

Integrating Retrieval-Augmented Generation (RAG) in Large Language Models (LLMs): An Approach for Faster and Accurate Search

Dr. Ashok Kumar¹, Pankaj Kumar², Amit Kumar³ and Varun Kumar⁴

¹Assistant Professor, Faculty of Library and Information Science, School of Social Sciences, Indira Gandhi National Open University, New Delhi. ashokkr@ignou.ac.in

²Assistant Professor, Department of Library and Information Science, School of Arts, Communication and Languages, Hemvati Nandan Bahuguna Garhwal University, Uttarakhand. informpankaj1994@gmail.com

³Library and Information Assistant, Visvesvaraya National Institute of Technology, Nagpur. amitkumar@vnit.ac.in

⁴Research Scholar, Faculty of Library and Information Science, School of Social Sciences, Indira Gandhi National Open University, New Delhi. 225034709@ignou.ac.in

How to cite this article: Ashok Kumar, Pankaj Kumar, Amit Kumar, Varun Kumar (2024). Integrating Retrieval-Augmented Generation (RAG) in Large Language Models (LLMs): An Approach for Faster and Accurate Search. *Library Progress International*, 44(2), 1070-1078.

ABSTRACT

The integration of Large Language Models (LLMs) with Retrieval-Augmented Generation (RAG) has been one of the largest breakthroughs that information retrieval has seen in the last few years. The paper discussed developments from the earliest keyword-based systems to today's neural models and focused on the impact LLMs have brought to precision and efficiency in search. LLMs, such as (Generative Pre-trained Transformer) GPT-3 and (Bidirectional Encoder Representations from Transformers) BERT, have improved the ability to understand and generate contextually relevant responses, surpassing previous IR methods. RAG models represent an advanced kind of retrieval, combining the understanding of a context by an LLM with external knowledge bases that could possibly retrieve information faster and more accurately. Challenges will arise, such as computational demands, bias, and factual accuracy; ongoing research will pursue optimization. This paper discusses all core methodologies and strategies and provides a perspective on the potential and limitations of using RAG with LLMs to revolutionize the information retrieval landscape.

KEYWORDS

Large Language Models (LLMs), Information Retrieval (IR), Retrieval-Augmented Generation (RAG), (Generative Pre-trained Transformer) GPT and (Bidirectional Encoder Representations from Transformers) BERT.

1. Introduction

The evolution of Information Retrieval Systems has undergone a transformative shift over the past few decades, driven by advances in computational linguistics and artificial intelligence. Before the large language models, IR primarily incorporated term-based approaches that mostly amounted to matching keywords from query user interfaces to the documents. However, neural models have allowed a much deeper understanding of context and semantics that has changed the game in IR (Zhu et al., 2024). Keyword-based systems are fundamental but limited by dependence on exact matches and failure to develop a rich understanding of context and semantics. This meant incomplete or irrelevant search results, hence demanding other approaches for efficiency and accuracy improvement in searching. The new developments are considered quantum computing, which is said to improve the accuracy of searching with algorithms such as Shor's and Grover's, two of the most popular quantum algorithms that demonstrate the positive impacts of quantum computing compared to classical computing, which handles big data much better (Mehta et al., 2024). In the earlier times of information retrieval, it uses Regular Expressions (RegX), which are strings of characters that define the search pattern. Its usefulness in strings, data traction is because of the ease of complicated query patterns and flexibility in text processing (Witten, Moffat,

& Bell, 1999). However, these systems were only limited in the semantic meanings they could not understand, thus always limiting the overall effectiveness of divergent and ambiguous queries handling. With machine learning and NLP, the sectors were finally able to take one big leap forward. Earlier, the original aim of models to be related to machine learning was using the techniques of supervised learning to help in retrieving text through statistical methods applied on ranking and retrieving information (Manning, Raghavan, & Schütze, 2008). The actual revolutionaries however were large language models.

For example, OpenAI's GPT-3 boasts the power of capability for text-like human and its followers, which finds its basis on deep architectures for understanding and generating human-like text from huge percolations of text data (Brown et al. 2020). In fact, LLMs differ from the earlier systems in that they can understand context and infer meaning and more relevant and coherent responses. The impacts of LLMs on information retrieval have already been profound. Combining contextual understanding and generation capabilities, LLMs have significantly made searches more accurate and efficient. This is stated about their performance in many NLP tasks, such as search and retrieval, where they have been found outperforming traditional keyword-based methods and far simpler machine learning models (Devlin et al., 2018).

Large Language Models (LLMs) have been demonstrated to be able to understand complex queries and return contextually relevant results. This sets a new benchmark for search technologies. Some of the most significant breakthroughs have emerged in the form of Retrieval-Augmented Generation (RAG) that relies on strengths both retrieval as well as generation approaches for augmenting the information retrieval ability, and precision (Gao et al., 2024). Such RAG models could then be used to avail the expansive contextual knowledge of an LLM to inject a retrieval mechanism that may, in turn, retrieve the correct documents or information for good results from large corpus.

As a result, the RAG model can then successfully work in both phases of retrieving information: first, extracting relevant information, and then composing logical answers that entail an intimate flow within the contextual framework set (Wu et al., 2024). It can thus be used when attempting to surmount disadvantages encountered with more traditional approaches to searching, such as usually failing to balance between retrieval speed and answer accurately (Gao et al., 2024). Despite these advancements, the limitations associated with LLMs cannot be ignored. Substantial computational demands exist while training and running these models, which creates severe efficiency and resource management concerns (Lan et al., 2019). The researchers have looked into several methods to improve these concerns, including model architecture optimization and incorporating techniques such as regular expressions to better enhance search performance.

2. Literature Review

The Large Language Models have greatly affected the world of IR, making several aspects of the search and retrieval system much better than those seen earlier. The whole ability to understand and generate text in the likeness of human languages has seen new paradigms regarding information indexing, searching, and retrieving.

Model Architectures and Advancements: The recent developments in large language model architectures have improved information retrieval systems concerning efficiency and effectiveness. When Devlin introduced BERT (Bidirectional Encoder Representations from Transformers) in 2018, context's determinism was finally unbounded, which revolutionized IR and significantly improved document retrieval and relevance assessment (Devlin et al., 2019). Similar things are going on with the Open AI GPT series, in particular, GPT-3, which can output coherent and contextual responses relevant to an application in query expansion and answer generation inside an information retrieval system (Brown et al., 2020).

Retrieval Performance and Evaluation: Better performances of LLMs have been found in Information Retrieval tasks such as query-document matching and relevance prediction. One of the main challenges in the retrieval task is advanced by the use of pseudo-relevant queries and LLM labels; they have helped in choosing the most effective dense retrievers-especially in conditions of domain shifts and unlabeled data (Khrantsova et al., 2024). Integrations with human-in-the-loop methodologies enhance the quality of LLM responses and, indeed, models such as GPT-4 have shown better performance in applications such as HR support, Education, Fraud and content moderation (Afzal et al., 2024).

Query Expansion and Contextual Understanding: (Wang et al., 2023) proposed a simple yet effective query expansion method, termed as query2doc, to enhance both sparse and dense retrieval systems utilizing large language models (LLMs), thereafter augmenting the query with created pseudo-documents. Kilbas et al. (2024) presented a technique that increased the contextualization size of LLMs through linear interpolation of positional embeddings, allowing such models as ruGPT-3.5 to take up to 6144 tokens, thus significantly increasing their robustness in processing long texts. (Luo et al., 2024) proposed extensible embeddings, allowing flexible extension of the context with a gain in model efficiency and compatibility.

Applications and Uses: LLMs have been implemented in a wide range of IR applications such as conversational search and document summarization. Therefore, LLMs are used for the development of the Chatbot that will answer questions from employees correctly by providing better response quality through human-in-the-loop techniques with state-of-the-art retrieval mechanisms (Afzal et al., 2024). In this way, the LLM can retrieve information buried underneath complex text data, such as legal and financial texts, using architectures dependent on RAG models (Hikov & Murphy, 2024).

3. ROLE OF LLMs IN INFORMATION RETRIEVAL

Large Language Models (LLMs) significantly impact the domain of information retrieval due to higher precision and effectiveness in retrieving information on the topic. They made possible the procedure of Retrieval-Augmented Generation (RAG), the procedure of merging LLM capabilities with external knowledge storages, which makes the response to factual questions more reliable (Sharma et al., 2024; Wang et al., 2024). However, there are still challenges; for example, systematic biases towards high-resource languages may perpetuate information silos and introduce marginalized views from low-resource contexts (Sharma et al., 2024). Even though the LLM revolutionizes access to information, it brings risks about misinformation and cost-effectiveness, so there is a need to collaborate towards solving those problems (Liu et al., 2024). The embedding of OCR along with LLMs has further improved the IR operations to transform document types into machine-readable form, which manifests synergy-driven developments along these lines and these can be embraced (Pakhale, 2023). In general, LLMs have undoubtedly increased IR capabilities, but it should be clearly looked at their limitations and biases to ensure proper usage.

4. IMPACT OF LLMs ON SPEED AND ACCURACY

A few techniques significantly improve the efficiency and accuracy of IR, such as embedding self-attention mechanisms within the mechanism by which LLMs learn to pay selective attention to information relevant to these models, remove redundant data, and greatly enhance processing efficiency and accuracy in various contexts of data (Li et al., 2024). Furthermore, frameworks rely on LLM's semantic knowledge capabilities to achieve highly accurate retrieval performance over 98.8% for user queries-using effective integration into external knowledge bases (Wang et al., 2024). For example, LLMs, Such as GPT-4, have significantly transformed this domain by facilitating generative retrieval and enhancing user interactions via Natural Language Processing (NLP) (Ai et al., 2023). Simultaneously, the real-time dynamic information feed integrated into LLMs avoids content hallucination and adds to contextual relevance in general, therefore building accuracy (Ouyang et al., 2023). The integration of Retrieval-Augmented Generation methodologies with context compression enhances the efficiency of computation while letting large language models handle complex quantities of contextual information without loss in output quality (Jiang et al., 2024). These advancements really express the potential of LLMs to transform information retrieval methods into faster and more reliable applications.

5. BUILDING LLMs FOR INFORMATION RETRIEVAL

This process involves leveraging the advanced capabilities of LLMs to enhance the efficiency and accuracy of retrieving relevant information from vast datasets and also integrates various strategies, including pre-training, fine-tuning, and the use of innovative architectures and frameworks. This Fig. 1 outlines a framework for building Large Language Models (LLMs) using Retrieval-Augmented Generation (RAG). It shows how a Retriever encodes, indexes, and retrieves data, which is then, combined with inputs through different **Fusion Techniques-Query-based, Logits-based, and Latent Fusion**. These fused inputs are passed to various Generators like GPT, and others to produce outputs. The diagram also highlights specific methods and models used within each fusion technique to improve the generation of contextually accurate responses (Wu et al., 2024). The following sections explore key aspects of building LLMs for IR.

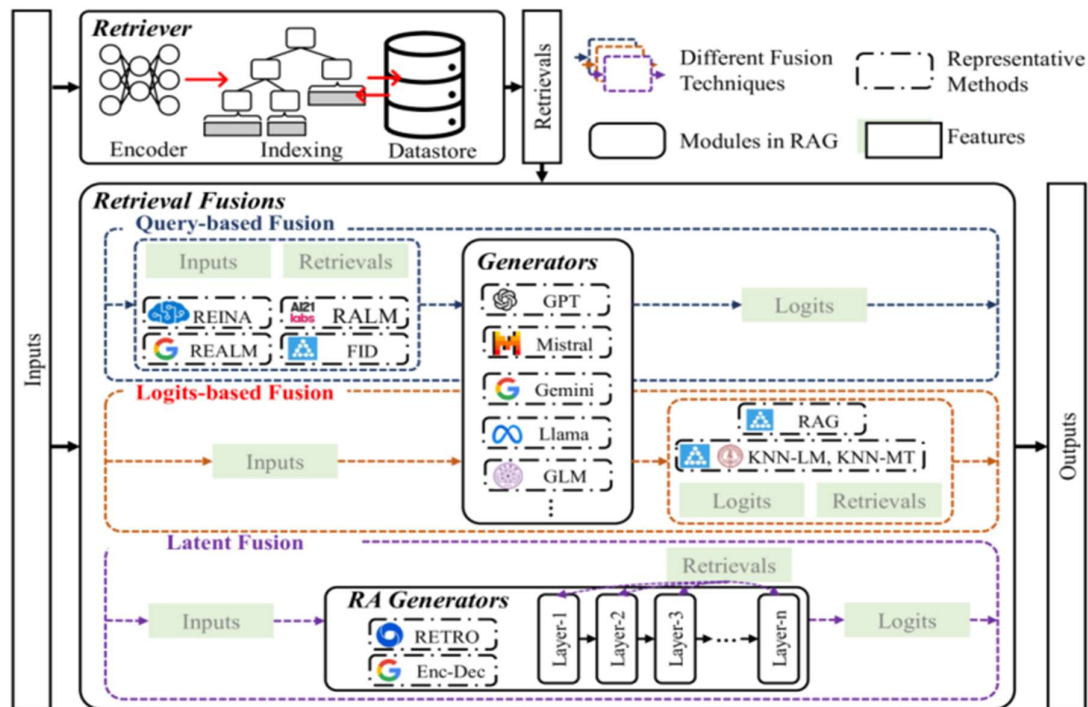


Fig. 1: The overview of Retrieval-Augmented Generation for Large Language Models. (Wu et al., 2024)

5.1 Pre-training and Fine-tuning Strategies

Pre-training on Diverse Data: They initially pre-trained on large-scale datasets that range from structured to unstructured. The pre-training step provides a foundation for learning massive data representation and primary language understanding. Fine-tuning of such models on specific tasks like annotating an image or queries for structured data enhances their capability to process and integrate information coming from various sources, hence improving the efficiency and accuracy of IR (Li et al., 2024).

Task-specific Adaptation: This includes task-specific layers that would involve semantic analysis in text as well as relational extraction in structured data, where the LLMs can attend to relevant information to dismiss the irrelevant data, which is the thrust in the IR application when dealing with mixed data types (Li et al., 2024).

5.2 Reactive and Self-Retrieval Pipelines

Reactive LLM Pipelines: The Motion framework represents the ultimate case of how reactive LLM pipelines let one inject new data such as user feedback and historical context within a prompt to enrich information retrieval: this process bumps up the quality of responses and adjusts to changing informational needs (Shankar & Parameswaran, 2024).

Self-Retrieval Architecture: The Self-Retrieval model inputs the whole corpus into the LLM, thus making retrieval a generation process that self-assesses. This type of end-to-end construction would dramatically increase the capabilities of applications involving LLMs such as Retrieval-Augmented Generation (Tang et al., 2024).

5.3 Multilingual and Domain-specific Challenges

Multilingual Capabilities and Bias: The bias of multilingual LLMs toward high-resource languages would then mean establishing dominant views and further marginalizing low-resource languages. These need to be counter-balanced in LLMs for the attainment of information parity in multilingual IR systems (Sharma et al., 2024).

Domain-specific Applications: In the healthcare industry, for example, the integration of LLMs into current models of information search transforms interpretation and retrieval of complex inquiries from databases to provide contextually relevant responses. This further enhances efficiency and accuracy in retrieval of information domains (Emdad & Rahman, 2024).

6. INNOVATIVE RETRIEVAL TECHNIQUES

Retrieval-Augmented Generation (RAG): The PG-RAG or Prompted Graph Retrieval-Augmented Generation. It indeed enhances the framework of information retrieval, bringing along the powers of LLMs within a graph-based structure and presents the potential of how LLMs could independently generate mental indices from unprocessed data, and then autonomously build mental indices aggregating it into a pseudo-graph database for structured retrieval. Indeed, it does work exceptionally well in both single-document as well as multi-document tasks, which often underlines the promise that RAG may offer to IR systems (Liang et al., 2024).

Cross-Domain Sequential Recommendation: The User Retrieval and Domain-Specific Large Language Model (URLLM) is the novel approach that combines mechanisms of user retrieval with large language models and domain-specific adaptations. This framework incorporates cross-domain recommendation in LLMs by integrating user retrieval and domain-specific generation. Therefore, this methodology receives varied information and provides space for the cross-application of domain knowledge so that recommendations may increase their accuracy (Shen et al., 2024).

While LLMs take great advancements in information retrieval, several issues remain that tend to include managing biases, cost-effectiveness, and misinformation handling. Moreover, academia and industry can be required to bring out improvements in these issues and fine-tune the IR systems following LLMs (Liu et al., 2024). This would include new decoding strategies, for example, Diver is a novel decoding strategy designed to promote the updating of the performance of language models and particularly IR tasks in mutual information verification. Interesting avenues for performance improvement in IR lie in this direction (Lu et al., 2024). Indeed, LLMs carry much hope; however, along with them come challenges in terms of bias, resource distribution, and a need for constant adaptation according to recent data and user feedbacks.

As shown in Fig. 2, how Retrieval-Augmented Generation (RAG) improves responses by incorporating significant external information into AI-generated responses. When a user searches the quires, the input is analyzed by a Large Language Model (LLM) that in the absence of Retrieval-Augmented Generation (RAG), delivers a generic response. Through RAG, the system acquires and catalogues pertinent papers, which then enhance the AI's answer, rendering it more knowledgeable and contextually precise. The conclusive response, produced by integrating the obtained data with the LLM, is more comprehensive and corresponds closely to the specific inquiry. (Gao et al., 2024).

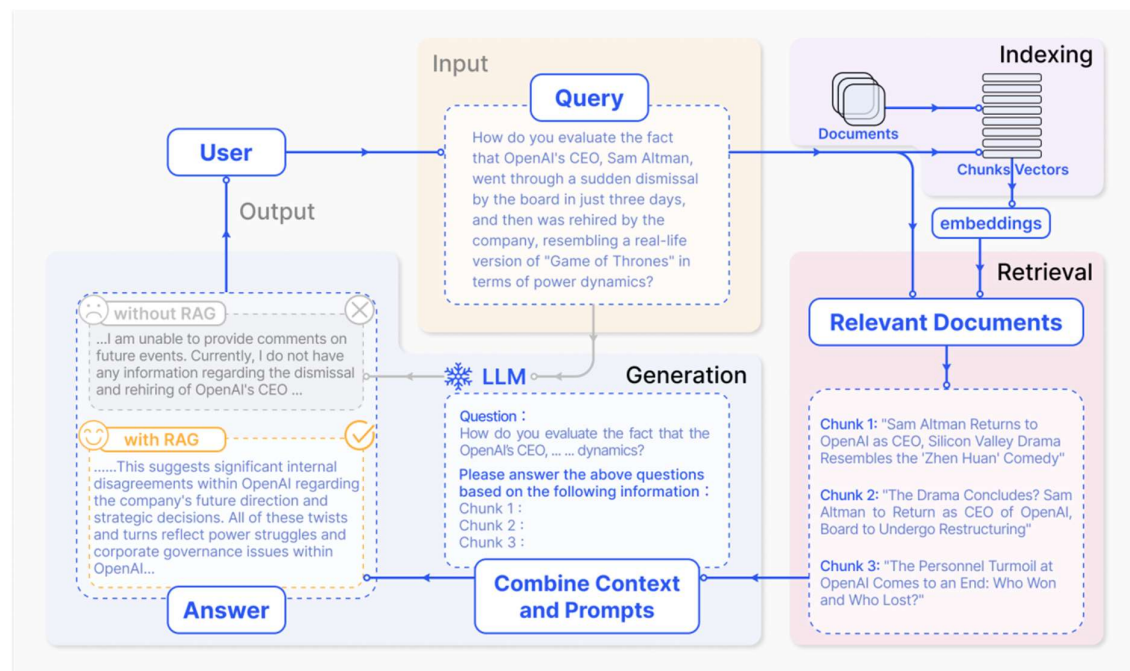


Fig. 2: An illustrative case of employing the RAG (Retrieve and Generate) approach for question answering. (Gao et al., 2024)

7. ADVANTAGES OF USING RAG FOR FASTER RETRIEVAL

Retrieval-Augmented Generation (RAG) provides substantial benefits for accelerated retrieval in large language models (LLMs) by combining external knowledge with model functionalities. The method aims to consistently provide enhanced efficiency and quality in information generation and retrieval. The subsequent sections specifically outline the advantages gained by employing RAG for rapid retrievals:

Sparse Context Selection: Sparse RAG introduces a mechanism of cutting the computation cost by encoding retrieved documents in parallel, cutting out latency from long-range attention. This makes LLMs only focus on high-relevance documents for their optimal speed and quality of generation (Zhu et al., 2024).

Efficient Retrieval Techniques: This results in a high accuracy of the recall of the enterprise RAG use of atomic units and synthetic queries, so the efficiency of the LLM is increased. The method enhances the rate of recall; therefore, there is a possibility of getting relevant information within a short time (Raina & Gales, 2024). RAG also brings in a new evaluation technique that significantly reduces computational costs and enhances runtime efficiency, making retrieval faster and economical (Salemi & Zamani, 2024).

Context Compression: COCOM, or Context Compression Method, this is one of the methods to make the data representation in Natural Language Processing more effective. Shortening the length of the contextual input reduces the time to generate a response. This compression is equivalent to the time for decoding and that for the quality of the response; it provides 5.69 times better performance than existing techniques (Rau et al., 2024).

Knowledge Caching: RAGCache models use an efficient caching system that stores the key pieces of information extracted from retrieved documents. Caching data points in such a manner minimizes redundant computation during similar requests. Intermediate states of retrieved knowledge are cached, in which case the retrieval process optimization is achieved. In this context, it brings up to four times speedup in Time to First Token (TTFT) and up to 2.1 times throughput speedup. This adaptive caching mechanism reduces latency by executing the retrieval and inference processes concurrently (Jin et al., 2024).

8. KEY CHALLENGES FACED BY LLMs

- **Bias and Ethical Concerns:** LLMs may reflect certain types of human-like biases that exist while they are in training; thus, the behavior tends to be often degenerated and biased. Such bias might hamper efforts by the model to fetch and generate unbiased information since it may consciously favor some views over others (Schramowski et al., 2022).
- **Spurious Correlations:** The LLMs learn statistical dependencies between variables unrelated due to bias while selecting the dataset. The wrong relations by the model therefore bring adverse impacts on the retrieved and generated information's reliability. (Subhash, 2023).
- **Adversarial Manipulation:** The LLMs lack interpretability in relation to results obtained from the processed information, that is, users can get manipulated or altered preferences. That is challenging in that one would want retrieved information not only to be relevant but also trustworthy (Arora & Arora, 2023).
- **Hallucination and Factuality Issues:** LLMs sometimes produce wrong or hallucinated information, thereby making their deployment not effective in the retrieval of accurate and reliable information. That is actually the case especially where there is the need for real-time accuracy for dynamic and interactive applications (Xie et al., 2023).

9. Result and Discussion

Combining the Retrieval-Augmented Generation model with LLMs has incredibly improved the accuracy, speed, and relevance while retrieving information. Moreover, LLMs such as GPT-3 and BERT can be much more descriptive than the keyword method, while RAG brings in some external knowledge sources to give more accurate answers. It is actually a balance between speed and quality, thus enhancing the user experience. However, the RAG systems require continuous improvement over different domains and languages, and two of the essential techniques, namely, the selection of context and knowledge caching improve the performance and reduce the waiting time.

10. Limitations

Large Language Models have a significant influence on information retrieval but at the same time carry many limitations. It spreads "hallucinations" of misleading information, making it suspicious to consider the correctness of the information presented because it cannot crosscheck with much reliability the credibility of sources (Liu et al., 2024). The evaluation of LLMs is challenging because the methods of assessment may differ and benchmarks are less standardized, especially in the case of more complex structures such as tables (Pang et al., 2024). LLMs are also biased more towards strong-resource languages sidelining low-resource languages, and hence, dominant cultural views get perpetuated (Sharma et al., 2024). They also face a problem in providing coherent information from conflicting sources and in more complicated tasks of reasoning which further constrains them in IR applications. With their ability to quickly shift attention to a random inconsequential fact, LLMs easily churn out by irrelevant information (Wu et al., 2024).

11. Conclusion

The integration of Retrieval-Augmented Generation (RAG) with Large Language Models (LLMs) may lead to an improvement in how fast and accurately information would be retrieved as well as by responses. Unlike this, the approach relies on external sources of knowledge to augment the factual accuracy and contextual relevance of large language model outputs. It is nevertheless very challenging to just see just how much RAG systems might be efficient and scalable. Recent work that has explored different mechanisms with regards to enhancing these systems includes but is not limited to retrieval speed, relevance in knowledge, and scalability of the system. There are also strict computational requirements and hallucinations possible from these models, and thus the reason these systems are still under active development. The IR system promises to do great things as LLMs can add to higher search accuracy and better user experience. They can also make the information easier to obtain so that any person can find their need.

REFERENCES

- [1] Afzal, A., Kowsik, A., Fani, R., & Matthes, F. (2024). Towards Optimizing and Evaluating a Retrieval Augmented QA Chatbot using LLMs with Human-in-the-Loop. In E. Dragut, Y. Li, L. Popa, S. Vucetic, & S. Srivastava (Eds.), *Proceedings of the Fifth Workshop on Data Science with Human-in-the-Loop (DaSH 2024)* (pp. 4–16). Association for Computational Linguistics. <https://doi.org/10.18653/v1/2024.dash-1.2>
- [2] Ai, Q., Bai, T., Cao, Z., Chang, Y., Chen, J., Chen, Z., Cheng, Z., Dong, S., Dou, Z., Feng, F., Gao, S., Guo, J., He, X., Lan, Y., Li, C., Liu, Y., Lyu, Z., Ma, W., Ma, J., Zhu, X. (2023). *Information Retrieval Meets Large Language Models: A Strategic Report from Chinese IR Community* (No. arXiv:2307.09751). arXiv. <https://doi.org/10.48550/arXiv.2307.09751>
- [3] Arora, A., & Arora, A. (2023). The promise of large language models in health care. *The Lancet*, 401(10377), 641. [https://doi.org/10.1016/S0140-6736\(23\)00216-7](https://doi.org/10.1016/S0140-6736(23)00216-7)
- [4] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Amodei, D. (2020a). *Language Models are Few-Shot Learners* (No. arXiv:2005.14165). arXiv. <https://doi.org/10.48550/arXiv.2005.14165>
- [5] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding* (No. arXiv:1810.04805). arXiv. <https://doi.org/10.48550/arXiv.1810.04805>
- [6] Emdad, F. B., & Rahman, M. I. (2024). *Potential Renovation of Information Search Process with the Power of Large Language Model for Healthcare* (No. arXiv:2407.01627). arXiv. <https://doi.org/10.48550/arXiv.2407.01627>
- [7] Gao, Y., Xiong, Y., Gao, X., Jia, K., Pan, J., Bi, Y., Dai, Y., Sun, J., Wang, M., & Wang, H. (2024). *Retrieval-Augmented Generation for Large Language Models: A Survey* (No. arXiv:2312.10997; Version 5). arXiv. <https://doi.org/10.48550/arXiv.2312.10997>

- [8] Herzog, H.-J., Vogel, G., & Schubert, U. (2002). LLM – a nonhydrostatic model applied to high-resolving simulations of turbulent fluxes over heterogeneous terrain. *Theoretical and Applied Climatology*, 73(1), 67–86. <https://doi.org/10.1007/s00704-002-0694-4>
- [9] Hikov, A., & Murphy, L. (2024). Information retrieval from textual data: Harnessing large language models, retrieval augmented generation and prompt engineering. *Journal of AI, Robotics & Workplace Automation*. <https://hstalks.com/article/8575/information-retrieval-from-textual-data-harnessing/>
- [10] Jiang, P., Fan, R., & Yong, Y. (2024). *Retrieval Augmented Generation via Context Compression Techniques for Large Language Models*. <https://doi.org/10.31219/osf.io/ua6j5>
- [11] Jin, C., Zhang, Z., Jiang, X., Liu, F., Liu, X., Liu, X., & Jin, X. (2024). *RAGCache: Efficient Knowledge Caching for Retrieval-Augmented Generation* (No. arXiv:2404.12457). arXiv. <https://doi.org/10.48550/arXiv.2404.12457>
- [12] Khrantsova, E., Zhuang, S., Baktashmotlagh, M., & Zuccon, G. (2024). Leveraging LLMs for Unsupervised Dense Retriever Ranking. *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 1307–1317. <https://doi.org/10.1145/3626772.3657798>
- [13] Kilbas, I., Gribanov, D., Mukhin, A., Paringer, R., & Kupriyanov, A. (2024). Expanding the Context of Large Language Models Via Linear Interpolation of Positional Embeddings. *2024 X International Conference on Information Technology and Nanotechnology (ITNT)*, 1–4. <https://doi.org/10.1109/ITNT60778.2024.10582292>
- [14] Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., & Soricut, R. (2019). *ALBERT: A Lite BERT for Self-supervised Learning of Language Representations*. <https://doi.org/10.48550/ARXIV.1909.11942>
- [15] Li, B., Jiang, G., Li, N., & Song, C. (2024). *Research on Large-scale Structured and Unstructured Data Processing based on Large Language Model* (No. 2024071364). Preprints. <https://doi.org/10.20944/preprints202407.1364.v1>
- [16] Liang, X., Niu, S., li, Z., Zhang, S., Song, S., Wang, H., Yang, J., Xiong, F., Tang, B., & Xi, C. (2024). *Empowering Large Language Models to Set up a Knowledge Retrieval Indexer via Self-Learning* (No. arXiv:2405.16933). arXiv. <https://doi.org/10.48550/arXiv.2405.16933>
- [17] Liu, Z., Zhou, Y., Zhu, Y., Lian, J., Li, C., Dou, Z., Lian, D., & Nie, J.-Y. (2024). Information Retrieval Meets Large Language Models. *Companion Proceedings of the ACM Web Conference 2024*, 1586–1589. <https://doi.org/10.1145/3589335.3641299>
- [18] Lu, J., Wang, C., & Zhang, J. (2024). *Diver: Large Language Model Decoding with Span-Level Mutual Information Verification* (No. arXiv:2406.02120). arXiv. <https://doi.org/10.48550/arXiv.2406.02120>
- [19] Luo, K., Liu, Z., Xiao, S., & Liu, K. (2024). *BGE Landmark Embedding: A Chunking-Free Embedding Method For Retrieval Augmented Long-Context Large Language Models* (No. arXiv:2402.11573). arXiv. <https://doi.org/10.48550/arXiv.2402.11573>
- [20] Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to Information Retrieval* (1st ed.). Cambridge University Press. <https://doi.org/10.1017/CBO9780511809071>
- [21] Mehta, M., D'souza, J., Karia, M., Kadam, V., Lad, M., & Therese, S. S. (2024). From Classical to Quantum: Evolution of Information Retrieval Systems. In S. Lanka, A. Sarasa-Cabezuelo, & A. Tugui (Eds.), *Trends in Sustainable Computing and Machine Intelligence* (pp. 299–312). Springer Nature. https://doi.org/10.1007/978-981-99-9436-6_21
- [22] Ouyang, Q., Wang, S., & Wang, B. (2023). *Enhancing Accuracy in Large Language Models Through Dynamic Real-Time Information Injection* (No. 2023121987). Preprints. <https://doi.org/10.20944/preprints202312.1987.v1>
- [23] Pakhale, K. (2023). Large Language Models and Information Retrieval. *IJFMR - International Journal For Multidisciplinary Research*, 5(6). <https://doi.org/10.36948/ijfmr.2023.v05i06.8841>

- [24] Pang, C., Cao, Y., Yang, C., & Luo, P. (2024). *Uncovering Limitations of Large Language Models in Information Seeking from Tables*. arXiv. <https://doi.org/10.48550/ARXIV.2406.04113>
- [25] Raina, V., & Gales, M. (2024). *Question-Based Retrieval using Atomic Units for Enterprise RAG* (No. arXiv:2405.12363). arXiv. <https://doi.org/10.48550/arXiv.2405.12363>
- [26] Rau, D., Wang, S., Déjean, H., & Clinchant, S. (2024). *Context Embeddings for Efficient Answer Generation in RAG* (No. arXiv:2407.09252). arXiv. <https://doi.org/10.48550/arXiv.2407.09252>
- [27] Salemi, A., & Zamani, H. (2024). *Evaluating Retrieval Quality in Retrieval-Augmented Generation* (No. arXiv:2404.13781). arXiv. <https://doi.org/10.48550/arXiv.2404.13781>
- [28] Schramowski, P., Turan, C., Andersen, N., Rothkopf, C. A., & Kersting, K. (2022). Large pre-trained language models contain human-like biases of what is right and wrong to do. *Nature Machine Intelligence*, 4(3), 258–268. <https://doi.org/10.1038/s42256-022-00458-8>
- [29] Shankar, S., & Parameswaran, A. G. (2024). Building Reactive Large Language Model Pipelines with Motion. *Companion of the 2024 International Conference on Management of Data*, 520–523. <https://doi.org/10.1145/3626246.3654734>
- [30] Sharma, N., Murray, K., & Xiao, Z. (2024). *Faux Polyglot: A Study on Information Disparity in Multilingual Large Language Models* (No. arXiv:2407.05502). arXiv. <https://doi.org/10.48550/arXiv.2407.05502>
- [31] Shen, T., Wang, H., Zhang, J., Zhao, S., Li, L., Chen, Z., Lian, D., & Chen, E. (2024). *Exploring User Retrieval Integration towards Large Language Models for Cross-Domain Sequential Recommendation* (No. arXiv:2406.03085). arXiv. <https://doi.org/10.48550/arXiv.2406.03085>
- [32] Subhash, V. (2023). *Can Large Language Models Change User Preference Adversarially?* (No. arXiv:2302.10291). arXiv. <https://doi.org/10.48550/arXiv.2302.10291>
- [33] Tang, Q., Chen, J., Yu, B., Lu, Y., Fu, C., Yu, H., Lin, H., Huang, F., He, B., Han, X., Sun, L., & Li, Y. (2024). *Self-Retrieval: Building an Information Retrieval System with One Large Language Model* (No. arXiv:2403.00801). arXiv. <https://doi.org/10.48550/arXiv.2403.00801>
- [34] Wang, L., Yang, N., & Wei, F. (2023). *Query2doc: Query Expansion with Large Language Models* (No. arXiv:2303.07678). arXiv. <https://doi.org/10.48550/arXiv.2303.07678>
- [35] Wang, M., Zhang, Y., Zhao, Q., Yang, J., & Zhang, H. (2024). *Redefining Information Retrieval of Structured Database via Large Language Models* (No. arXiv:2405.05508). arXiv. <https://doi.org/10.48550/arXiv.2405.05508>
- [36] Witten, I. H., Moffat, A., & Bell, T. C. (1999). *Managing gigabytes: compressing and indexing documents and images* (2nd ed). Morgan Kaufmann Publishers.
- [37] Wu, S., Xiong, Y., Cui, Y., Wu, H., Chen, C., Yuan, Y., Huang, L., Liu, X., Kuo, T.-W., Guan, N., & Xue, C. J. (2024, July 18). *Retrieval-Augmented Generation for Natural Language Processing: A Survey*. <https://arxiv.org/html/2407.13193v1>
- [38] Xie, T., Wan, Y., Huang, W., Zhou, Y., Liu, Y., Linghu, Q., Wang, S., Kit, C., Grazian, C., Zhang, W., & Hoex, B. (2023). *Large Language Models as Master Key: Unlocking the Secrets of Materials Science with GPT*. <https://doi.org/10.48550/ARXIV.2304.02213>
- [39] Zhu, Y., Gu, J.-C., Sikora, C., Ko, H., Liu, Y., Lin, C.-C., Shu, L., Luo, L., Meng, L., Liu, B., & Chen, J. (2024). *Accelerating Inference of Retrieval-Augmented Generation via Sparse Context Selection* (No. arXiv:2405.16178). arXiv. <https://doi.org/10.48550/arXiv.2405.16178>
- [40] Zhu, Y., Yuan, H., Wang, S., Liu, J., Liu, W., Deng, C., Chen, H., Dou, Z., & Wen, J.-R. (2024). *Large Language Models for Information Retrieval: A Survey* (No. arXiv:2308.07107). arXiv. <https://doi.org/10.48550/arXiv.2308.07107>