated by fixing the modules in the 'generation' node.

When the preceding node (A) is difficult to evaluate and the subsequent node (B) that is easier to evaluate, with the number of modules to evaluate being $m$ and $n$ respectively, a comprehensive experiment would require $m \times n$ combinations. However, using the above method, we perform $m$ trials at stage A and $n$ trials at stage B, thereby reducing the number of required experiments to only $m + n$ combinations.

## 3.2  Data

This study utilizes the ARAGOG dataset(Eibich et al. (2024)), a tailored dataset derived from the AI ArXiv collection, accessible via Hugging Face (Briggs (2023)). The dataset consists of 423 selected research papers centered around the themes of AI and LLMs, sourced from arXiv. This selection offers a comprehensive foundation for constructing a database to test the RAG techniques and creating a set of evaluation data to assess their effectiveness.

### 3.2.1 RAG Database Construction

For the study, a subset of 13 research papers were selected for their potential to generate specific, technical questions suitable for evaluating Retrieval-Augmented Generation (RAG) systems. To better simulate a real-world vector database environment, where noise and irrelevant documents are present, the database was expanded to include the full dataset of 423 papers available. The papers were chunked using a chunk size of 512 tokens and an overlap of 50 tokens.

### 3.2.2 Evaluation Data Preparation

The evaluation dataset comprises 107 question-answer (QA) pairs generated with the assistance of GPT-4. The generation process was guided by specific criteria to ensure that the questions were challenging, technically precise, and reflective of potential user inquiries sent to an RAG system. Each QA pair was then reviewed by humans to validate its relevance and accuracy, ensuring that the evaluation data accurately measures the RAG techniques' performance in real-world applications. The QA dataset is available in this paper's associated Github repository.
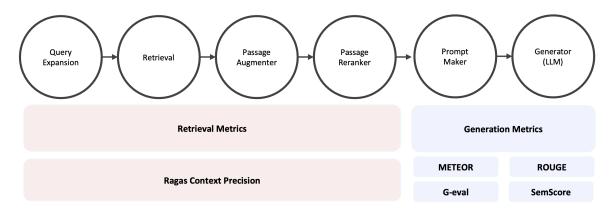
## 3.3 Metrics



Figure 5: Metrics used at each stages

Our experiment employs LLM-based retrieval metrics to find an appropriate retrieval structure at the Retrieval stage for QA problems that are not mapped to readily available knowledge snippets.

### 3.3.1 Retrieval Metric

Ragas (Retrieval Augmented Generation Assessment)(Es et al. (2023)) is a framework designed for reference-free evaluation tools. In our work, we employ the Ragas Context Precision metric.

$$ContextPrecision@K = \frac{\sum_{k=1}^{K}(Precision@k \times v_k)}{truepositives@K}$$

$$Precision@k = \frac{truepositives@k}{(truepositives@k + falsepositives@k)}$$

$K$ is the total number of retrieved passages(prediction class: $positive$). And, $v_k \in \{0,1\}$ is the relevance indicator at rank $k$. K is a parameter that can be set at each retrieval stage. In 3, the parameter 'top_k' related to ragas_context_precision refers to this K value. We utilize GPT-4 turbo to evaluate each retrieved passage and determine its actual class—whether the passage is relevant ($true$) or irrelevant ($false$) to the given query and the target generation output.

### 3.3.2 Generation Metric

One of the challenges in evaluating the generation of outputs by large language models (LLMs) is the lack of a single metric that encompasses all perspectives. Therefore, we use a normalized mean of four metrics to determine the best generation model at each stage.

To intuitively consider the use of strong references and specialized terminology, we employ n-gram based metrics such as ROUGE and METEOR.

To account for semantic similarity with the answer, we adopt the SemScore(Aynetdinov and Akbik (2024)). Using a well-trained embedding model, we compute the cosine similarity in the embedding space between a target text and a generated text. In the experiment, OpenAI's $text - embedding - ada - 002$ was employed as the embedding model for SemScore.

Additionally, we choose G-Eval(Liu et al. (2023b))(GPT4 turbo) to evaluate generation quality comprehensively. G-Eval is an evaluation framework utilizing a Chain-of-Thought(CoT) technique based large language model. This framework employs LLMs to generate scores ranging from 1 to 5. Through prompt engineering, it allows for evaluations from various perspectives. In this experiment, we adopted the perspectives of coherence, consistency, fluency, and relevance, using OpenAI's GPT-4 turbo (gpt-4-0125-preview) model. In this paper, we utilize the average scores for these four aspects as the g_eval score.

## 3.4 Candidate modules

In our study, we focused on implementing and evaluating advanced RAG methods (Gao et al. (2024a)). The experimental setup included several distinct stages. The stages comprised Query Expansion, Retrieval, Passage Augmentation, Passage Reranking, Prompt Creation, and Text Generation. For each stage, we conducted evaluations to select the best-performing module. The output from the best-performing module would be used as the input for the subsequent stage. Below, we detail the modules tested, the metrics used for evaluation, and additional pertinent information for each stage. Note that modules named with the prefix 'pass' indicate that the module produces the same output as its input.

### 3.4.1 Query Expansion

The 'top k' of ragas context precision set to 10. And the modules tested included in the Query Expansion stage:

- `pass_query_expansion` : A module that outputs the same input.

- `query_decompose` (LLM: OpenAI(gpt-3.5-turbo), temperature : [0.2, 1.0])

- `hyde` (LLM: OpenAI(gpt-3.5-turbo), max_token : 64)

### 3.4.2 Retrieval

The 'top k' of ragas context precision set to 10. And the modules tested included in the Retrieval stage:

- `bm25` (tokenizer: GPT-2)

- `vectordb` (OpenAI_embed_3_large)

- `hybrid_rrf` (rrf_k : [3, 5, 10])

- `hybrid_cc` (cc : [0.3, 0.7])

- `hybrid_dbsf` (dbsf : [0.7, 0.3])

### 3.4.3 Passage Augmenter

The 'top k' of ragas context precision set to 15. And the modules tested included in the Retrieval stage:

- `pass_passage_augmenter`
- `prev_next_augmenter` (mode: both)

### 3.4.4 Passage Reranker

The 'top k' of ragas context precision set to 5. And the modules tested included in the Passage Reranking stage:

- `pass_reranker`
- `tart`
- `monot5`
- `upr`
- `rankgpt`
- `colbert_reranker`
- `sentence_transformer_reranker`
- `flag_embedding_reranker`
- `flag_embedding_llm_reranker`

### 3.4.5 Prompt Maker

For the Prompt Maker stage, we employed multiple metrics, including METEOR, ROUGE, and sem_score (OpenAI), alongside 'g_eval' (GPT-4-0125-preview), averaged for comprehensive evaluation. Since a fixed module in the generator node is required to evaluate the 'prompt_maker', we used 'llama_index' implemented with OpenAI's GPT-3.5 Turbo (temperature 0.0) as a fixed generator module. The modules tested were:

- `f_string`
- `long_context_reorder`

### 3.4.6 Generator

The modules tested included in the Generator stage:

- `llama_index_llm` (OpenAI GPT-3.5 Turbo, temperature 0.0)

Each of these stages were thoroughly evaluated to ensure that the best-performing modules were selected based on their respective metrics, ensuring that the subsequent stages received the most effective input possible. The following sections detail the results and analysis of each stage's experiments. Additionally, both METEOR and execution time are provided to offer more detailed evaluation results but are not used for module selection.

# 4 Results

This study systematically evaluates different advanced RAG techniques using Retrieval Metrics and Generation Metrics. A comparative analysis is presented using bar plots to visualize the distribution of these metrics, and the results are interpreted with tables.

## 4.1 Retrieval Metric

To evaluate the performance of our retrieval system, we utilized the Ragas Context Precision as the primary retrieval metric. This metric was applied across a set of 107 queries, yielding 107 individual precision scores for each experimental module. For each module, we computed the mean Ragas Context Precision score from the 107 individual scores. These mean values provide a summary measure of each module's retrieval effectiveness.

To facilitate a comparative analysis of the modules, we visualized the mean Ragas Context Precision scores using bar plots. This visualization allows for a clear comparison of the retrieval performance across the different modules.

### 4.1.1 Query Expansion

The bar plots for Context Precision (Figure 6) indicate varied performance across Query Expansion techniques. For our evaluation, we utilized the Ragas Context Precision metric with a top-k value of 10. Specifically, we calculated the Ragas Context Precision@10 score to assess the retrieval performance. This metric evaluates the precision of the top 10 retrieved passages in terms of their relevance to the given query. We then compared these scores across different retrieval methods to determine their effectiveness. This approach allows us to quantify and compare the retrieval accuracy of various models, providing a robust measure of their performance.
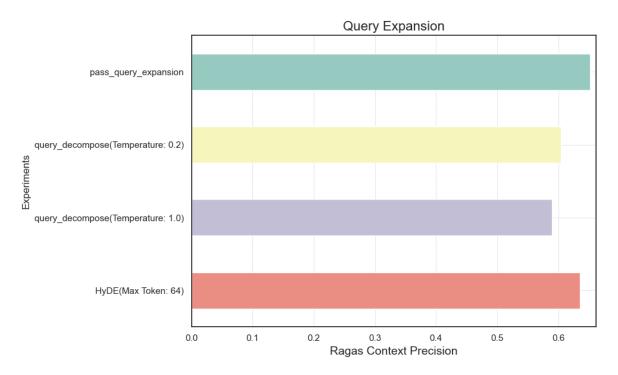


Figure 6: Barplots of Ragas Retrieval Precision by Query Expansion Experiments.

| Query Expansion Modules | Execution Time(seconds) | Ragas Context Precision@10 |
|---|---|---|
| **Pass Query Expansion** | 0.000017 | **0.651694** |
| Query Decompose(Temp: 0.2) | 0.200412 | 0.603911 |
| Query Decompose(Temp: 1.0) | 0.172025 | 0.589451 |
| HyDE(Max Token: 64) | 0.375629 | 0.634954 |

Table 2: Table of Execution Time and Ragas Context Precision by Query Expansion Experiments.

The highest score was achieved by pass query expansion, which uses the base query without Query Expansion. We can see that Query Expansion can improve search performance on certain data, but on other data, it can make retrieval performance worse than the base query, resulting in lower overall RAG performance. Techniques that utilize hypothetical document embedding (HyDE) show lower precision than pass and fail to improve retrieval performance. Decompose performs worse than pass and Hypothetical Document Embedding (HyDE) at both temperatures.

### 4.1.2 Retrieval

The bar plots for Context Precision (Figure 7) indicate varied performance across Retrieval techniques. We compared these scores across different retrieval methods to determine their effectiveness. This approach allows us to quantify and compare the retrieval accuracy of various models, providing a robust measure of their performance.

In our comparative experiment between BM25 and VectorDB, the BM25 algorithm demonstrated superior performance relative to VectorDB. The retrieval performance of VectorDB and BM25 differs across various datasets. Empirical experiments are necessary to determine the better approach, as semantic retrieval methods(VectorDB) sometimes outperform lexical approaches(BM25), and vice versa. In this particular dataset, BM25 demonstrated superior performance compared to traditional VectorDB methods. Consequently, we configured the hybrid retrieval system with weights of 0.7 for BM25 and 0.3 for VectorDB to form hybrid modules. This weighting scheme was chosen to leverage the strengths of BM25 while incorporating the benefits of VectorDB, thereby optimizing the overall retrieval performance. This approach also aimed to reduce the computational cost of the experiment. However, it is important to note that alternative weight configurations can be explored to potentially enhance the performance further. Future experiments could involve varying the hybrid weights to identify the optimal balance between BM25 and VectorDB contributions.
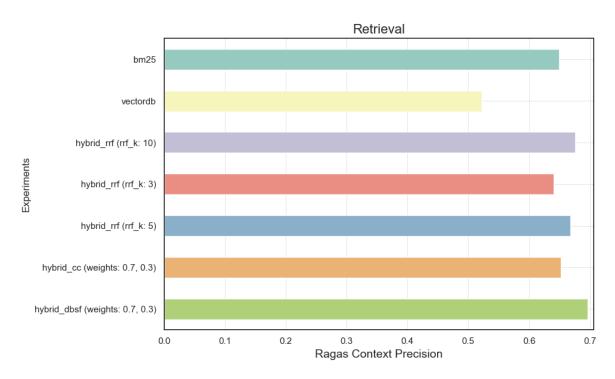
Figure 7: Barplots of Ragas Retrieval Precision by Retrieval Experiments.

| Retrieval Modules | Execution Time(seconds) | Ragas Context Precision@10 |
| --- | --- | --- |
| Bm25 | 0.274728 | 0.649015 |
| VectorDB | 0.496673 | 0.522239 |
| Hybrid RRF (RRF-k: 10) | 0.771401 | 0.676157 |
| Hybrid RRF (RRF-k: 3) | 0.771401 | 0.640295 |
| Hybrid RRF (RRF-k: 5) | 0.771401 | 0.668342 |
| Hybrid CC (Weights: 0.7, 0.3) | 0.771401 | 0.652625 |
| **Hybrid DBSF (Weights: 0.7, 0.3)** | 0.771401 | **0.696401** |

Table 3: Table of Execution Time and Ragas Context Precision by Retrieval Experiments

The highest score was achieved by hybrid DBSF (weights: 0.7, 0.3), which uses the Distribution Based Score Fusion algorithm. The next best performance was observed with the Hybrid Reciprocal Rank Fusion (RRF) method, specifically with an RRF-k value of 10. Our analysis indicates a positive correlation between the RRF-k value and retrieval performance; as the k value increased, so did the performance metrics. At an RRF-k value of 3, the precision was lower than that of the baseline BM25. However, increasing the k value to 10 resulted in the second-highest performance among all tested configurations. Notably, an RRF-k value of 5 alone yielded higher precision values than Hybrid CC. The Hybrid CC algorithm demonstrated superior performance compared to the traditional BM25 algorithm. This suggests that, for this particular dataset, hybrid retrieval methods generally achieve higher precision than conventional retrieval methods such as BM25 and VectorDB (Vector Similarity Search).

### 4.1.3   Passage Augmenter

Passage Augmenter increases the number of Retrieved Passages, so we set top-k to 15. The mode of the prev-next augmenter is set to 'both', thereby incorporating both the previous and next paragraphs. In the previous retrieval step, 10 paragraphs per query are initially obtained. To enhance the context, these paragraphs are augmented by including the previous and next paragraphs, resulting in a total of 30 passages per query. The passage augmenter, configured with a top-k parameter of 15, subsequently selects the 15 highest-scoring passages from these 30 passages. These selected passages are then forwarded to the subsequent processing stage, the passage reranker.
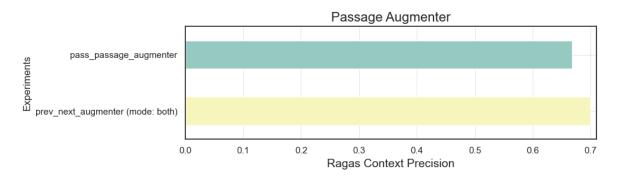


Figure 8: Bar plots of Ragas Retrieval Precision by Passage Augmenter Experiments.

| Passage Augmenter Modules | Execution Time(seconds) | Ragas Context Precision@15 |
| --- | --- | --- |
| Pass Passage Augmenter | 0.003849 | 0.667531 |
| **Prev Next Augmenter** | 0.792790 | **0.699620** |

Table 4: Table of Execution Time and Ragas Context Precision by Passage Augmenter Experiments

In our experiments, the Prev Next Augmenter demonstrated superior performance compared to the Pass Passage Augmenter. This outcome is contingent on the specific corpus data used in our study. The data revealed that the preceding and succeeding passages of the retrieved passage contain valuable contextual information, which enhances the retrieval performance. This finding suggests that incorporating context from adjacent passages is more effective in leveraging contextual information than the Pass Passage Augmenter, thereby improving the overall retrieval accuracy.

### 4.1.4   Passage Reranker

The bar plots for Context Precision (Figure 9) indicate varied performance across Passage Reranker techniques. We compared these scores across different retrieval methods to determine their effectiveness. This approach allows us to quantify and compare the retrieval accuracy of various models, providing a robust measure of their performance.

In the previous passage augmenter step, 15 paragraphs per query are obtained. However, we decided that it would be expensive to include all 15 passages in the prompt. Therefore, as a practical approach, we decided to include only 5 passages in the prompt and set the top-k parameter in the passage reranker to 5. Thus, the passage reranker calculates the scores of 15 passages per query from the previous passage augmenter. It then selects only the five passages with the highest scores and sends them to the next step, the prompt maker.
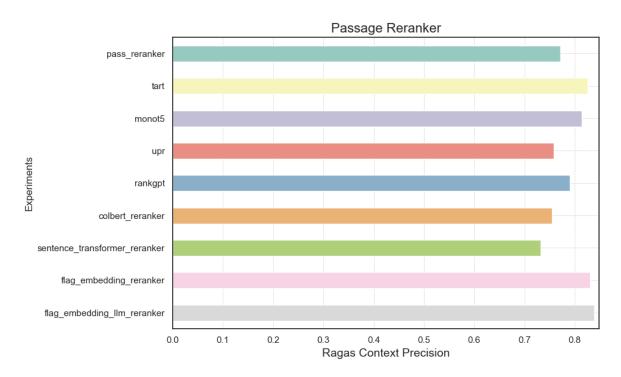
Figure 9: Bar plots of Ragas Retrieval Precision by Passage Reranker Experiments.

| Passage Reranker Modules | Execution Time(seconds) | Ragas Context Precision@5 |
|---|---|---|
| Pass Reranker | 0.000020 | 0.770846 |
| Tart | 0.282748 | 0.826207 |
| Monot5 | 1.473688 | 0.814006 |
| UPR | 0.601418 | 0.758684 |
| RankGPT | 0.219637 | 0.790732 |
| Colbert Reranker | 0.071596 | 0.755244 |
| Sentence Transformer Reranker | 0.020938 | 0.732386 |
| Flag Embedding Reranker | 0.288357 | 0.830218 |
| **Flag Embedding LLM Reranker** | 1.910619 | **0.838253** |

Table 5: Table of Execution Time and Ragas Context Precision by Passage Reranker Experiments

The highest performance in our evaluation was achieved by the LLM Reranker utilizing Flag Embedding. The second highest performing module was also a Flag Embedding Reranker. These results indicate that Flag Embedding rerankers are particularly effective on this dataset. Specifically, the Flag Embedding LLM Reranker, which is based on a LLM, outperformed the Flag Embedding Reranker, which is based on a language model (LM). This suggests that LLM-based rerankers are more effective than LM-based rerankers for this dataset. However, it is noteworthy that LLM-based rerankers are not universally superior. For instance, RankGPT, another LLM-based reranker, ranked fifth out of the nine experimental modules. This demonstrates variability in performance among different LLM-based rerankers.

## 4.2 Generation Metric

To evaluate the answers generated by RAG, we employed four distinct generation metrics. These metrics were applied to each experimental module across 107 queries, resulting in 107 individual accuracy scores per module. The mean value of these 107 scores was then calculated to obtain the Generation Metric Score for each module.

To facilitate a comparative analysis of the modules, we visualized the scores of the four generation metrics using bar plots. Among these metrics, METEOR, ROUGE, and Sem Score are scaled between 0 and 1, whereas G-Eval is scaled between 1 and 5. Due to the differing scales of these metrics, we plotted two separate graphs. This approach ensures a clear and accurate comparison of the performance of RAG's answers across the different experimental modules.

### 4.2.1 Prompt Maker

We utilized identical prompts for the two experimental modules, namely the f-string and the long context reorder modules. By using the same prompts across both modules, we ensured consistency in the experimental conditions, allowing for a more accurate comparison of their performance. The prompts employed in these experiments were as follows:

> **Prompt**
>
> You are an expert Q&A system that is trusted around the world for your factual accuracy. Always answer the query using the provided context information, and not prior knowledge. Ensure your answers are fact-based and accurately reflect the context provided.
> Some rules to follow:
> 1. Never directly reference the given context in your answer.
> 2. Avoid statements like 'Based on the context, ...' or 'The context information ...' or anything along those lines.
> 3. Focus on succinct answers that provide only the facts necessary, do not be verbose. Your answers should be max two sentences, up to 250 characters.
> ————————————
> (context str)
> ————————————
>
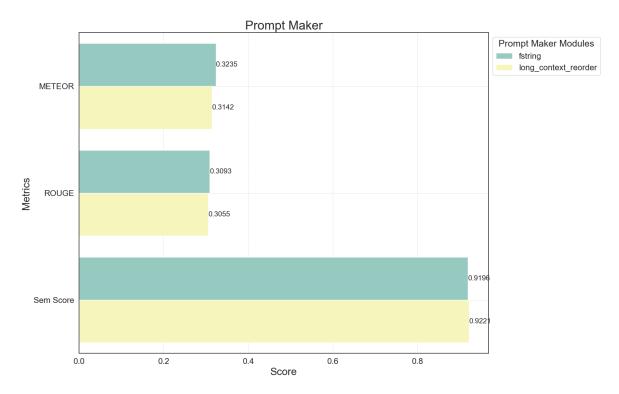> Given the context information and not prior knowledge, answer the query.
> Query: (query str)
> Answer:

Figure 10: Barplots of METEOR, ROUGE and Sem Score by Prompt Maker Experiments.

| Prompt Maker Modules | METEOR | ROUGE | Sem Score |
|---|---|---|---|
| **F-string** | **0.3235** | **0.3093** | 0.9196 |
| Long Context Reorder | 0.3142 | 0.3055 | **0.9221** |

Table 6: Table of METEOR, ROUGE and Sem Score by Prompt Maker Experiments

For the METEOR and ROUGE metrics, the f-string module achieved slightly higher scores compared to the long context reorder module. Similarly, for the n-gram based metric, the f-string module also demonstrated superior performance.

Conversely, the long context reorder module scored marginally better on the Sem Score metric, which is based on semantic similarity. However, the differences in performance between the two modules are minimal, as illustrated in the bar plots (Figure 10).
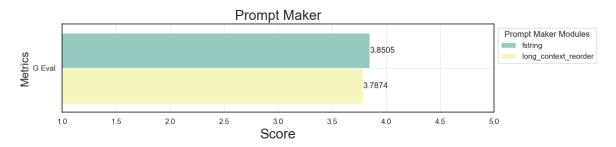


Figure 11: Barplots of G Eval by Prompt Maker Experiments.

| Prompt Maker Modules | G Eval |
| --- | --- |
| **F-string** | **3.8505** |
| Long Context Reorder | 3.7874 |

Table 7: Table of G Eval by Prompt Maker Experiments

In the evaluation using the G-Eval metric, which is judged by an LLM, the f-string module outperformed the long context reorder module.

### 4.2.2 Generator

The LLM model used in our experiments was fixed to GPT-3.5-turbo with a temperature setting of 0.0. As the experiment was not designed to compare the performance between different LLM models, we do not include graphs comparing performance results across various LLM models.

Below are the final results obtained using the selected model and pipeline:

| LLM | METEOR | ROUGE | Sem Score | G Eval |
| --- | --- | --- | --- | --- |
| **gpt-3.5-Turbo** | **0.3246** | **0.3054** | **0.9186** | **3.8037** |

Table 8: Table of Generation metrics by Generator Experiments

In this experiment, both the Prompt Maker and the generator module utilized the same model (GPT-3.5-turbo) and identical temperature settings (0.0). the performance values of the generator module were directly influenced by the outcomes of the Prompt Maker. As a result, Although the settings for both modules were identical, the inherent variability of the LLM resulted in different metric values. While the inherent variability in LLM outputs led to slight variations in the Generation Metric values, these differences were minimal. This indicates that both modules performed consistently under the same conditions, making the performance values of the two modules directly comparable.

## 5   Discussion

The lower performance with the query expansion method can be attributed to the fact that the queries in the configured evaluation dataset were not multi-hop. Hybrid DBSF showed the highest performance among retrieval methods, and the flag embedding llm reranker showed the highest performance among rerankers. The use of the prev-next augmenter as a passage augmenter slightly improved retrieval performance. Also, not using the long context reorder resulted in slightly better performance than using it.

At reranker node, some rerankers, such as UPR, ColBERT, and Sentence Transformer, performed worse than the baseline Pass Reranker without reranking. This indicates that, depending on the dataset, reranking can sometimes degrade performance rather than improve it. The lowest performance was observed with the Sentence Transformer Reranker. The ARAGOG dataset combines multiple papers to form a single query, which may require knowledge from several passages to generate an answer. However, UPR generates queries based on a single passage, likely leading to its lower performance. ColBERT is embedding based re-ranking method. In retrieval stage result, the $ragas\_context\_precision$@10 of BM25(sparse retrieval) is 0.649015 and, the $ragas\_context\_precision$@10 of VectorDB(dense retrieval). The ColBERT method is conceptually similar to Dense retrieval. It may be suggested that the language model has limited semantic understanding of the domain-specific dataset. especially, Sentence Transformer Reranker use $ms - marco - MiniLM - L - 2 - v2$ model. Since it has smaller model parameter size than the models used in other rerankers, semantic understanding issues arising from the domain-specific data may be more pronounced.

# 6 Limitations

- **Normalization Methods**: There are various normalization methods available when executing hybrid cc retrieval, such as TMM. (Bruch et al. (2023)) In this paper, we computed hybrid retrieval using only two different normalization methods. The performance of other normalization methods can vary.

- **Small search space of hybrid retrieval parameter**: In the hybrid retrieval method, both cc and rrf, the hyper-parameters can affect the performance of retrieval. However, in this paper, we evaluated only a few hyper-parameters due to the high cost of retrieval metrics.

- **Expensive to reproduce**: The high computational requirements of RAGAS retrieval metrics make it difficult to reproduce experiments, posing a barrier to the validation and verification of results.

- **Lack of Meta-Evaluation**: There is no meta-evaluation process to quantify how much AutoRAG improves the RAG pipeline performance compared to other methods.

- **Inherent LLM variability**: LLMs are known to produce slightly different outputs even when provided with the same input, due to their stochastic nature. Therefore, the results and interpretations of the pipeline may not be entirely robust.

# 7 Conclusion

In this paper, we use AutoRAG, an automated RAG optimization framework. It identifies the best module for each node using a greedy approach, and the combination of these modules is expected to achieve performance close to the optimal RAG combination.

Using AutoRAG, we were able to automatically optimize the RAG system for the dataset from ARAGOG(Eibich et al. (2024)).

To facilitate future RAG optimization and RAG system evaluation using various datasets and additional RAG techniques, we have released the AutoRAG framework as code on the Github repository. By utilizing AutoRAG, it will be possible to conduct relative evaluations of RAG techniques not covered in this paper, as well as validations using datasets with a larger number of QA pairs from various domains.

# 8 Future Work

- **Evaluation of AutoRAG Optimization Capabilities**: Evaluating the AutoRAG methodology itself is not an easy task. However, it is important to assess the performance of the AutoML framework itself. A performance comparison with methodologies like AutoRAG-HP (Fu et al. (2024)) would also be necessary.

- **Experiments on More Diverse Datasets**: Since automatic testing is possible using AutoRAG, experiments on more datasets should be conducted. This will help understand the characteristics of datasets from various domains and gather information on suitable RAG techniques.

- **Experiments on Additional RAG Modules**: In RAG systems, it is crucial to set appropriate chunking strategies and use suitable parsing techniques for documents. Additionally, techniques like Modular RAG(Gao et al. (2024b)) are being applied. Modular RAG is controlled by multiple components for entire RAG process, then the RAG pipeline can be more flexible and complex than simple linear RAG pipeline. Optimization and experiments can be conducted, including these RAG techniques.

# References

A. Asai, T. Schick, P. Lewis, X. Chen, G. Izacard, S. Riedel, H. Hajishirzi, and W. tau Yih. Task-aware retrieval with instructions, 2022. URL https://arxiv.org/abs/2211.09260.

A. Aynetdinov and A. Akbik. Semscore: Automated evaluation of instruction-tuned llms based on semantic textual similarity, 2024. URL https://arxiv.org/abs/2401.17072.

P. Bajaj, D. Campos, N. Craswell, L. Deng, J. Gao, X. Liu, R. Majumder, A. McNamara, B. Mitra, T. Nguyen, M. Rosenberg, X. Song, A. Stoica, S. Tiwary, and T. Wang. Ms marco: A human generated machine reading comprehension dataset, 2018. URL https://arxiv.org/abs/1611.09268.

J. Briggs. Hugging face, 2023. URL https://huggingface.co/datasets/jamescalam/ai-arxiv. Oct, 10, 2023.

D. Brown, S. Jain, V. Novotný, and nlp4whp. "rank bm25", 2022. URL https://doi.org/10.5281/zenodo.6106156.

S. Bruch, S. Gai, and A. Ingber. An analysis of fusion functions for hybrid retrieval. *ACM Trans. Inf. Syst.*, 42(1), aug 2023. ISSN 1046-8188. doi: 10.1145/3596512. URL https://doi.org/10.1145/3596512.

A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, and V. Stoyanov. Unsupervised cross-lingual representation learning at scale. *CoRR*, abs/1911.02116, 2019. URL http://arxiv.org/abs/1911.02116.

G. V. Cormack, C. L. A. Clarke, and S. Büttcher. Reciprocal rank fusion outperforms condorcet and individual rank learning methods. *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, 2009. URL https://api.semanticscholar.org/CorpusID:12408211.

M. Eibich, S. Nagpal, and A. Fred-Ojala. Aragog: Advanced rag output grading, 2024. URL https://arxiv.org/abs/2404.01037.

S. Es, J. James, L. Espinosa-Anke, and S. Schockaert. Ragas: Automated evaluation of retrieval augmented generation, 2023. URL https://arxiv.org/abs/2309.15217.

J. Fu, X. Qin, F. Yang, L. Wang, J. Zhang, Q. Lin, Y. Chen, D. Zhang, S. Rajmohan, and Q. Zhang. Autorag-hp: Automatic online hyper-parameter tuning for retrieval-augmented generation, 2024. URL https://arxiv.org/abs/2406.19251.

L. Gao, X. Ma, J. Lin, and J. Callan. Precise zero-shot dense retrieval without relevance labels, 2022. URL https://arxiv.org/abs/2212.10496.

Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, M. Wang, and H. Wang. Retrieval-augmented generation for large language models: A survey, 2024a. URL https://arxiv.org/abs/2312.10997.

Y. Gao, Y. Xiong, M. Wang, and H. Wang. Modular rag: Transforming rag systems into lego-like reconfigurable frameworks, 2024b. URL https://arxiv.org/abs/2407.21059.

V. Karpukhin, B. Oğuz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, and W. tau Yih. Dense passage retrieval for open-domain question answering, 2020. URL https://arxiv.org/abs/2004.04906.

O. Khattab and M. Zaharia. Colbert: Efficient and effective passage search via contextualized late interaction over bert, 2020. URL https://arxiv.org/abs/2004.12832.

J. Lin. The neural hype and comparisons against weak baselines, 2019. URL https://paperswithcode.com/paper/the-neural-hype-and-comparisons-against-weak.

N. F. Liu, K. Lin, J. Hewitt, A. Paranjape, M. Bevilacqua, F. Petroni, and P. Liang. Lost in the middle: How language models use long contexts, 2023a. URL https://arxiv.org/abs/2307.03172.

Y. Liu, D. Iter, Y. Xu, S. Wang, R. Xu, and C. Zhu. G-eval: Nlg evaluation using gpt-4 with better human alignment, 2023b. URL https://arxiv.org/abs/2303.16634.

R. Nogueira, Z. Jiang, and J. Lin. Document ranking with a pretrained sequence-to-sequence model, 2020. URL https://arxiv.org/abs/2003.06713.

J. Pereira, R. Fidalgo, R. Lotufo, and R. Nogueira. Visconde: Multi-document qa with gpt-3 and neural reranking, 2022. URL https://arxiv.org/abs/2212.09656.

C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL http://jmlr.org/papers/v21/20-074.html.

C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer, 2023. URL https://arxiv.org/abs/1910.10683.

S. Robertson and H. Zaragoza. The probabilistic relevance framework: Bm25 and beyond. *Found. Trends Inf. Retr.*, 3(4):333–389, apr 2009. ISSN 1554-0669. doi: 10.1561/1500000019. URL https://doi.org/10.1561/1500000019.

D. S. Sachan, M. Lewis, M. Joshi, A. Aghajanyan, W. tau Yih, J. Pineau, and L. Zettlemoyer. Improving passage retrieval with zero-shot question generation, 2023.

K. Santhanam, O. Khattab, J. Saad-Falcon, C. Potts, and M. Zaharia. Colbertv2: Effective and efficient retrieval via lightweight late interaction, 2022.

W. Sun, L. Yan, X. Ma, S. Wang, P. Ren, Z. Chen, D. Yin, and Z. Ren. Is chatgpt good at search? investigating large language models as re-ranking agents, 2023. URL https://arxiv.org/abs/2304.09542.

G. Team and T. et al. Gemma: Open models based on gemini research and technology, 2024. URL https://arxiv.org/abs/2403.08295.

X. Zhang, X. Ma, P. Shi, and J. Lin. Mr. tydi: A multi-lingual benchmark for dense retrieval, 2021. URL https://arxiv.org/abs/2108.08787.