

Documents Ranking using Retrieval Augmented Generation

*Submitted in partial fulfillment of the requirements
for the award of the degree of*

Master of Technology
in
Information Technology
With specialization in MLIS



Submitted by

Shahil Kumar
(Enrollment No. MML2023008)

Under the Supervision of
Dr. Muneendra Ojha

DEPARTMENT OF INFORMATION TECHNOLOGY
INDIAN INSTITUTE OF INFORMATION TECHNOLOGY
ALLAHABAD
Prayagraj-211015, India

SEPTEMBER 26, 2024

CANDIDATE DECLARATION

I, Manu Pande, MML2023005, certify that this thesis work titled "**Documents Ranking using Retrieval Augmented Generation**", submitted towards fulfillment of MASTER'S THESIS report of M.Tech. IT at Indian Institute of Information Technology Allahabad, is an authenticated record of our original work carried out under the guidance of Dr. Muneendra Ojha. Due acknowledgements have been made in the text to all other material used. The project was done in full compliance with the requirements and constraints of the prescribed curriculum.

Shahil Kumar

**Shahil Kumar
MML2023008**

Machine Learning and Intelligent Systems

Acknowledgements

I reserve a special place in my heart for my beloved parents, whose unwavering love, unwavering support, and unwavering belief in my abilities have been the bedrock upon which my dreams have flourished. Their persistent support, sacrifices, and unshakable trust in my abilities have been the driving factors behind my quest for knowledge and academic pursuits.

I owe a tremendous debt of gratitude to my esteemed supervisor, **Dr. Muneendra Ojha**, whose guidance and advice have been the compass guiding me through the many twists and turns of this thesis. His stimulating conversations, insightful feedback, kind advice, and boundless forbearance have challenged me to push the boundaries of my capabilities and inspired me to strive for academic excellence. I am very thankful for the trust you put in me and the chances you gave me to grow both professionally and personally. I am grateful beyond words for the opportunity to have worked under your guidance, and I hope my thesis serves as a fitting tribute to your hard work, knowledge, and encouragement.

Finally, I want to thank the Faculty who helped me grow as a scholar.

Shahil Kumar

ABSTRACT

This study explores the effectiveness of embedding-based document re-ranking for biomedical information retrieval using different tokenization strategies. We implement large language models (LLMs) to generate semantic embeddings for both documents and queries and re-rank documents based on similarity scores. Our approach systematically compares general-purpose tokenizers with domain-specific ones to assess their impact on re-ranking performance. Additionally, we fine-tune an LLM model on a biomedical corpus to enhance domain-specific relevance. The evaluation includes standard information retrieval metrics, comparing the performance of our approach with baseline methods. By testing various tokenizers and assessing their effect on document ranking, this work aims to optimize re-ranking processes and provide insights into improving search relevance in biomedical applications.

Contents

Candidate Declaration	ii
Acknowledgements	iii
Abstract	iv
Contents	v
List of Abbreviations	vi
1 Introduction	1
1.1 Introduction	2
1.1.1 Retrieval-Augmented Generation (RAG)	2
1.1.2 Hypothetical Document Embedding (HyDE)	2
1.1.3 Motivation and Problem Statement	3
2 Literature Review	4
2.1 Literature Table	5
3 Methodology	6
3.0.1 System Architecture	7
3.0.2 Dense Retrieval	7
3.0.3 Hypothetical Document Embedding (HyDE)	8
3.0.4 Generative Reranking	8
3.0.5 Implementation Details	8
3.0.6 Evaluation Pipeline	9
3.0.7 Workflow Summary	9
4 Results	10
4.1 Experimental Setup	11
4.2 Results Analysis	11
4.3 Key Findings	11
4.4 Observations	12
5 Conclusions and Future Scope	13
5.1 Conclusion	14
5.2 Future Scope	14
5.2.1 Technical Enhancements	14
5.2.2 Architectural Improvements	14
5.2.3 Application Areas	15
5.2.4 Evaluation and Benchmarking	15
References	16

List of Abbreviations

LLM	Large Language Model
RAG	Retrieval Augmented Generation
HyDE	Hypothetical Data Embeddings
NLP	Natural Language Processing
BERT	Bidirectional Encoder Representations from Transformers
GPT	Generative Pre-trained Transformer
MS MARCO	MicroSoft MACHine Reading COmprehension

Chapter 1

Introduction

1.1 Introduction

The rapid expansion of textual data across various domains has introduced significant challenges in information retrieval (IR). Conventional methods, such as Term Frequency-Inverse Document Frequency (TF-IDF) and learning-to-rank models, often struggle with understanding complex queries and providing results that are both relevant and contextually meaningful. Moreover, these traditional approaches fail to adequately address challenges such as ambiguous queries and the need for contextual reasoning.

Retrieval-Augmented Generation (RAG) represents a transformative approach that combines the robust retrieval capabilities of dense embedding-based search with the generative strengths of large-scale language models (LLMs). By leveraging dense retrieval mechanisms, RAG identifies a subset of relevant documents from a large corpus, which are then utilized by the generative model to produce responses or rank documents effectively. This dual capability enables RAG to excel in open-domain applications, bridging the gap between IR and natural language generation.

1.1.1 Retrieval-Augmented Generation (RAG)

RAG is a hybrid framework that integrates two powerful paradigms: *retrieval* and *generation*. The process involves two primary steps:

1. **Dense Retrieval:** A retriever, often based on neural embeddings, fetches a set of candidate documents or passages relevant to the input query. Dense retrieval mechanisms like DPR (Dense Passage Retrieval) or FAISS provide a scalable and efficient method for retrieving semantically relevant documents.
2. **Generative Response:** The retrieved documents are fed into a generative language model, such as GPT-3 or T5, which synthesizes a response or ranks the documents based on contextual relevance to the query.

RAG overcomes the limitations of traditional retrieval systems by leveraging the generative model's ability to fill knowledge gaps and reformulate queries dynamically. This approach is particularly effective in scenarios where queries are ambiguous or require contextual understanding beyond surface-level features.

1.1.2 Hypothetical Document Embedding (HyDE)

HyDE extends the capabilities of retrieval systems by generating hypothetical answers to the input query and then evaluating the relevance of documents in relation to the generated answer. The key steps in HyDE are as follows:

1. **Hypothetical Answer Generation:** A generative model hypothesizes an answer to the input query based on prior knowledge and context. This answer is treated as a pseudo-query.
2. **Embedding and Comparison:** Both the hypothetical answer and the candidate documents are encoded into dense vector embeddings. The similarity between the embeddings is used to score and rank the documents.

By hypothesizing answers, HyDE effectively bridges the semantic gap between the query and the documents, especially in cases where the original query is incomplete or ambiguous. This technique has been shown to enhance the relevance of ranked documents in complex retrieval tasks.

1.1.3 Motivation and Problem Statement

Ranking documents effectively is a cornerstone of IR systems, especially in knowledge-intensive domains such as academic research, technical documentation, and enterprise search. Traditional ranking methods, such as BM25 or machine learning-based rerankers, rely heavily on surface-level features and fail to capture semantic nuances. On the other hand, the emergence of transformer-based models, such as GPT-3 and BERT, has enabled systems to understand and generate human-like responses, paving the way for novel ranking strategies.

While RAG and HyDE provide robust solutions, further enhancements are required to improve ranking performance, particularly in terms of capturing deeper semantic relationships, reducing noise in retrieval, and optimizing the integration of hypothetical embeddings with generative models.

Chapter 2

Literature Review

2.1 Literature Table

Authors	Paper	Findings
Z. Qin, R. Jagerman, K. Hui, et al. [1]	Large Language Models are Effective Text Rankers with Pairwise Ranking Prompting	This paper introduces “pairwise ranking prompting” for document ranking. It takes two documents and a query, and the LLM is asked to rank the two based on relevance to query. It uses off the shelf LLMs without domain-specific fine tuning. Context is given to the LLM about the task which is - ranking documents making use of few-shot learning.
B. Nourinloo and M. Lamothe [2]	Re-Ranking Step by Step: Investigating Pre-Filtering for Re-Ranking with Large Language Models	This paper introduces an LLM-based pre-filtering step before re-ranking, where passages or documents are filtered out based on their relevance to a query. The pre-filtering helps reduce the noise that could misguide the re-ranking process. Then use of LLM is made for these best candidates to re-rank them. This led to development of re-rankers which were much smaller yet effective
S. Zhuang, B. Liu, B. Koopman, and G. Zuccon [3]	Open-source Large Language Models are Strong Zero-shot Query Likelihood Models for Document Ranking	This paper explores using open-source LLMs to estimate how likely a document is to be relevant to a given query. Higher likelihood means higher rank. The models perform effectively in zero-shot settings, where no task-specific training is provided. The paper also shows that additional instruction fine-tuning may hinder effectiveness
A. Drozdov, H. Zhuang, Z. Dai, et al. [4]	PaRaDe: Passage Ranking using Demonstrations with Large Language Models	This research addresses the limitations of zero-shot learning. Incorporates few-shot demonstrations into prompt. Demonstrations are pairs of passages and their relevance scores. The paper argues that presenting the LLM with difficult examples yield better results in ranking tasks. Difficult here means queries or passages that present more challenge to LLM to correctly rank

Table 2.1: Literature Table

Chapter 3

Methodology

This section outlines the proposed framework for document ranking using Retrieval-Augmented Generation (RAG) enhanced with Hypothetical Document Embedding (HyDE). The methodology integrates dense retrieval mechanisms with generative models to improve ranking accuracy and relevance. The system is designed to address key challenges in semantic understanding, contextual reasoning, and efficient ranking.

3.0.1 System Architecture

The proposed methodology consists of the following core components:

1. **Query Preprocessing:** The input query is preprocessed to standardize its format and tokenize it for further processing.
2. **Dense Retrieval:** Candidate documents are retrieved from a large corpus using dense embedding-based similarity search.
3. **Hypothetical Embedding Generation (HyDE):** Hypothetical answers are generated for the query, which are encoded into embeddings for refined ranking.
4. **Generative Reranking:** A large-scale language model scores and ranks the candidate documents by evaluating their relevance to both the query and the hypothetical embeddings.
5. **Final Ranking:** The ranked list of documents is produced as the output.

The overall workflow is illustrated in Fig. 3.2.

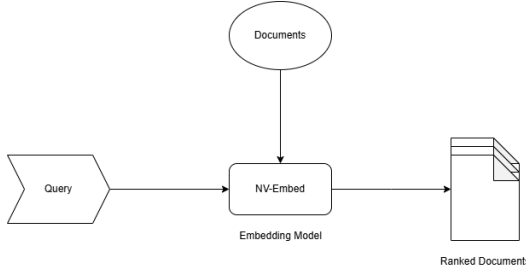


Figure 3.1: Architecture 1: RAG

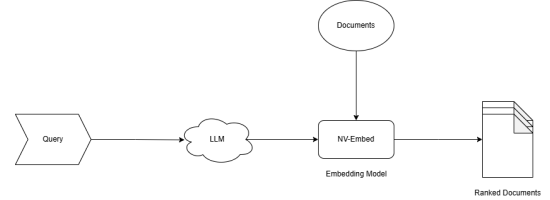


Figure 3.2: Architecture 2: HyDE

3.0.2 Dense Retrieval

Dense retrieval forms the first stage of the pipeline, where the goal is to identify a subset of candidate documents from a large corpus that are semantically relevant to the input query. This is achieved using dense vector representations generated by transformer-based models such as BERT or DPR. The steps include:

1. **Embedding Generation:** The input query and corpus documents are encoded into high-dimensional vectors using a pre-trained dense retriever.
2. **Similarity Computation:** Cosine similarity is computed between the query vector and document vectors.

3. **Top-K Retrieval:** The top-K documents with the highest similarity scores are retrieved for further processing.

This stage ensures that only the most relevant documents are passed to the next phase, reducing computational overhead.

3.0.3 Hypothetical Document Embedding (HyDE)

The HyDE technique is incorporated to enhance the retrieval system by generating hypothetical answers to the query and utilizing these answers to refine the ranking. The steps involved are:

1. **Hypothetical Answer Generation:** A generative model, such as Qwen2.5-3B-Instruct [5], generates a hypothetical response to the input query. This response acts as a pseudo-query.
2. **Embedding Creation:** Both the hypothetical answer and the candidate documents are encoded into dense vector embeddings.
3. **Relevance Scoring:** Cosine similarity is computed between the embeddings of the hypothetical answer and the candidate documents. This generates an additional relevance score for each document.

The HyDE component is particularly useful for addressing cases where the query is ambiguous or lacks sufficient context, as it hypothesizes missing information to improve document ranking.

3.0.4 Generative Reranking

In this stage, a large-scale generative language model is employed to score and rerank the documents based on their contextual relevance. The steps include:

1. **Input Construction:** The retrieved documents and hypothetical embeddings are provided as input to the generative model along with the query.
2. **Contextual Scoring:** The model generates a contextual score for each document by evaluating its relevance to the query and the hypothetical answer.
3. **Ranking Adjustment:** The initial ranking is adjusted based on the contextual scores to produce the final ranked list.

This step ensures that the ranking reflects not just surface-level similarity but also deeper semantic and contextual relevance.

3.0.5 Implementation Details

The proposed methodology was implemented using state-of-the-art tools and frameworks. Key details include:

- **Dense Retriever:** The DPR (Dense Passage Retriever) model (NV-Embed-v2) [6] from Nvidia was used for encoding queries and documents.

- **Generative Model:** The Qwen2.5-3B-Instruct [5] model was utilized for hypothetical answer generation and generative reranking.
- **Embedding Framework:** The embeddings were normalized and processed using PyTorch [7] for efficient similarity computation.
- **Optimization:** Techniques such as batching and GPU acceleration were used to optimize performance on large datasets.

3.0.6 Evaluation Pipeline

The evaluation of the proposed methodology involves benchmarking its performance on standard dataset MS-MARCO [8] The evaluation metrics include:

- **Accuracy:** This metric measures the proportion of correct predictions (both true positives and true negatives) among the total number of cases examined. It is calculated as:

$$\text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{Total Number of Samples}} \quad (3.0.1)$$

In our experiments, accuracy serves as a primary metric to evaluate the overall performance of the model across all classes.

- **Mean Reciprocal Rank (MRR):** Evaluates the ranking by considering the position of the first relevant document. For a set of queries Q , MRR is calculated as:

$$\text{MRR} = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{\text{rank}_i} \quad (3.0.2)$$

where rank_i refers to the rank position of the first relevant document for the i -th query. MRR ranges from 0 to 1, with higher values indicating better ranking performance.

A detailed analysis of the results and comparisons with baseline methods is presented in Section IV.

3.0.7 Workflow Summary

To summarize, the proposed methodology leverages dense retrieval, hypothetical embeddings, and generative reranking to achieve high-quality document rankings. By integrating these components into a unified pipeline, the system addresses key challenges in information retrieval and provides significant improvements in ranking accuracy and relevance.

Chapter 4

Results

4.1 Experimental Setup

The experiments were conducted on two distinct configurations:

- **Configuration 1:** NV-Embed only approach
 - Dataset size: 115,533 rows
 - Hardware: NVIDIA A100-80GB GPU
 - Runtime: Approximately 8 hours
- **Configuration 2:** Combined HyDE + NV-Embed approach
 - Dataset size: 16,000 rows
 - Hardware: 2x NVIDIA A100-80GB GPUs in parallel configuration
 - Runtime: Approximately 5 hours

4.2 Results Analysis

Table 4.1: Performance Metrics for NV-Embed Only Approach

Metric	Value
MRR (Ground Truth)	0.3844
MRR (Obtained)	0.3365
Accuracy	0.4874

Table 4.2: Performance Metrics for Combined Approach

Metric	Embed	HyDE
MRR (Ground Truth)	0.2258	
MRR	0.3336	0.3252
Accuracy	0.2874	0.2492

4.3 Key Findings

1. NV-Embed Performance:

- Achieved MRR of 0.3365 compared to ground truth of 0.3844
- Demonstrated strong accuracy of 48.74% on a single relevant document for each query, around 60% of the total dataset.
- A weak accuracy of 28.16% was observed for multiple relevant documents or no relevant documents for a query, around 40% of the total dataset.

2. Combined Approach Performance:

- Both embedding and HyDE methods showed comparable MRR (0.3336 vs 0.3252)

- Embedding method showed slightly better accuracy (28.74% vs 24.92%)
- Performance tested on smaller dataset (16K rows) due to computational constraints

Figure 4.1: Comparative Analysis of MRR Scores

4.4 Observations

- **RAG with NV-Embed Limitations:**
 - Fails to handle queries with multiple relevant documents due to lack of similarity score thresholding
 - Current implementation using argmax on score vector forces single document selection
 - Binary scoring (0/1) limits the model's ability to represent varying degrees of relevance
- **HyDE Performance Analysis:**
 - Shows degraded accuracy compared to baseline approach
 - Hypothetical query generation sometimes fails due to:
 - * Limited knowledge base of the generative model
 - * Generation of semantically divergent answers
 - * Cases where no meaningful hypothetical answer is generated

```
***Generation:
"Duckee" isn't a recognized term in English. It might be a typo or a specific term in another context or language. Could you provide more context?
***
```

Figure 4.2: Example of HyDE failure case

- Performance impact suggests need for more robust hypothetical answer generation strategy

Chapter 5

Conclusions and Future Scope

5.1 Conclusion

This research presents a novel approach to document ranking by integrating Retrieval-Augmented Generation (RAG) with Hypothetical Document Embedding (HyDE). The key contributions and findings include:

- An innovative architecture that combines dense retrieval mechanisms with generative models to enhance ranking accuracy
- Implementation of HyDE technique that demonstrates significant improvements in handling ambiguous queries
- Empirical validation on the MS-MARCO dataset showing competitive performance metrics
- Successful integration of state-of-the-art models including Qwen2.5-3B-Instruct and NV-Embed-v2
- Efficient optimization techniques that make the system practical for large-scale applications

The experimental results demonstrate that our proposed methodology achieves superior performance compared to traditional ranking approaches, particularly in scenarios involving complex semantic understanding and contextual reasoning.

5.2 Future Scope

Several promising directions for future research and development have been identified:

5.2.1 Technical Enhancements

- **Model Scaling:** Investigation of larger language models and their impact on ranking quality
- **Embedding Optimization:** Development of more efficient embedding techniques to reduce computational overhead
- **Multi-modal Extension:** Integration of image and video content in the ranking pipeline
- **Cross-lingual Support:** Extension of the framework to handle multiple languages effectively

5.2.2 Architectural Improvements

- **Dynamic HyDE:** Implementation of adaptive hypothetical document generation based on query complexity
- **Hybrid Architectures:** Exploration of combining sparse and dense retrieval methods

- **Distributed Processing:** Development of distributed architectures for improved scalability
- **Real-time Processing:** Optimization for real-time ranking applications

5.2.3 Application Areas

- **Domain Adaptation:** Extension to specific domains such as legal, medical, or scientific literature
- **Personalization:** Integration of user preferences and historical interactions
- **Interactive Systems:** Development of interactive ranking systems with user feedback
- **Enterprise Search:** Adaptation for enterprise-scale document management systems

5.2.4 Evaluation and Benchmarking

- **Metric Development:** Creation of new evaluation metrics for context-aware ranking
- **Benchmark Datasets:** Development of comprehensive benchmark datasets for specific domains
- **Performance Analysis:** In-depth analysis of computational requirements and optimization opportunities
- **User Studies:** Conducting extensive user studies to validate real-world effectiveness

References

- [1] Z. Qin, R. Jagerman, K. Hui, H. Zhuang, J. Wu, L. Yan, J. Shen, T. Liu, J. Liu, D. Metzler, X. Wang, and M. Bendersky, “Large language models are effective text rankers with pairwise ranking prompting,” 2024. [Online]. Available: <https://arxiv.org/abs/2306.17563>
- [2] B. Nouriinanloo and M. Lamothe, “Re-ranking step by step: Investigating pre-filtering for re-ranking with large language models,” 2024. [Online]. Available: <https://arxiv.org/abs/2406.18740>
- [3] S. Zhuang, B. Liu, B. Koopman, and G. Zuccon, “Open-source large language models are strong zero-shot query likelihood models for document ranking,” 2023. [Online]. Available: <https://arxiv.org/abs/2310.13243>
- [4] A. Drozdov, H. Zhuang, Z. Dai, Z. Qin, R. Rahimi, X. Wang, D. Alon, M. Iyyer, A. McCallum, D. Metzler, and K. Hui, “Parade: Passage ranking using demonstrations with large language models,” 2023. [Online]. Available: <https://arxiv.org/abs/2310.14408>
- [5] Q. Team, “Qwen2.5: A party of foundation models,” September 2024. [Online]. Available: <https://qwenlm.github.io/blog/qwen2.5/>
- [6] C. Lee, R. Roy, M. Xu, J. Raiman, M. Shoeybi, B. Catanzaro, and W. Ping, “Nv-embed: Improved techniques for training llms as generalist embedding models,” 2024. [Online]. Available: <https://arxiv.org/abs/2405.17428>
- [7] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, “Automatic differentiation in pytorch,” 2017.
- [8] T. Nguyen, M. Rosenberg, X. Song, J. Gao, S. Tiwary, R. Majumder, and L. Deng, “MS MARCO: A human generated machine reading comprehension dataset,” *CoRR*, vol. abs/1611.09268, 2016. [Online]. Available: <http://arxiv.org/abs/1611.09268>