

Large Language Models for Information Retrieval: A Survey

Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng,
Haonan Chen, Zhicheng Dou, and Ji-Rong Wen

Abstract—As a primary means of information acquisition, information retrieval (IR) systems, such as search engines, have integrated themselves into our daily lives. These systems also serve as components of dialogue, question-answering, and recommender systems. The trajectory of IR has evolved dynamically from its origins in term-based methods to its integration with advanced neural models. While the neural models excel at capturing complex contextual signals and semantic nuances, thereby reshaping the IR landscape, they still face challenges such as data scarcity, interpretability, and the generation of contextually plausible yet potentially inaccurate responses. This evolution requires a combination of both traditional methods (such as term-based sparse retrieval methods with rapid response) and modern neural architectures (such as language models with powerful language understanding capacity). Meanwhile, the emergence of large language models (LLMs), typified by ChatGPT and GPT-4, has revolutionized natural language processing due to their remarkable language understanding, generation, generalization, and reasoning abilities. Consequently, recent research has sought to leverage LLMs to improve IR systems. Given the rapid evolution of this research trajectory, it is necessary to consolidate existing methodologies and provide nuanced insights through a comprehensive overview. In this survey, we delve into the confluence of LLMs and IR systems, including crucial aspects such as query rewriters, retrievers, rerankers, and readers. Additionally, we explore promising directions, such as search agents, within this expanding field.

Index Terms—Large Language Models; Information Retrieval; Query Rewrite; Rerank; Reader; Fine-tuning; Prompting



1 INTRODUCTION

INFORMATION access is one of the fundamental daily needs of human beings. To fulfill the need for rapid acquisition of desired information, various information retrieval (IR) systems have been developed [1–4]. Prominent examples include search engines such as Google, Bing, and Baidu, which serve as IR systems on the Internet, adept at retrieving relevant web pages in response to user queries, and provide convenient and efficient access to information on the Internet. It is worth noting that IR extends beyond web page retrieval. In dialogue systems (chatbots) [1, 5–8], such as Microsoft Xiaoice [2], Apple Siri,¹ and Google Assistant,² IR systems play a crucial role in retrieving appropriate responses to user input utterances, thereby producing natural and fluent human-machine conversations. Similarly, in question-answering systems [3, 9], IR systems are employed to select relevant clues essential for addressing user questions effectively. In image search engines [4], IR systems excel at returning images that align with user input queries. Given the exponential growth of information, research and industry have become increasingly interested in the development of effective IR systems.

The core function of an IR system is retrieval, which aims to determine the relevance between a user-issued query and the content to be retrieved, including various types of information such as texts, images, music, and more. For the scope of this survey, we concentrate solely on review-

ing those text retrieval systems, in which query-document relevance is commonly measured by their matching score.³ Given that IR systems operate on extensive repositories, the efficiency of retrieval algorithms becomes of paramount importance. To improve the user experience, the retrieval performance is enhanced from both the upstream (query reformulation) and downstream (reranking and reading) perspectives. As an upstream technique, query reformulation is designed to refine user queries so that they are more effective at retrieving relevant documents [10, 11]. With the recent surge in the popularity of conversational search, this technique has received increasing attention. On the downstream side, reranking approaches are developed to further adjust the document ranking [12–14]. In contrast to the retrieval stage, reranking is performed only on a limited set of relevant documents, already retrieved by the retriever. Under this circumstance, the emphasis is placed on achieving higher performance rather than keeping higher efficiency, allowing for the application of more complex approaches in the reranking process. Additionally, reranking can accommodate other specific requirements, such as personalization [15–18] and diversification [19–22]. Following the retrieval and reranking stages, a reading component is incorporated to summarize the retrieved documents and deliver a concise document to users [23, 24]. While traditional IR systems typically require users to gather and organize relevant information themselves; however, the reading component is an integral part of new IR systems such as New

All authors are from Gaoling School of Artificial Intelligence and School of Information, Renmin University of China.

Contact e-mail: yutaozhu94@gmail.com, dou@ruc.edu.cn

1. Apple Siri, <https://www.apple.com/siri/>
2. Google Assistant, <https://assistant.google.com/>

3. The term “document” will henceforth refer to any text-based content subject to retrieve, including both long articles and short passages.

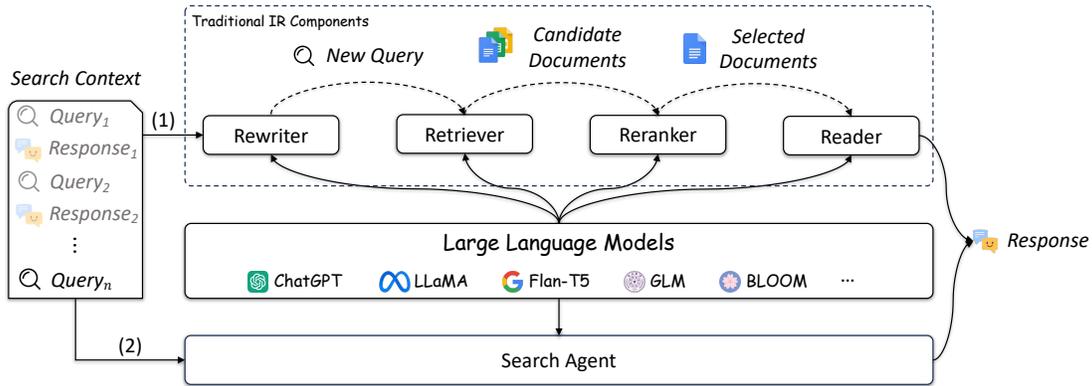


Fig. 1. Overview of existing studies that apply LLMs into IR. (1) LLMs can be used to enhance traditional IR components, such as query rewriter, retriever, reranker, and reader. (2) LLMs can also be used as search agents to perform multiple IR tasks.

Bing,⁴ streamlining users’ browsing experience and saving valuable time.

The trajectory of IR has traversed a dynamic evolution, transitioning from its origins in term-based methods to the integration of neural models. Initially, IR was anchored in term-based methods [25] and Boolean logic, focusing on keyword matching for document retrieval. The paradigm gradually shifted with the introduction of vector space models [26], unlocking the potential to capture nuanced semantic relationships between terms. This progression continued with statistical language models [27, 28], refining relevance estimation through contextual and probabilistic considerations. The influential BM25 algorithm [29] played an important role during this phase, revolutionizing relevance ranking by accounting for term frequency and document length variations. The most recent chapter in IR’s journey is marked by the ascendancy of neural models [3, 30–32]. These models excel at capturing intricate contextual cues and semantic nuances, reshaping the landscape of IR. However, these neural models still face challenges such as data scarcity, interpretability, and the potential generation of plausible yet inaccurate responses. Thus, the evolution of IR continues to be a journey of balancing traditional strengths (such as the BM25 algorithm’s high efficiency) with the remarkable capability (such as semantic understanding) brought about by modern neural architectures.

Large language models (LLMs) have recently emerged as transformative forces across various research fields, such as natural language processing (NLP) [33–35], recommender systems [36–39], finance [40], and even molecule discovery [41]. These cutting-edge LLMs are primarily based on the Transformer architecture and undergo extensive pre-training on diverse textual sources, including web pages, research articles, books, and codes. As their scale continues to expand (including both model size and data volume), LLMs have demonstrated remarkable advances in their capabilities. On the one hand, LLMs have exhibited unprecedented proficiency in language understanding and generation, resulting in responses that are more human-like and better align with human intentions. On the other hand, the larger LLMs have shown impressive emergent abilities

when dealing with complex tasks [42], such as generalization and reasoning skills. Notably, LLMs can effectively apply their learned knowledge and reasoning abilities to tackle new tasks with just a few task-specific demonstrations or appropriate instructions [43, 44]. Furthermore, advanced techniques, such as in-context learning, have significantly enhanced the generalization performance of LLMs without requiring fine-tuning on specific downstream tasks [34]. This breakthrough is particularly valuable, as it reduces the need for extensive fine-tuning while attaining remarkable task performance. Powered by prompting strategies such as chain-of-thought, LLMs can generate outputs with step-by-step reasoning, navigating complex decision-making processes [45]. Leveraging the impressive power of LLMs can undoubtedly improve the performance of IR systems. By incorporating these sophisticated language models, IR systems can provide users with more accurate responses, ultimately reshaping the landscape of information access and retrieval.

Initial efforts have been made to utilize the potential of LLMs in the development of novel IR systems. Notably, in terms of practical applications, New Bing is designed to improve the users’ experience of using search engines by extracting information from disparate web pages and condensing it into concise summaries that serve as responses to user-generated queries. In the research community, LLMs have proven useful within specific modules of IR systems (such as retrievers), thereby enhancing the overall performance of these systems. Due to the rapid evolution of LLM-enhanced IR systems, it is essential to comprehensively review their most recent advancements and challenges.

Our survey provides an insightful exploration of the intersection between LLMs and IR systems, covering key perspectives such as query rewriters, retrievers, rerankers, and readers (as shown in Figure 1).⁵ We also include some recent studies that leverage LLMs as search agents to perform various IR tasks. This analysis enhances our understanding of LLMs’ potential and limitations in advancing the IR field.

5. As yet, there has not been a formal definition for LLMs. In this paper, we mainly focus on models with more than 1B parameters. We also notice that some methods do not rely on such strictly defined LLMs, but due to their representativeness, we still include an introduction to them in this survey.

4. New Bing, <https://www.bing.com/new>

For this survey, we create a Github repository by collecting the relevant papers and resources about LLM4IR.⁶ We will continue to update the repository with newer papers. This survey will also be periodically updated according to the development of this area. We notice that there are several surveys for PLMs, LLMs, and their applications (e.g., AIGC or recommender systems) [46–52]. Among these, we highly recommend the survey of LLMs [52], which provides a systematic and comprehensive reference to many important aspects of LLMs. Compared with them, we focus on the techniques and methods for developing and applying LLMs for IR systems. In addition, we notice a perspective paper discussing the opportunity of IR when meeting LLMs [53]. It would be an excellent supplement to this survey regarding future directions.

The remaining part of this survey is organized as follows: Section 2 introduces the background for IR and LLMs. Section 3, 4, 5, 6 respectively review recent progress from the four perspectives of query rewriter, retriever, reranker, and reader, which are four key components of an IR system. Then, Section 8 discusses some potential directions in future research. Finally, we conclude the survey in Section 9 by summarizing the major findings.

2 BACKGROUND

2.1 Information Retrieval

Information retrieval (IR), as an essential branch of computer science, aims to efficiently retrieve information relevant to user queries from a large repository. Generally, users interact with the system by submitting their queries in textual form. Subsequently, IR systems undertake the task of matching and ranking these user-supplied queries against an indexed database, thereby facilitating the retrieval of the most pertinent results.

The field of IR has witnessed significant advancement with the emergence of various models over time. One such early model is the Boolean model, which employs Boolean logic operators to combine query terms and retrieve documents that satisfy specific conditions [25]. Based on the “bag-of-words” assumption, the vector space model [26] represents documents and queries as vectors in term-based space. Relevance estimation is then performed by assessing the lexical similarity between the query and document vectors. The efficiency of this model is further improved through the effective organization of text content using the inverted index. Moving towards more sophisticated approaches, statistical language models have been introduced to estimate the likelihood of term occurrences and incorporate context information, leading to more accurate and context-aware retrieval [27, 54]. In recent years, the neural IR [30, 55, 56] paradigm has gained considerable attention in the research community. By harnessing the powerful representation capabilities of neural networks, this paradigm can capture semantic relationships between queries and documents, thereby significantly enhancing retrieval performance.

Researchers have identified several challenges with implications for the performance and effectiveness of IR systems, such as query ambiguity and retrieval efficiency. In

light of these challenges, researchers have directed their attention toward crucial modules within the retrieval process, aiming to address specific issues and effectuate corresponding enhancements. The pivotal role of these modules in ameliorating the IR pipeline and elevating system performance cannot be overstated. In this survey, we focus on the following four modules, which have been greatly enhanced by LLMs.

Query Rewriter is an essential IR module that seeks to improve the precision and expressiveness of user queries. Positioned at the early stage of the IR pipeline, this module assumes the crucial role of refining or modifying the initial query to align more accurately with the user’s information requirements. As an integral part of query rewriting, query expansion techniques, with pseudo relevance feedback being a prominent example, represent the mainstream approach to achieving query expression refinement. In addition to its utility in improving search effectiveness across general scenarios, the query rewriter finds application in diverse specialized retrieval contexts, such as personalized search and conversational search, thus further demonstrating its significance.

Retriever, as discussed here, is typically employed in the early stages of IR for document recall. The evolution of retrieval technologies reflects a constant pursuit of more effective and efficient methods to address the challenges posed by ever-growing text collections. In numerous experiments on IR systems over the years, the classical “bag-of-words” model BM25 [29] has demonstrated its robust performance and high efficiency. In the wake of the neural IR paradigm’s ascendancy, prevalent approaches have primarily revolved around projecting queries and documents into high-dimensional vector spaces, and subsequently computing their relevance scores through inner product calculations. This paradigmatic shift enables a more efficient understanding of query-document relationships, leveraging the power of vector representations to capture semantic similarities.

Reranker, as another crucial module in the retrieval pipeline, primarily focuses on fine-grained reordering of documents within the retrieved document set. Different from the retriever, which emphasizes the balance of efficiency and effectiveness, the reranker module places a greater emphasis on the quality of document ranking. In pursuit of enhancing the search result quality, researchers delve into more complex matching methods than the traditional vector inner product, thereby furnishing richer matching signals to the reranker. Moreover, the reranker facilitates the adoption of specialized ranking strategies tailored to meet distinct user requirements, such as personalized and diversified search results. By integrating domain-specific objectives, the reranker module can deliver tailored and purposeful search results, enhancing the overall user experience.

Reader has evolved as a crucial module with the rapid development of LLM technologies. Its ability to comprehend real-time user intent and generate dynamic responses based on the retrieved text has revolutionized the presentation of IR results. In comparison to presenting a list of candidate

6. <https://github.com/RUC-NLPIR/LLM4IR-Survey>

documents, the reader module organizes answer texts more intuitively, simulating the natural way humans access information. To enhance the credibility of generated responses, the integration of references into generated responses has been an effective technique of the reader module.

Furthermore, researchers explore unifying the above modules to develop a novel LLM-driven search model known as the **Search Agent**. The search agent is distinguished by its simulation of an automated search and result understanding process, which furnishes users with accurate and readily comprehensible answers. WebGPT [24] serves as a pioneering work in this category, which models the search process as a sequence of actions of an LLM-based agent within a search engine environment, autonomously accomplishing the whole search pipeline. By integrating the existing search stack, search agents have the potential to become a new paradigm in future IR.

2.2 Large Language Models

Language models (LMs) are designed to calculate the generative likelihood of word sequences by taking into account the contextual information from preceding words, thereby predicting the probability of subsequent words. Consequently, by employing certain word selection strategies (such as greedy decoding or random sampling), LMs can proficiently generate natural language texts. Although the primary objective of LMs lies in text generation, recent studies [57] have revealed that a wide array of natural language processing problems can be effectively reformulated into a text-to-text format, thus rendering them amenable to resolution through text generation. This has led to LMs becoming the *de facto* solution for the majority of text-related problems.

The evolution of LMs can be categorized into four primary stages, as discussed in prior literature [52]. Initially, LMs were rooted in statistical learning techniques and were termed **statistical language models**. These models tackled the issue of word prediction by employing the Markov assumption to predict the subsequent word based on preceding words. Thereafter, neural networks, particularly recurrent neural networks (RNNs), were introduced to calculate the likelihood of text sequences and establish **neural language models**. These advancements made it feasible to utilize LMs for representation learning beyond mere word sequence modeling. ELMo [58] first proposed to learn contextualized word representations through pre-training a bidirectional LSTM (biLSTM) network on large-scale corpora, followed by fine-tuning on specific downstream tasks. Similarly, BERT [59] proposed to pre-train a Transformer [60] encoder with a specially designed Masked Language Modeling (MLM) task and Next Sentence Prediction (NSP) task on large corpora. These studies initiated a new era of **pre-trained language models** (PLMs), with the “pre-training then fine-tuning” paradigm emerging as the prevailing learning approach. Along this line, numerous generative PLMs (*e.g.*, GPT-2 [33], BART [61], and T5 [57]) have been developed for text generation problems including summarization, machine translation, and dialogue generation. Recently, researchers have observed that increasing the scale of PLMs (*e.g.*, model size or data amount) can

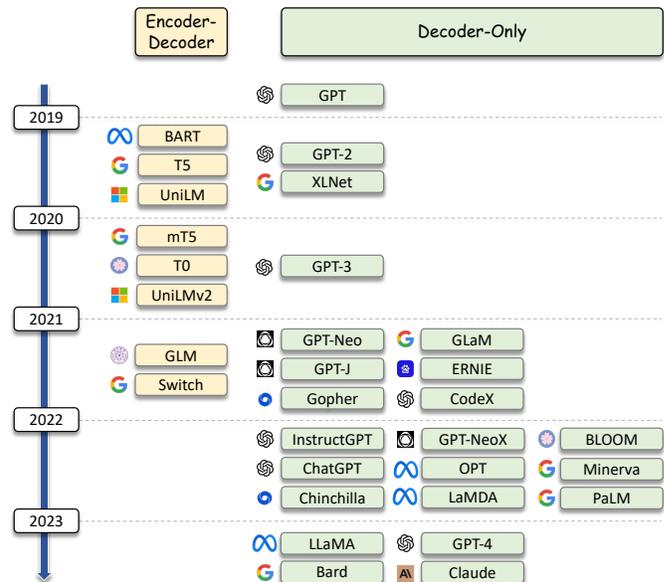


Fig. 2. The evolution of LLMs (encoder-decoder and decoder-only structures).

consistently improve their performance on downstream tasks (a phenomenon commonly referred to as the *scaling law* [62, 63]). Moreover, large-sized PLMs exhibit promising abilities (termed *emergent abilities* [42]) in addressing complex tasks, which are not evident in their smaller counterparts. Therefore, the research community refers to these large-sized PLMs as **large language models** (LLMs).

As shown in Figure 2, existing LLMs can be categorized into two groups based on their architectures: encoder-decoder [57, 61, 64–69] and decoder-only [33–35, 70–80] models. The encoder-decoder models incorporate an encoder component to transform the input text into vectors, which are then employed for producing output texts. For example, T5 [57] is an encoder-decoder model that converts each natural language processing problem into a text-to-text form and resolves it as a text generation problem. In contrast, decoder-only models, typified by GPT, rely on the Transformer decoder architecture. It uses a self-attention mechanism with a diagonal attention mask to generate a sequence of words from left to right. Building upon the success of GPT-3 [34], which is the first model to encompass over 100B parameters, several noteworthy models have been inspired, including GPT-J, BLOOM [78], OPT [75], Chinchilla [81], and LLaMA [35]. These models follow the similar Transformer decoder structure as GPT-3 and are trained on various combinations of datasets.

Owing to their vast number of parameters, fine-tuning LLMs for specific tasks, such as IR, is often deemed impractical. Consequently, two prevailing methods for applying LLMs have been established: in-context learning (ICL) and parameter-efficient fine-tuning. ICL is one of the emergent abilities of LLMs [34] empowering them to comprehend and furnish answers based on the provided input context, rather than relying merely on their pre-training knowledge. This method requires only the formulation of the task description and demonstrations in natural language, which are then fed as input to the LLM. Notably, parameter tuning is not

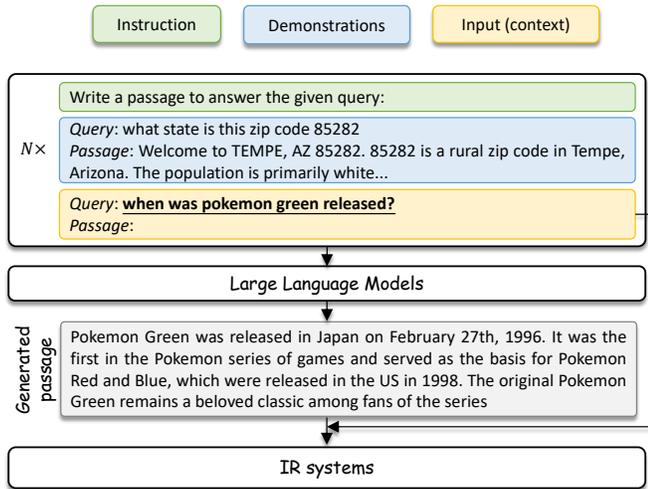


Fig. 3. An example of LLM-based query rewriting for ad-hoc search. The example is cited from the Query2Doc paper [86]. LLMs are used to generate a passage to supplement the original query, where $N = 0$ and $N > 0$ correspond to zero-shot and few-shot scenarios.

required for ICL. Additionally, the efficacy of ICL can be further augmented through the adoption of chain-of-thought prompting, involving multiple demonstrations (describe the chain of thought examples) to guide the model’s reasoning process. ICL is the most commonly used method for applying LLMs to IR. Parameter-efficient fine-tuning [82–84] aims to reduce the number of trainable parameters while maintaining satisfactory performance. LoRA [82], for example, has been widely applied to open-source LLMs (e.g., LLaMA and BLOOM) for this purpose. Recently, QLoRA [85] has been proposed to further reduce memory usage by leveraging a frozen 4-bit quantized LLM for gradient computation. Despite the exploration of parameter-efficient fine-tuning for various NLP tasks, its implementation in IR tasks remains relatively limited, representing a potential avenue for future research.

3 QUERY REWRITER

Query rewriting in modern IR systems is essential for improving search query effectiveness and accuracy. It reformulates users’ original queries to better match search results, alleviating issues like vague queries or vocabulary mismatches between the query and target documents. This task goes beyond mere synonym replacement, requiring an understanding of user intent and query context, particularly in complex searches like conversational queries. Effective query rewriting enhances search engine performance.

Traditional methods for query rewriting improve retrieval performance by expanding the initial query with information from highly-ranked relevant documents. Mainly-used methods include relevance feedback [87–92], word-embedding based methods [93, 94] etc. However, the limited ability of semantic understanding and comprehension of user search intent limits their performance in capturing the full scope of user intent.

Recent advancements in LLMs present promising opportunities to boost query rewriting capabilities. On one hand,

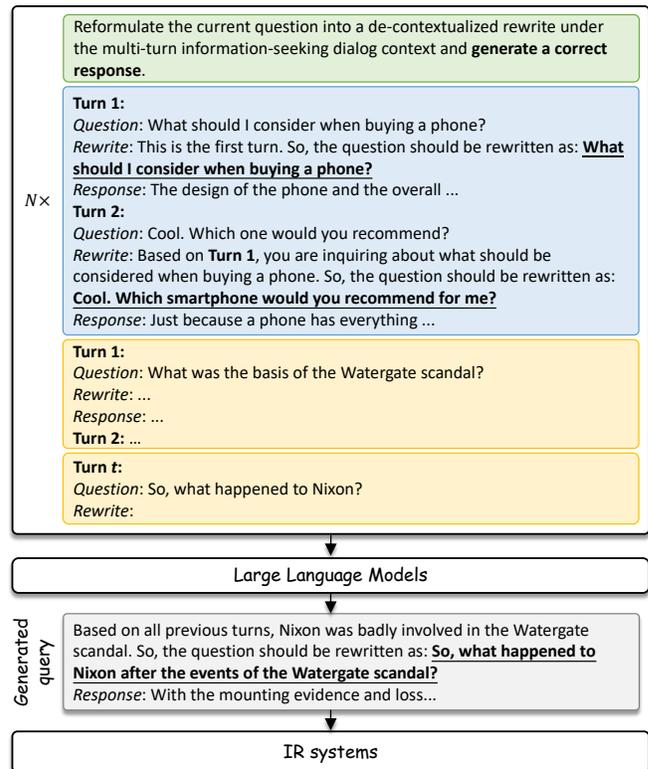


Fig. 4. An example of LLM-based query rewriting for conversational search. The example is cited from LLMCS [95]. The LLM is used to generate a query based on the demonstrations and previous search context. Additional responses are required to be generated for improving the query understanding. $N = 0$ and $N > 0$ correspond to zero-shot and few-shot scenarios.

given the context and subtleties of a query, LLMs can provide more accurate and contextually relevant rewrites. On the other hand, LLMs can leverage their extensive knowledge to generate synonyms and related concepts, enhancing queries to cover a broader range of relevant documents, thereby effectively addressing the vocabulary mismatch problem. In the following sections, we will introduce the recent works that employ LLMs in query rewriting.

3.1 Rewriting Scenario

Query rewriting typically serves two scenarios: ad-hoc retrieval, which mainly addresses vocabulary mismatches between queries and candidate documents, and conversational search, which refines queries based on evolving conversations. The upcoming section will delve into the role of query rewriting in these two domains and explore how LLMs enhance this process.

3.1.1 Ad-hoc Retrieval

In ad-hoc retrieval, queries are often short and ambiguous. In such scenarios, the main objectives of query rewriting include adding synonyms or related terms to address vocabulary mismatches and clarifying ambiguous queries to more accurately align with user intent. From this perspective, LLMs have inherent advantages in query rewriting.

Primarily, LLMs have a deep understanding of language semantics, allowing them to capture the meaning of queries more effectively. Besides, LLMs can leverage their extensive training on diverse datasets to generate contextually relevant synonyms and expand queries, ensuring broader and more precise search result coverage. Additionally, studies have shown that LLMs’ integration of external factual corpora [96–99] and thoughtful model design [100] further enhance their accuracy in generating effective query rewrites, especially for specific tasks.

Currently, there are many studies leveraging LLMs to rewrite queries in adhoc retrieval. We introduce the typical method Query2Doc [86] as an example. As shown in Figure 3, Query2Doc prompts the LLMs to generate a relevant passage according to the original query (“when was pokemon green released?”). Subsequently, the original query is expanded by incorporating the generated passage. The retriever module uses this new query to retrieve a list of relevant documents. Notably, the generated passage contains additional detailed information, such as “Pokemon Green was released in Japan on February 27th”, which effectively mitigates the “vocabulary mismatch” issue to some extent.

In addition to addressing the “vocabulary mismatch” problem [96–99, 101, 102], other works utilize LLMs for different challenges in ad-hoc retrieval. For instance, Prompt-Case [103] leverages LLMs in legal case retrieval to simplify complex queries into more searchable forms. This involves using LLMs to identify legal facts and issues, followed by a prompt-based encoding scheme for effective language model encoding.

3.1.2 Conversational Search

Query rewrites in conversational search play a pivotal role in enhancing the search experience. Unlike traditional queries in ad-hoc retrieval, conversational search involves a dialogue-like interaction, where the context and user intent evolve with each interaction. In conversational search, query rewriting involves understanding the entire conversation’s context, clarifying any ambiguities, and personalizing responses based on user history. The process includes dynamic query expansion and refinement based on dialogue information. This makes conversational query rewriting a sophisticated task that goes beyond traditional search, focusing on natural language understanding and user-centric interaction.

In the era of LLMs, leveraging LLMs in conversational search tasks offers several advantages. First, LLMs possess strong contextual understanding capabilities, enabling them to better comprehend users’ search intent within the context of multi-turn conversations between users and the system. Second, LLMs exhibit powerful generation abilities, allowing them to simulate dialogues between users and the system, thereby facilitating more robust search intent modeling.

The LLMCS framework [95] is a pioneering approach that employs LLMs to effectively extract and understand user search intent within conversational contexts. As illustrated in their work, LLMCS uses LLMs to produce both query rewrites and extensive hypothetical system responses from various perspectives. These outputs are combined

into a comprehensive representation that effectively captures the user’s full search intent. The experimental results show that including detailed hypothetical responses with concise query rewrites markedly improves search performance by adding more plausible search intent. Ye et al. [104] claims that human query rewrite may lack sufficient information for optimal retrieval performance. It defines four essential properties for well-formed LLM-generated query rewrites. Results show that their method informative query rewrites can yield substantially improved retrieval performance compared to human rewrites.

Besides, LLMs can be used as a data expansion tool in conversational dense retrieval. Attributed to the high cost of producing hand-written dialogues, data scarcity presents a significant challenge in the domain of conversational search. To address this problem, CONVERSER [105] employs LLMs to generate synthetic passage-dialogue pairs through few-shot demonstrations. Furthermore, it efficiently trains a dense retriever using a minimal dataset of six in-domain dialogues, thus mitigating the issue of data sparsity.

3.2 Rewriting Knowledge

Query rewriting typically necessitates additional corpora for refining initial queries. Considering that LLMs incorporate world knowledge in their parameters, they are naturally capable of rewriting queries. We refer to these methods, which rely exclusively on the intrinsic knowledge of LLMs, as LLM-only methods. While LLMs encompass a broad spectrum of knowledge, they may be inadequate in specialized areas. Furthermore, LLMs can introduce concept drift, leading to noisy relevance signals. To address this issue, some methods incorporate domain-specific corpora to provide more detailed and relevant information in query rewriting. We refer to methods enhanced by domain-specific corpora to boost LLM performance as corpus-enhanced LLM-based methods. In this section, we will introduce these two methods in detail.

3.2.1 LLM-only methods

LLMs are capable of storing knowledge within their parameters, making it a natural choice to capitalize on this knowledge for the purpose of query rewriting. As a pioneering work in LLM-based query rewriting, HyDE [101] generates a hypothetical document by LLMs according to the given query and then uses a dense retriever to retrieve relevant documents from the corpus that are relevant to the generated document. Query2doc [86] generates pseudo documents via prompting LLMs with few-shot demonstrations, and then expands the query with the generated pseudo document. Furthermore, the influence of different prompting methods and various model sizes on query rewriting has also been investigated [102]. To better accommodate the frozen retriever and the LLM-based reader, a small language model is employed as the rewriter that is trained using reinforcement learning techniques with the rewards provided by the LLM-based reader [100]. GFF [106] presents a “Generate, Filter, and Fuse” method for query expansion. It employs an LLM to create a set of related keywords via a reasoning chain. Then, a self-consistency filter is used to identify the most important keywords, which are

concatenated with the original queries for the downstream reranking task.

It is worth noting that though the designs of these methods are different, all of them rely on the world knowledge stored in LLMs without additional corpora.

3.2.2 Corpus-enhanced LLM-based methods

Although LLMs exhibit remarkable capabilities, the lack of domain-specific knowledge may lead to the generation of hallucinatory or irrelevant queries. To address this issue, recent studies [96–99] have proposed a hybrid approach that enhances LLM-based query rewriting methods with an external document corpus.

Why incorporate a document corpus? The integration of a document corpus offers several notable advantages. Firstly, it boosts relevance by using relevant documents to refine query generation, reducing irrelevant content and improving contextually appropriate outputs. Second, enhancing LLMs with up-to-date information and specialized knowledge in specific fields enables them to effectively deal with queries that are both current and specific to certain domains.

How to incorporate a document corpus? Thanks to the flexibility of LLMs, various paradigms have been proposed to incorporate a document corpus into LLM-based query rewriting, which can be summarized as follows.

- *Late fusion of LLM-based re-writing and pseudo relevance feedback (PRF) retrieval results.* Traditional PRF methods leverage relevant documents retrieved from a document corpus to rewrite queries, which restricts the query to the information contained in the target corpus. On the contrary, LLM-based rewriting methods provide external context not present in the corpus, which is more diverse. Both approaches have the potential to independently enhance retrieval performance. Therefore, a straightforward strategy for combining them is using a weighted fusion method for retrieval results [99].

- *Combining retrieved relevant documents in the prompts of LLMs.* In the era of LLMs, incorporating instructions within the prompts is the most flexible method for achieving specific functionalities. QUILL [97] and CAR [107] illustrate how retrieval augmentation of queries can provide LLMs with context that significantly enhances query understanding. LameR [108] takes this further by using LLM expansion to improve the simple BM25 retriever, introducing a retrieve-rewrite-retrieve framework. Experimental results reveal that even basic term-based retrievers can achieve comparable performance when paired with LLM-based rewriters. Additionally, InteR [98] proposes a multi-turn interaction framework between search engines and LLMs. This enables search engines to expand queries using LLM-generated insights, while LLMs refine prompts using relevant documents sourced from the search engines.

- *Enhancing factuality of generative relevance feedback (GRF) by pseudo relevance feedback (PRF).* Although generative documents are often relevant and diverse, they exhibit hallucinatory characteristics. In contrast, traditional documents are generally regarded as reliable sources of factual information. Motivated by this observation, GRM [96] proposes a novel technique known as relevance-aware sample estimation (RASE). RASE leverages relevant documents retrieved from

TABLE 1. Partial Examples of different prompting methods in query rewriting.

Methods	Prompts
<i>Zero-shot</i>	
HyDE [101]	Please write a passage to answer the question. Question: {#Question} Passage:
LameR [108]	Give a question {#Question} and its possible answering passages: A. {#Passage 1} B. {#Passage 2} C. {#Passage 3} ... Please write a correct answering passage.
<i>Few-shot</i>	
Query2Doc [101]	Write a passage that answers the given query: Query: {#Query 1} Passage: {#Passage 1} ... Query: {#Query} Passage:
<i>Chain-of-Thought</i>	
CoT [102]	Answer the following query based on the context: Context: {#PRF doc 1} {#PRF doc 2} {#PRF doc 3} Query: {#Query} Give the rationale before answering

the collection to assign weights to generated documents. In this way, GRM ensures that relevance feedback is not only diverse but also maintains a high degree of factuality.

3.3 Rewriting Approaches

There are three main approaches used for leveraging LLMs in query rewriting: *prompting methods*, *fine-tuning*, and *knowledge distillation*. Prompting methods involve using specific prompts to direct LLM output, providing flexibility and interpretability. Fine-tuning adjusts pre-trained LLMs on specific datasets or tasks to improve domain-specific performance, mitigating the general nature of LLM world knowledge. Knowledge distillation, on the other hand, transfers LLM knowledge to lightweight models, simplifying the complexity associated with retrieval augmentation. In the following section, we will introduce these three methods in detail.

3.3.1 Prompting

Prompting in LLMs refers to the technique of providing a specific instruction or context to guide the model’s generation of text. The prompt serves as a conditioning signal and influences the language generation process of the model. Existing prompting strategies can be roughly categorized into three groups: zero-shot prompting, few-shot prompting, and chain-of-thought (CoT) prompting [45].

- *Zero-shot prompting.* Zero-shot prompting involves instructing the model to generate texts on a specific topic without any prior exposure to training examples in that domain or topic. The model relies on its pre-existing knowledge and language understanding to generate coherent and contextually relevant expanded terms for original queries. Experiments show that zero-shot prompting is a simple yet effective method for query rewriting [98, 99, 102, 108–110].

- *Few-shot prompting.* Few-shot prompting, also known as in-context learning, involves providing the model with a limited set of examples or demonstrations related to the

desired task or domain [86, 102, 109, 110]. These examples serve as a form of explicit instruction, allowing the model to adapt its language generation to the specific task or domain at hand. Query2Doc [86] prompts LLMs to write a document that answers the query with some demo query-document pairs provided by the ranking dataset, such as MSMARCO [111] and NQ [112]. This work experiments with a single prompt. To further study the impact of different prompt designing, recent works [102] have explored eight different prompts, such as prompting LLMs to generate query expansion terms instead of entire pseudo documents and CoT prompting. There are some illustrative prompts in Table 1. This work conducts more experiments than Query2Doc, but the results show that the proposed prompt is less effective than Query2Doc.

- *Chain-of-thought prompting.* CoT prompting [45] is a strategy that involves iterative prompting, where the model is provided with a sequence of instructions or partial outputs [102, 109]. In conversational search, the process of query re-writing is multi-turn, which means queries should be refined step-by-step with the interaction between search engines and users. This process is naturally coincided with CoT process. As shown in 4, users can conduct the CoT process through adding some instructions during each turn, such as “Based on all previous turns, xxx”. While in ad-hoc search, there is only one-round in query re-writing, so CoT could only be accomplished in a simple and coarse way. For example, as shown in Table 1, researchers add “Give the rationale before answering” in the instructions to prompt LLMs think deeply [102].

3.3.2 Fine-tuning

Fine-tuning is an effective approach for adapting LLMs to specific domains. This process usually starts with a pre-trained language model, like GPT-3, which is then further trained on a dataset tailored to the target domain. This domain-specific training enables the LLM to learn unique patterns, terminology, and context relevant to the domain, which is able to improve its capacity to produce high-quality query rewrites.

BEQUE [113] leverages LLMs for rewriting queries in e-commerce product searches. It designs three Supervised Fine-Tuning (SFT) tasks: quality classification of e-commerce query rewrites, product title prediction, and CoT query rewriting. To our knowledge, it is the first model to directly fine-tune LLMs, including ChatGLM [68, 114], ChatGLM2.0 [68, 114], Baichuan [115], and Qwen [116], specifically for the query rewriting task. After the SFT stage, BEQUE uses an offline system to gather feedback on the rewrites and further aligns the rewriters with e-commerce search objectives through an object alignment stage. Online A/B testing demonstrates the effectiveness of the method.

3.3.3 Knowledge Distillation

Although LLM-based methods have demonstrated significant improvements in query rewriting tasks, their practical implementation for online deployment is hindered by the substantial latency caused by the computational requirements of LLMs. To address this challenge, knowledge distillation has emerged as a prominent technique in the

TABLE 2. Summary of existing LLM-enhanced query rewriting methods. “Docs” and “KD” stand for document corpus and knowledge distillation, respectively.

Methods	Target	Data	Generation
HyDE [97]	Ad-hoc	LLMs	Prompting
Jagerman et al. [102]	Ad-hoc	LLMs	Prompting
Query2Doc [86]	Ad-hoc	LLMs	Prompting
Ma et al. [100]	Ad-hoc	LLMs	Finetuning
PromptCase [103]	Ad-hoc	LLMs	Prompting
GRF+PRF [99]	Ad-hoc	LLMs + Docs	Prompting
GRM [96]	Ad-hoc	LLMs + Docs	Prompting
Inter [98]	Ad-hoc	LLMs + Docs	Prompting
LameR [108]	Ad-hoc	LLMs + Docs	Prompting
CAR [107]	Ad-hoc	LLMs + Docs	Prompting
QUILL [97]	Ad-hoc	LLMs + Docs	KD & Finetuning
LLMCS [95]	Conversational	LLMs	Prompting
CONVERSER [105]	Conversational	LLMs	Prompting
Ye et al. [104]	Conversational	LLMs	Prompting

industry. In the QUILL [97] framework, a two-stage distillation method is proposed. This approach entails utilizing a retrieval-augmented LLM as the professor model, a vanilla LLM as the teacher model, and a lightweight BERT model as the student model. The professor model is trained on two extensive datasets, namely Orcas-I [117] and EComm [97], which are specifically curated for query intent understanding. Subsequently, a two-stage distillation process is employed to transfer knowledge from the professor model to the teacher model, followed by knowledge transfer from the teacher model to the student model. Empirical findings demonstrate that this knowledge distillation methodology surpasses the simple scaling up of model size from base to XXL, resulting in even more substantial improvements. In a recently proposed “rewrite-retrieve-read” framework [100], an LLM is first used to rewrite the queries by prompting, followed by a retrieval-augmented reading process. To improve framework effectiveness, a trainable rewriter, implemented as a small language model, is incorporated to further adapt search queries to align with both the frozen retriever and the LLM reader’s requirements. The rewriter’s refinement involves a two-step training process. Initially, supervised warm-up training is conducted using pseudo data. Then, the retrieve-then-read pipeline is described as a reinforcement learning scenario, with the rewriter’s training acting as a policy model to maximize pipeline performance rewards.

3.4 Limitations

While LLMs offer promising capabilities for query rewriting, they also meet several challenges. Here, we outline two main limitations of LLM-based query rewriters.

3.4.1 Concept Drifts

When using LLMs for query rewriting, they may introduce unrelated information, known as concept drift, due to their extensive knowledge base and tendency to produce detailed and redundant content. While this can enrich the query, it also risks generating irrelevant or off-target results.

This phenomenon has been reported in several studies [107, 113, 118] These studies highlight the need for a balanced approach in LLM-based query rewriting, ensuring

that the essence and focus of the original query are maintained while leveraging the LLM’s ability to enhance and clarify the query. This balance is crucial for effective search and IR applications.

3.4.2 Correlation between Retrieval Performance and Expansion Effects

Recently, a comprehensive study [119] conduct experiments on various expansion techniques and downstream ranking models, which reveals a notable negative correlation between retriever performance and the benefits of expansion. Specifically, while expansion tends to enhance the scores of weaker models, it generally hurts stronger models. This observation suggests a strategic approach: employ expansions with weaker models or in scenarios where the target dataset substantially differs in format from the training corpus. In other cases, it is advisable to avoid expansions to maintain clarity of the relevance signal.

4 RETRIEVER

In an IR system, the retriever serves as the first-pass document filter to collect broadly relevant documents for user queries. Given the enormous amounts of documents in an IR system, the retriever’s efficiency in locating relevant documents is essential for maintaining search engine performance. Meanwhile, a high recall is also important for the retriever, as the retrieved documents are then fed into the ranker to generate final results for users, which determines the ranking quality of search engines.

In recent years, retrieval models have shifted from relying on statistic algorithms [29] to neural models [3, 31]. The latter approaches exhibit superior semantic capability and excel at understanding complicated user intent. The success of neural retrievers relies on two key factors: *data* and *model*. From the data perspective, a large amount of high-quality training data is essential. This enables retrievers to acquire comprehensive knowledge and accurate matching patterns. Furthermore, the intrinsic quality of search data, *i.e.*, issued queries and document corpus, significantly influences retrieval performance. From the model perspective, a strongly representational neural architecture allows retrievers to effectively store and apply knowledge obtained from the training data.

Unfortunately, there are some long-term challenges that hinder the advancement of retrieval models. First, user queries are usually short and ambiguous, making it difficult to precisely understand the user’s search intents for retrievers. Second, documents typically contain lengthy content and substantial noise, posing challenges in encoding long documents and extracting relevant information for retrieval models. Additionally, the collection of human-annotated relevance labels is time-consuming and costly. It restricts the retrievers’ knowledge boundaries and their ability to generalize across different application domains. Moreover, existing model architectures, primarily built on BERT [59], exhibit inherent limitations, thereby constraining the performance potential of retrievers. Recently, LLMs have exhibited extraordinary abilities in language understanding, text generation, and reasoning. This has motivated researchers to use these abilities to tackle the aforementioned challenges

and aid in developing superior retrieval models. Roughly, these studies can be categorized into two groups, *i.e.*, (1) leveraging LLMs to generate search data, and (2) employing LLMs to enhance model architecture.

4.1 Leveraging LLMs to Generate Search Data

In light of the quality and quantity of search data, there are two prevalent perspectives on how to improve retrieval performance via LLMs. The first perspective revolves around search data refinement methods, which concentrate on reformulating input queries to precisely present user intents. The second perspective involves training data augmentation methods, which leverage LLMs’ generation ability to enlarge the training data for dense retrieval models, particularly in zero- or few-shot scenarios.

4.1.1 Search Data Refinement

Typically, input queries consist of short sentences or keyword-based phrases that may be ambiguous and contain multiple possible user intents. Accurately determining the specific user intent is essential in such cases. Moreover, documents usually contain redundant or noisy information, which poses a challenge for retrievers to extract relevance signals between queries and documents. Leveraging the strong text understanding and generation capabilities of LLMs offers a promising solution to these challenges. As yet, research efforts in this domain primarily concentrate on employing LLMs as query rewriters, aiming to refine input queries for more precise expressions of the user’s search intent. Section 3 has provided a comprehensive overview of these studies, so this section refrains from further elaboration. In addition to query rewriting, an intriguing avenue for exploration involves using LLMs to enhance the effectiveness of retrieval by refining lengthy documents. This intriguing area remains open for further investigation and advancement.

4.1.2 Training Data Augmentation

Due to the expensive economic and time costs of human-annotated labels, a common problem in training neural retrieval models is the lack of training data. Fortunately, the excellent capability of LLMs in text generation offers a potential solution. A key research focus lies in devising strategies to leverage LLMs’ capabilities to generate pseudo-relevant signals and augment the training dataset for the retrieval task.

Why do we need data augmentation? Previous studies of neural retrieval models focused on supervised learning, namely training retrieval models using labeled data from specific domains. For example, MS MARCO [111] provides a vast repository, containing a million passages, more than 200,000 documents, and 100,000 queries with human-annotated relevance labels, which has greatly facilitated the development of supervised retrieval models. However, this paradigm inherently constrains the retriever’s generalization ability for out-of-distribution data from other domains. The application spectrum of retrieval models varies from natural question-answering to biomedical IR, and it is expensive to annotate relevance labels for data from different domains. As a result, there is an emerging need for zero-shot

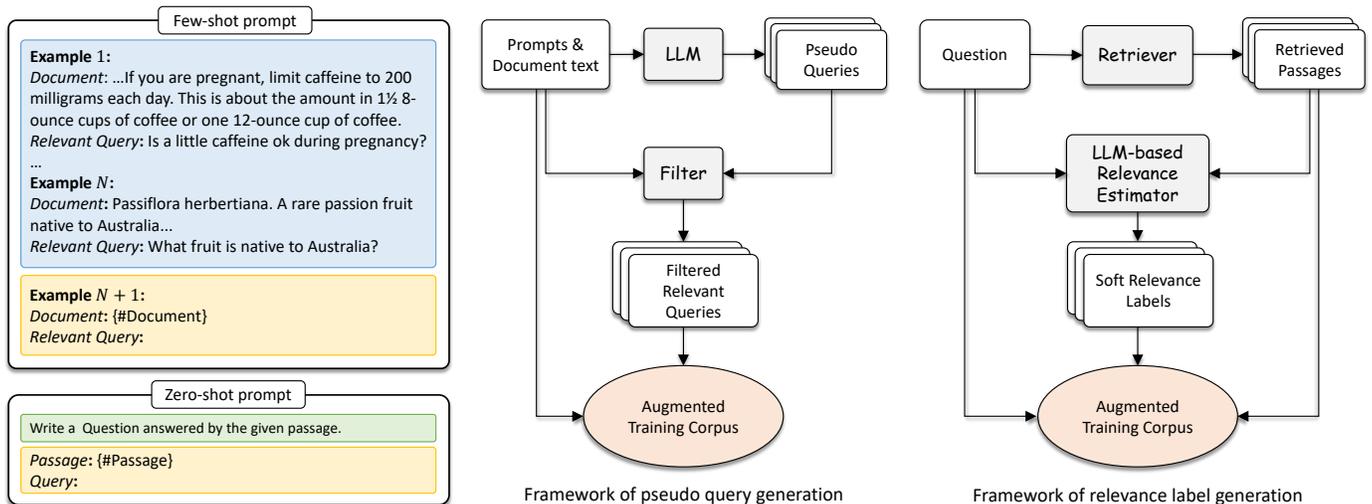


Fig. 5. Two typical frameworks for LLM-based data augmentation in the retrieval task (right), along with their prompt examples (left). Note that the methods of relevance label generation do not treat questions as inputs but regard their generation probabilities conditioned on the retrieved passages as soft relevance labels.

TABLE 3. The comparison of existing data augmentation methods powered by LLMs for training retrieval models.

Methods	# Examples	Generator	Synthetic Data	Filter Method	LLMs' tuning
InPairs [120]	3	Curie	Relevant query	Generation probability	Fixed
Ma et al. [121]	0-2	Alpaca-LLaMA & tk-Instruct	Relevant query	-	Fixed
InPairs-v2 [122]	3	GPT-J	Relevant query	Relevance score from fine-tuned monoT5-3B	Fixed
PROMPTAGATOR [123]	0-8	FLAN	Relevant query	Round-trip filtering	Fixed
TQGen [124]	0	T0	Relevant query	Generation probability	Fixed
UDAPDR [125]	0-3	GPT3 & FLAN-T5-XXL	Relevant query	Round-trip filtering	Fixed
SPTAR [126]	1-2	LLaMA-7B & Vicuna-7B	Relevant query	BM25 filtering	Soft Prompt tuning
ART [127]	0	T5-XL & T5-XXL	Soft relevance labels	-	Fixed

and few-shot learning models to address this problem [128]. A common practice to improve the models' effectiveness in a target domain without adequate label signals is through data augmentation.

How to apply LLMs for data augmentation? In the scenario of IR, it is easy to collect numerous documents. However, the challenging and costly task lies in gathering real user queries and labeling the relevant documents accordingly. Considering the strong text generation capability of LLMs, many researchers [120, 122] suggest using LLM-driven processes to create pseudo queries or relevance labels based on existing collections. These approaches facilitate the construction of relevant query-document pairs, enlarging the training data for retrieval models. According to the type of generated data, there are two mainstream approaches that complement the LLM-based data augmentation for retrieval models, *i.e.*, pseudo query generation and relevance label generation. Their frameworks are visualized in Figure 5. Next, we will give an overview of the related studies.

- *Pseudo query generation.* Given the abundance of documents, a straightforward idea is to use LLMs for generating their corresponding pseudo queries. One such illustration is presented by inPairs [120], which leverages the in-context learning capability of GPT-3. This method employs a collection of query-document pairs as demonstrations. These pairs are combined with a document and presented as input

to GPT-3, which subsequently generates possible relevant queries for the given document. By combining the same demonstration with various documents, it is easy to create a vast pool of synthetic training samples and support the fine-tuning of retrievers on specific target domains. Recent studies [121] have also leveraged open-sourced LLMs, such as Alpaca-LLaMA and tk-Instruct, to produce sufficient pseudo queries and applied curriculum learning to pre-train dense retrievers. To enhance the reliability of these synthetic samples, a fine-tuned model (*e.g.*, a monoT5-3B model fine-tuned on MSMARCO [122]) is employed to filter the generated queries. Only the top pairs with the highest estimated relevance scores are kept for training. This “generating-then-filtering” paradigm can be conducted iteratively in a round-trip filtering manner, *i.e.*, by first fine-tuning a retriever on the generated samples and then filtering the generated samples using this retriever. Repeating these EM-like steps until convergence can produce high-quality training sets [123]. Furthermore, by adjusting the prompt given to LLMs, they can generate queries of different types. This capability allows for a more accurate simulation of real queries with various patterns [124].

In practice, it is costly to generate a substantial number of pseudo queries through LLMs. Balancing the generation costs and the quality of generated samples has become an urgent problem. To tackle this, UDAPDR [125] is proposed, which first produces a limited set of synthetic queries using

LLMs for the target domain. These high-quality examples are subsequently used as prompts for a smaller model to generate a large number of queries, thereby constructing the training set for that specific domain. It is worth noting that the aforementioned studies primarily rely on fixed LLMs with frozen parameters. Empirically, optimizing LLMs' parameters can significantly improve their performance on downstream tasks. Unfortunately, this pursuit is impeded by the prohibitively high demand for computational resources. To overcome this obstacle, SPTAR [126] introduces a soft prompt tuning technique that only optimizes the prompts' embedding layer during the training process. This approach allows LLMs to better adapt to the task of generating pseudo-queries, striking a favorable balance between training cost and generation quality.

In addition to the above studies, pseudo query generation methods are also introduced in other application scenarios, such as conversational dense retrieval [105] and multilingual dense retrieval [129].

- *Relevance label generation.* In some downstream tasks of retrieval, such as question-answering, the collection of questions is also sufficient. However, the relevance labels connecting these questions with the passages of supporting evidence are very limited. In this context, leveraging the capability of LLMs for relevance label generation is a promising approach that can augment the training corpus for retrievers. A recent method, ART [127], exemplifies this approach. It first retrieves the top-relevant passages for each question. Then, it employs an LLM to produce the generation probabilities of the question conditioned on these top passages. After a normalization process, these probabilities serve as soft relevance labels for the training of the retriever.

Additionally, to highlight the similarities and differences among the corresponding methods, we present a comparative result in Table 3. It compares the aforementioned methods from various perspectives, including the number of examples, the generator employed, the type of synthetic data produced, the method applied to filter synthetic data, and whether LLMs are fine-tuned. This table serves to facilitate a clearer understanding of the landscape of these methods.

4.2 Employing LLMs to Enhance Model Architecture

Leveraging the excellent text encoding and decoding capabilities of LLMs, it is feasible to understand queries and documents with greater precision compared to earlier smaller-sized models [59]. Researchers have endeavored to utilize LLMs as the foundation for constructing advanced retrieval models. These methods can be grouped into two categories, *i.e.*, dense retrievers and generative retrievers.

4.2.1 Dense Retriever

In addition to the quantity and quality of the data, the representative capability of models also greatly influences the efficacy of retrievers. Inspired by the LLM's excellent capability to encode and comprehend natural language, some researchers [130–132] leverage LLMs as retrieval encoders and investigate the impact of model scale on retriever performance.

General Retriever. Since the effectiveness of retrievers primarily relies on the capability of text embedding, the evolution of text embedding models often has a significant impact on the progress of retriever development. In the era of LLMs, a pioneer work is made by OpenAI [130]. They view the adjacent text segments as positive pairs to facilitate the unsupervised pre-training of a set of text embedding models, denoted as *cpt-text*, whose parameter values vary from 300M to 175B. Experiments conducted on the MS MARCO [111] and BEIR [128] datasets indicate that larger model scales have the potential to yield improved performance in unsupervised learning and transfer learning for text search tasks. Nevertheless, pre-training LLMs from scratch is prohibitively expensive for most researchers. To overcome this limitation, some studies [131, 133] use pre-trained LLMs to initialize the bi-encoder of dense retriever. Specifically, GTR [133] adopts T5-family models, including T5-base, Large, XL, and XXL, to initialize and fine-tune dense retrievers. RepLLaMA [131] further fine-tunes the LLaMA model on multiple stages of IR, including retrieval and reranking. For the dense retrieval task, RepLLaMA appends an end-of-sequence token "`</s>`" to the input sequences, *i.e.*, queries or documents, and regards its output embeddings as the representation of queries or documents. The experiments confirm again that larger model sizes can lead to better performance, particularly in zero-shot settings. Notably, the researchers of RepLLaMA [131] also study the effectiveness of applying LLaMA in the reranking stage, which will be introduced in Section 5.1.3.

Task-aware Retriever. While the aforementioned studies primarily focus on using LLMs as text embedding models for downstream retrieval tasks, retrieval performance can be greatly enhanced when task-specific instructions are integrated. For example, TART [132] devises a task-aware retrieval model that introduces a task-specific instruction before the question. This instruction includes descriptions of the task's intent, domain, and desired retrieved unit. For instance, given that the task is question-answering, an effective prompt might be "Retrieve a Wikipedia text that answers this question. {question}". Here, "Wikipedia" (domain) indicates the expected source of retrieved documents, "text" (unit) suggests the type of content to retrieve, and "answers this question" (intent) demonstrates the intended relationship between the retrieved texts and the question. This approach can take advantage of the powerful language modeling capability and extensive knowledge of LLMs to precisely capture the user's search intents across various retrieval tasks. Considering the efficiency of retrievers, it first fine-tunes a TART-full model with cross-encoder architecture, which is initialized from LLMs (*e.g.*, T0-3B, Flan-T5). Then, a TART-dull model initialized from Contriever [134] is learned by distilling knowledge from the TART-full.

4.2.2 Generative Retriever

Traditional IR systems typically follow the "index-retrieval-rank" paradigm to locate relevant documents based on user queries, which has proven effective in practice. However, these systems usually consist of three separate modules: the index module, the retrieval module, and the reranking module. Therefore, optimizing these modules collectively

can be challenging, potentially resulting in sub-optimal retrieval outcomes. Additionally, this paradigm demands additional space for storing pre-built indexes, further burdening storage resources. Recently, model-based generative retrieval methods [135–137] have emerged to address these challenges. These methods move away from the traditional “index-retrieval-rank” paradigm and instead use a unified model to directly generate document identifiers (*i.e.*, DocIDs) relevant to the queries. In these model-based generative retrieval methods, the knowledge of the document corpus is stored in the model parameters, eliminating the need for additional storage space for the index. Existing methods have explored generating document identifiers through fine-tuning and prompting of LLMs [138, 139]

Fine-tuning LLMs. Given the vast amount of world knowledge contained in LLMs, it is intuitive to leverage them for building model-based generative retrievers. DSI [138] is a typical method that fine-tunes the pre-trained T5 models on retrieval datasets. The approach involves encoding queries and decoding document identifiers directly to perform retrieval. They explore multiple techniques for generating document identifiers and find that constructing semantically structured identifiers yields optimal results. In this strategy, DSI applies hierarchical clustering to group documents according to their semantic embeddings and assigns a semantic DocID to each document based on its hierarchical group. To ensure the output DocIDs are valid and do represent actual documents in the corpus, DSI constructs a trie using all DocIDs and utilizes a constraint beam search during the decoding process. Furthermore, this approach observes that the scaling law, which suggests that larger LMs lead to improved performance, is also applied to generative retrievers.

Prompting LLMs. In addition to fine-tuning LLMs for retrieval, it has been found that LLMs (*e.g.*, GPT-series models) can directly generate relevant web URLs for user queries with a few in-context demonstrations [139]. This unique capability of LLMs is believed to arise from their training exposure to various HTML resources. As a result, LLMs can naturally serve as generative retrievers that directly generate document identifiers to retrieve relevant documents for input queries. To achieve this, an LLM-URL [139] model is proposed. It utilizes the GPT-3 *text-davinci-003* model to yield candidate URLs. Furthermore, it designs regular expressions to extract valid URLs from these candidates to locate the retrieved documents.

To provide a comprehensive understanding of this topic, Table 4 summarizes the common and unique characteristics of the LLM-based retrievers discussed above.

4.3 Limitations

Though some efforts have been made for LLM-augmented retrieval, there are still many areas that require more detailed investigation. For example, a critical requirement for retrievers is fast response, while the main problem of existing LLMs is the huge model parameters and overlong inference time. Addressing this limitation of LLMs to ensure the response time of retrievers is a critical task. Moreover, even when employing LLMs to augment datasets (a context

TABLE 4. The comparison of retrievers that leverage LLMs as the foundation. “KD” is short for “Knowledge Distillation”.

Methods	Backbone	Architecture	LLM’s tuning
cpt-text [130]	GPT-series	Dense	Pre-training Fine-tuning
GTR [133]	T5	Dense	Pre-training & Fine-tuning
RepLLaMA [131]	LLAMA	Dense	Fine-tuning
TART-full [132]	T0 & Flan-T5	Dense	Fine-tuning & Prompting
TART-dual [132]	Contriever	Dense	KD & Prompting
DSI [138]	T5	Generative	Fine-tuning
LLM-URL [139]	GPT-3	Generative	Prompting

TABLE 5. Summary of existing LLM-based re-ranking methods. “Enc” and “Dec” denote encoder and decoder, respectively.

Paradigm	Type	Method
Supervised Rerankers	Enc-only	[140]
	Enc-dec	[13], [141], [142], [143]
	Dec-only	[131], [144], [145]
Unsupervised Rerankers	Pointwise	[146], [147], [148], [149], [150], [151]
	Listwise	[152], [153], [154]
	Pairwise	[155], [156]
Data Augmentation	-	[157], [158], [159], [160], [161], [162]

with lower inference time demands), the potential mismatch between LLM-generated texts and real user queries could impact retrieval effectiveness. Furthermore, as LLMs usually lack domain-specific knowledge, they need to be fine-tuned on task-specific datasets before applying them to downstream tasks. Therefore, developing efficient strategies to fine-tune these LLMs with numerous parameters emerges as a key concern.

5 RERANKER

Reranker, as the second-pass document filter in IR, aims to rerank a document list retrieved by the retriever (*e.g.*, BM25) based on the query-document relevance. Based on the usage of LLMs, the existing LLM-based reranking methods can be divided into three paradigms: utilizing LLMs as supervised rerankers, utilizing LLMs as unsupervised rerankers, and utilizing LLMs for training data augmentation. These paradigms are summarized in Table 5 and will be elaborated upon in the following sections. Recall that we will use the term *document* to refer to the text retrieved in general IR scenarios, including instances such as passages (*e.g.*, passages in MS MARCO passage ranking dataset [111]).

5.1 Utilizing LLMs as Supervised Rerankers

Supervised fine-tuning is an important step in applying pre-trained LLMs to a reranking task. Due to the lack of awareness of ranking during pre-training, LLMs cannot appropriately measure the query-document relevance and fully understand the reranking tasks. By fine-tuning LLMs on task-specific ranking datasets, such as the MS MARCO passage ranking dataset [111], which includes signals of

both relevance and irrelevance, LLMs can adjust their parameters to yield better performance in the reranking tasks. Based on the backbone model structure, we can categorize existing supervised rerankers as: (1) encoder-only, (2) encoder-decoder, and (3) decoder-only.

5.1.1 Encoder-only

The encoder-based rerankers represent a significant turning point in applying LLMs to document ranking tasks. They demonstrate how some pre-trained language models (e.g., BERT [59]) can be finetuned to deliver highly accurate relevance predictions. A representative approach is monoBERT [140], which transforms a query-document pair into a sequence “[CLS] *query* [SEP] *document* [SEP]” as the model input and calculates the relevance score by feeding the “[CLS]” representation into a linear layer. The reranking model is optimized based on the cross-entropy loss.

5.1.2 Encoder-Decoder

In this field, existing studies mainly formulate document ranking as a generation task and optimize an encoder-decoder-based reranking model [13, 141–143]. Specifically, given the query and the document, reranking models are usually fine-tuned to generate a single token, such as “true” or “false”. During inference, the query-document relevance score is determined based on the logit of the generated token. For example, a T5 model can be fine-tuned to generate classification tokens for relevant or irrelevant query-document pairs [13]. At inference time, a softmax function is applied to the logits of “true” and “false” tokens, and the relevance score is calculated as the probability of the “true” token. The following method [141] involves a multi-view learning approach based on the T5 model. This approach simultaneously considers two tasks: generating classification tokens for a given query-document pair and generating the corresponding query conditioned on the provided document. DuoT5 [142] considers a triple (q, d_i, d_j) as the input of the T5 model and is fine-tuned to generate token “true” if document d_i is more relevant to query q_i than document d_j , and “false” otherwise. During inference, for each document d_i , it enumerates all other documents d_j and uses global aggregation functions to generate the relevance score s_i for document d_i (e.g., $s_i = \sum_j p_{i,j}$, where $p_{i,j}$ represents the probability of generating “true” when taking (q, d_i, d_j) as the model input).

Although these generative loss-based methods outperform several strong ranking baselines, they are not optimal for reranking tasks. This stems from two primary reasons. First, it is commonly expected that a reranking model will yield a numerical relevance score for each query-document pair rather than text tokens. Second, compared to generation losses, it is more reasonable to optimize the reranking model using ranking losses (e.g., RankNet [163]). Recently, RankT5 [143] has directly calculated the relevance score for a query-document pair and optimized the ranking performance with “pairwise” or “listwise” ranking losses. An avenue for potential performance enhancement lies in the substitution of the base-sized T5 model with its larger-scale counterpart.

5.1.3 Decoder-only

Recently, there have been some attempts [131, 144, 145] to rerank documents by fine-tuning decoder-only models (such as LLaMA). For example, RankLLaMA [131] proposes formatting the query-document pair into a prompt “query: {*query*} document: {*document*} [EOS]” and utilizes the last token representation for relevance calculation. Besides, RankingGPT [144] has been proposed to bridge the gap between LLMs’ conventional training objectives and the specific needs of document ranking through two-stage training. The first stage involves continuously pretraining LLMs using a large number of relevant text pairs collected from web resources, helping the LLMs to naturally generate queries relevant to the input document. The second stage focuses on improving the model’s text ranking performance using high-quality supervised data and well-designed loss functions. Different from these pointwise rerankers [131, 144], Rank-without-GPT [145] proposes to train a listwise reranker that directly outputs a reranked document list. The authors first demonstrate that existing pointwise datasets (such as MS MARCO [111]), which only contain binary query-document labels, are insufficient for training efficient listwise rerankers. Then, they propose to use the ranking results of existing ranking systems (such as Cohere rerank API) as gold rankings to train a listwise reranker based on Code-LLaMA-Instruct.

5.2 Utilizing LLMs as Unsupervised Rerankers

As the size of LLMs scales up (e.g., exceeding 10 billion parameters), it becomes increasingly difficult to fine-tune the reranking model. Addressing this challenge, recent efforts have attempted to prompt LLMs to directly enhance document reranking in an unsupervised way. In general, these prompting strategies can be divided into three categories: pointwise, listwise, and pairwise methods. A comprehensive exploration of these strategies follows in the subsequent sections.

5.2.1 Pointwise methods

The pointwise methods measure the relevance between a query and a single document, and can be categorized into two types: relevance generation [146, 147] and query generation [148–150].

The upper part in Figure 6 (a) shows an example of relevance generation based on a given prompt, where LLMs output a binary label (“Yes” or “No”) based on whether the document is relevant to the query. Following [13], the query-document relevance score $f(q, d)$ can be calculated based on the log-likelihood of token “Yes” and “No” with a softmax function:

$$f(q, d) = \frac{\exp(S_Y)}{\exp(S_Y) + \exp(S_N)}, \quad (1)$$

where S_Y and S_N represent the LLM’s log-likelihood scores of “Yes” and “No” respectively. In addition to binary labels, Zhuang et al. [147] propose to incorporate fine-grained relevance labels (e.g., “highly relevant”, “somewhat relevant” and “not relevant”) into the prompt, which helps LLMs more effectively differentiate among documents with varying levels of relevance to a query.

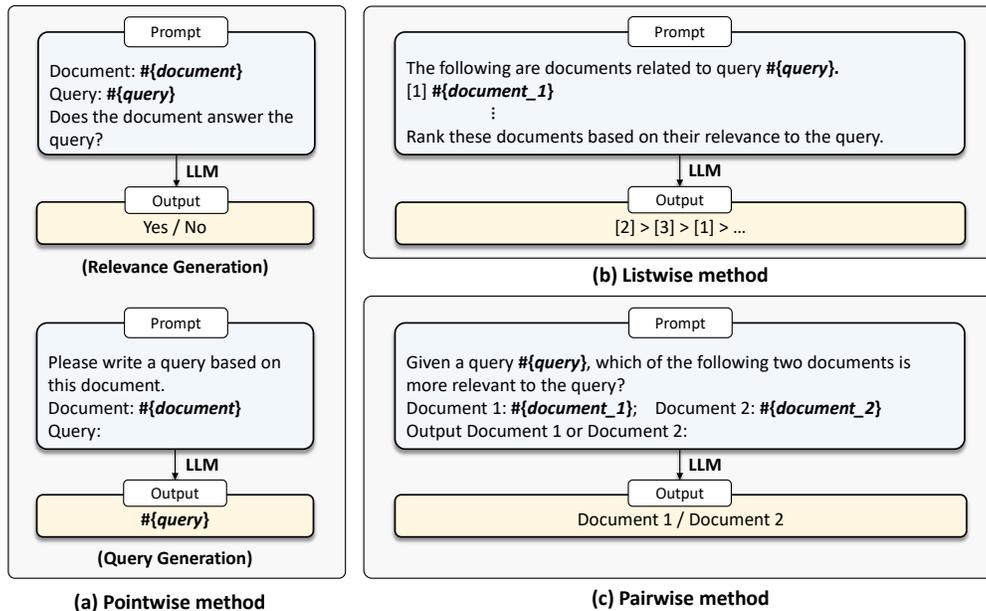


Fig. 6. Three types of unsupervised reranking methods: (a) pointwise methods that consist of relevance generation (upper) and query generation (lower), (b) listwise methods, and (c) pairwise methods.

As for the query generation shown in the lower part of Figure 6 (a), the query-document relevance score is determined by the average log-likelihood of generating the actual query tokens based on the document:

$$\text{score} = \frac{1}{|q|} \sum_i \log p(q_i | q_{<i}, d, \mathcal{P}), \quad (2)$$

where $|q|$ denotes the token number of query q , d denotes the document, and \mathcal{P} represents the provided prompt. The documents are then reranked based on their relevance scores. It has been proven that some LLMs (such as T0) yield significant performance in zero-shot document reranking based on the query generation method [148]. Recently, research [149] has also shown that the LLMs that are pre-trained without any supervised instruction fine-tuning (such as LLaMA) also yield robust zero-shot ranking ability.

Although effective, these methods primarily rely on a handcrafted prompt (e.g., “Please write a query based on this document”), which may not be optimal. As prompt is a key factor in instructing LLMs to perform various NLP tasks, it is important to optimize prompt for better performance. Along this line, a discrete prompt optimization method Co-Prompt [150] is proposed for better prompt generation in reranking tasks. Besides, PaRaDe [151] proposes a difficulty-based method to select few-shot demonstrations to include in the prompt, proving significant improvements compared with zero-shot prompts.

Note that these pointwise methods rely on accessing the output logits of LLMs to calculate the query-document relevance scores. As a result, they are not applicable to closed-sourced LLMs, whose API-returned results do not include logits.

5.2.2 Listwise Methods

Listwise methods [152, 153] aim to directly rank a list of documents (see Figure 6 (b)). These methods insert the

query and a document list into the prompt and instruct the LLMs to output the reranked document identifiers. Due to the limited input length of LLMs, it is not feasible to insert all candidate documents into the prompt. To alleviate this issue, these methods employ a sliding window strategy to rerank a subset of candidate documents each time. This strategy involves ranking from back to front using a sliding window, re-ranking only the documents within the window at a time.

Although listwise methods have yielded promising performance, they still suffer from some weaknesses. First, according to the experimental results [152], only the GPT-4-based method can achieve competitive performance. When using smaller parameterized language models (e.g., FLAN-UL2 with 20B parameters), listwise methods may produce very few usable results and underperform many supervised methods. Second, the performance of listwise methods is highly sensitive to the document order in the prompt. When the document order is randomly shuffled, listwise methods perform even worse than BM25 [152], revealing positional bias issues in the listwise ranking of LLMs. To alleviate this issue, Tang et al. [154] introduce a permutation self-consistency method, which involves shuffling the list in the prompt and aggregating the generated results to achieve a more accurate and unbiased ranking.

5.2.3 Pairwise Methods

In pairwise methods [155], LLMs are given a prompt that consists of a query and a document pair (see Figure 6 (c)). Then, they are instructed to generate the identifier of the document with higher relevance. To rerank all candidate documents, aggregation methods like AllPairs are used. AllPairs first generates all possible document pairs and aggregates a final relevance score for each document. To speed up the ranking process, efficient sorting algorithms, such as heap sort and bubble sort, are usually employed [155].

TABLE 6. The comparison between different methods. N denotes the number of documents to rerank. The Complexity, Logits, and Batch represent the computational complexity, whether accesses LLM’s logits, and whether allows batch inference respectively. k is the constant in sliding windows strategy. As for the Performance, we use NDCG@10 as a metric, and the results are calculated by reranking the top 100 documents retrieved by BM25 on TREC-DL2019 and TREC-DL2020. The best model is in bold while the second-best is marked with an underline. The results come from previous study [155]. *Since the parameters of ChatGPT have not been released, its model parameters are based on public estimates [164].

	Methods	LLM	Size	Properties			Performance	
				Complexity	Logits	Batching	TREC-DL19	TREC-DL20
Initial Retriever	BM25	-	-	-	-	-	50.58	47.96
Supervised	monoBERT [140]	BERT	340M	-	✓	✓	70.50	67.28
	monoT5 [13]	T5	220M	-	✓	✓	71.48	66.99
	RankT5 [143]	T5	3B	-	✓	✓	71.22	69.49
Unsupervised-Pointwise	Query Generation [148]	FLAN-UL2	20B	$O(N)$	✓	✓	58.95	60.02
	Relevance Generation [146]	FLAN-UL2	20B	$O(N)$	✓	✓	64.61	65.39
Unsupervised-Listwise	RankGPT _{3.5} [152]	gpt-3.5-turbo	154B*	$O(k * N)$			65.80	62.91
	RankGPT ₄ [152]	gpt-4	1T*	$O(k * N)$			75.59	<u>70.56</u>
Unsupervised-Pairwise	PRP-Allpair [155]	FLAN-UL2	20B	$O(N^2)$	✓	✓	<u>72.42</u>	70.68
	PRP-Heapsort [155]	FLAN-UL2	20B	$O(N * \log N)$	✓		71.88	69.43

These sorting algorithms utilize efficient data structures to compare document pairs selectively and elevate the most relevant documents to the top of the ranking list, which is particularly useful in top- k ranking. Experimental results show the state-of-the-art performance on the standard benchmarks using moderate-size LLMs (*e.g.*, Flan-UL2 with 20B parameters), which are much smaller than those typically employed in listwise methods (*e.g.*, GPT3.5).

Although effective, pairwise methods still suffer from high time complexity. To alleviate the efficiency problem, a setwise approach [156] has been proposed to compare a set of documents at a time and select the most relevant one from them. This approach allows the sorting algorithms (such as heap sort) to compare more than two documents at each step, thereby reducing the total number of comparisons and speeding up the sorting process.

5.2.4 Comparison and Discussion

In this part, we will compare different unsupervised methods from various aspects to better illustrate the strengths and weaknesses of each method, which is summarized in Table 6. We choose representative methods [146, 148, 152, 155] in pointwise, listwise and pairwise ranking, and include several supervised methods [13, 140, 143] mentioned in Section 5.1 for performance comparison.

The pointwise methods (Query Generation and Relevance Generation) judge the relevance of each query-document pair independently, thus offering lower time complexity and enabling batch inference. However, compared to other methods, it does not have an advantage in terms of performance. The listwise method yields significant performance especially when calling GPT-4, but suffers from expensive API cost and non-reproducibility [160]. Compared with the listwise method, the pairwise method shows competitive results based on a much smaller model FLAN-UL2 (20B). Stemming from the necessity to compare an extensive number of document pairs, its primary drawback is low efficiency.

5.3 Utilizing LLMs for Training Data Augmentation

Furthermore, in the realm of reranking, researchers have explored the integration of LLMs for training data augmentation [157–162]. For example, ExaRanker [157] generates explanations for retrieval datasets using GPT-3.5, and subsequently trains a seq2seq ranking model to generate relevance labels along with corresponding explanations for given query-document pairs. InPars-Light [158] is proposed as a cost-effective method to synthesize queries for documents by prompting LLMs. Contrary to InPars-Light [158], a new dataset ChatGPT-RetrievalQA [159] is constructed by generating synthetic documents based on LLMs in response to user queries.

Recently, many studies [160–162] have also attempted to distill the document ranking capability of LLMs into a specialized model. RankVicuna [160] proposes to use the ranking list of RankGPT_{3.5} [152] as the gold list to train a 7B parameter Vicuna model. RankZephyr [161] introduces a two-stage training strategy for distillation: initially applying the RankVicuna recipe to train Zephyr γ in the first stage, and then further finetuning it in the second stage with the ranking results from RankGPT₄. These two studies not only demonstrate competitive results but also alleviate the issue of ranking results non-reproducibility of black-box LLMs. Besides, researchers [162] have also tried to distill the ranking ability of a pairwise ranker, which is computationally demanding, into a simpler but more efficient pointwise ranker.

5.4 Limitations

Although recent research on utilizing LLMs for document reranking has made significant progress, it still faces some challenges. For example, considering the cost and efficiency, minimizing the number of calls to LLM APIs is a problem worth studying. Besides, while existing studies mainly focus on applying LLMs to open-domain datasets (such as MS-MARCO [111]) or relevance-based text ranking tasks, their adaptability to in-domain datasets [128] and non-standard ranking datasets [165] remains an area that demands more comprehensive exploration.

6 READER

With the impressive capabilities of LLMs in understanding, extracting, and processing textual data, researchers explore expanding the scope of IR systems beyond content ranking to answer generation. In this evolution, a reader module has been introduced to generate answers based on the document corpus in IR systems. By integrating a reader module, IR systems can directly present conclusive passages to users. Compared with providing a list of documents, users can simply comprehend the answering passages instead of analyzing the ranking list in this new paradigm. Furthermore, by repeatedly providing documents to LLMs based on their generating texts, the final generated answers can potentially be more accurate and information-rich than the original retrieved lists.

A naive strategy for implementing this function is to heuristically provide LLMs with documents relevant to the user queries or the previously generated texts to support the following generation. However, this passive approach limits LLMs to merely collecting documents from IR systems without active engagement. An alternative solution is to train LLMs to interact proactively with search engines. For example, LLMs can formulate their own queries instead of relying solely on user queries or generated texts for references. According to the way LLMs utilize IR systems in the reader module, we can categorize them into *passive readers* and *active readers*. Each approach has its advantages and challenges for implementing LLM-powered answer generation in IR systems. Furthermore, since the documents provided by upstream IR systems are sometimes too long to directly feed as input for LLMs, some compression modules are proposed to extractively or abtractively compress the retrieved contexts for LLMs to understand and generate answers for queries. We will present these reader and compressor modules in the following parts and briefly introduce the existing analysis work on retrieval-augmented generation strategy and their applications.

6.1 Passive Reader

To generate answers for users, a straightforward strategy is to supply the retrieved documents according to the queries or previously generated texts from IR systems as inputs to LLMs for creating passages [23, 166–171, 173, 175, 176, 178–180]. By this means, these approaches use the LLMs and IR systems separately, with LLMs functioning as passive recipients of documents from the IR systems. The strategies for utilizing LLMs within IR systems’ reader modules can be categorized into the following three groups according to the frequency of retrieving documents for LLMs.

6.1.1 Once-Retrieval Reader

To obtain useful references for LLMs to generate responses for user queries, an intuitive way is to retrieve the top documents based on the queries themselves in the beginning. For example, REALM [166] adopts this strategy by directly attending the document contents to the original queries to predict the final answers based on masked language modeling. RAG [167] follows this strategy but applies the generative language modeling paradigm. However, these two approaches only use language models with limited

parameters, such as BERT and BART. Recent approaches such as REPLUG [168] and Atlas [169] have improved them by leveraging LLMs such as GPTs, T5s, and LLaMAs for response generation. To yield better answer generation performances, these models usually fine-tune LLMs on QA tasks. However, due to the limited computing resources, many methods [170, 171, 179] choose to prompt LLMs for generation as they could use larger LMs in this way. Furthermore, to improve the quality of the generated answers, several approaches [172, 181] also try to train or prompt the LLMs to generate contexts such as citations or notes in addition to the answers to force LLMs to understand and assess the relevance of retrieved passages to the user queries. Some approaches [180] evaluate the importance of each retrieved reference using policy gradients to indicate which reference is more useful for generating. Besides, researchers explore instruction tuning LLMs such LLaMAs to improve their abilities to generate conclusive passages relying on retrieved knowledge [182, 183].

6.1.2 Periodic-Retrieval Reader

However, while generating long conclusive answers, it is shown [23, 173] that only using the references retrieved by the original user intents as in once-retrieval readers may be inadequate. For example, when providing a passage about “Barack Obama”, language models may need additional knowledge about his university, which may not be included in the results of simply searching the initial query. In conclusion, language models may need extra references to support the following generation during the generating process, where multiple retrieval processes may be required. To address this, solutions such as RETRO [23] and RALM [173] have emerged, emphasizing the periodic collection of documents based on both the original queries and the concurrently generated texts (triggering a retrieval every n generated tokens). In this manner, when generating the text about the university career of Barack Obama, the LLM can receive additional documents as supplementary materials. This need for additional references highlights the necessity for multiple retrieval iterations to ensure robustness in subsequent answer generation. Notably, RETRO [23] introduces a novel approach incorporating cross-attention between the generating texts and the references within the Transformer attention calculation, as opposed to directly embedding references into the input texts of LLMs. Since it involves additional cross-attention modules in the Transformer’s structure, RETRO trains this model from scratch. However, these two approaches mainly rely on the successive n tokens to separate generation and retrieve documents, which may not be semantically continuous and may cause the collected references noisy and useless. To solve this problem, some approaches such as IRCoT [175] also explore retrieving documents for every generated sentence, which is a more complete semantic structure. Furthermore, researchers find that the whole generated passages can be considered as conclusive contexts for current queries and can be used to find more relevant knowledge to generate more thorough answers. Consequently, many recent approaches [174, 184, 185] have also tried to extend this periodic-retrieval paradigm to iteratively using the whole generated passages to retrieve references to re-generate the

TABLE 7. The comparison of existing representative methods that have a passive reader module. REALM and RAG do not use LLMs, but their frameworks have been widely applied in many following approaches.

Methods	Backbone models	Where to incorporate retrieval	When to retrieve	How to use LLMs
REALM [166]	BERT	Input layer	In the beginning	Fine-tuning
RAG [167]	BART	Input layer	In the beginning	Fine-tuning
REPLUG [168]	GPT	Input layer	In the beginning	Fine-tuning
Atlas [169]	T5	Input layer	In the beginning	Fine-tuning
Lazaridou et al. [170]	Gopher	Input layer	In the beginning	Prompting
He et al. [171]	GPT	Input layer	In the beginning	Prompting
Chain-of-Note [172]	LLaMA	Input layer	In the beginning	Fine-tuning
RALM [173]	LLaMA & OPT & GPT	Input layer	During generation (every n tokens)	Prompting
RETRO [23]	Transformer	Attention layer	During generation (every n tokens)	Training from scratch
ITERGEN [174]	GPT	Input layer	During generation (every answer)	Prompting
IRCoT [175]	Flan-T5 & GPT	Input layer	During generation (every sentence)	Prompting
FLARE [176]	GPT	Input layer	During generation (aperiodic)	Prompting
Self-RAG [177]	LLaMA	Input layer	During generation (aperiodic)	Fine-tuning

answers, until the iterations reach a pre-defined limitation. Particularly, these methods can be regarded as special periodic-retrieval readers that retrieve passages when every answer is (re)-generated. Since the LLMs can receive more comprehensive and relevant references with the iterations increase, these methods that combine retrieval-augmented-generation and generation-augmented-retrieval strategies can generate more accurate answers but consume more computation costs.

6.1.3 Aperiodic-Retrieval Reader

In the above strategy, the retrieval systems supply documents to LLMs in a periodic manner. However, retrieving documents in a mandatory frequency may mismatch the retrieval timing and can be costly. Recently, FLARE [176] has addressed this problem by automatically determining the timing of retrieval according to the probability of generating texts. Since the probability can serve as an indicator of LLMs’ confidence during text generation [186, 187], a low probability for a generated term could suggest that LLMs require additional knowledge. Specifically, when the probability of a term falls below a predefined threshold, FLARE employs IR systems to retrieve references in accordance with the ongoing generated sentences, while removing these low-probability terms. FLARE adopts this strategy of prompting LLMs for answer generation solely based on the probabilities of generating terms, avoiding the need for fine-tuning while still maintaining effectiveness. Besides, self-RAG [177] tends to solve this problem by training LLMs such as LLaMA to generate specific tokens when they need additional knowledge to support following generations. Another critical model is introduced to judge whether the retrieved references are beneficial for generating.

We summarize representative passive reader approaches in Table 7, considering various aspects such as the backbone language models, the insertion point for retrieved references, the timing of using retrieval models, and the tuning strategy employed for LLMs.

6.2 Active Reader

However, the passive reader-based approaches separate IR systems and generative language models. This signifies that LLMs can only submissively utilize references provided by IR systems and are unable to interactively engage with the

IR systems in a manner akin to human interaction such as issuing queries to seek information.

To allow LLMs to actively use search engines, Self-Ask [188] and DSP [189] try to employ few-shot prompts for LLMs, triggering them to search queries when they believe it is required. For example, in a scenario where the query is “When was the existing tallest wooden lattice tower built?”, these prompted LLMs can decide to search a query “What is the existing tallest wooden lattice tower” to gather necessary references as they find the query cannot be directly answered. Once acquired information about the tower, they can iteratively query IR systems for more details until they determine to generate the final answers instead of asking questions. Notably, these methods involve IR systems to construct a single reasoning chain for LLMs. MRC [190] further improves these methods by prompting LLMs to explore multiple reasoning chains and subsequently combining all generated answers using LLMs.

6.3 Compressor

Existing LLMs, especially open-sourced ones, such as LLaMA and Flan-T5, have limited input lengths (usually 4,096 or 8,192 tokens). However, the documents or web pages retrieved by upstream IR systems are usually long. Therefore, it is difficult to concatenate all the retrieved documents and feed them into LLMs to generate answers. Though some approaches manage to solve these problems by aggregating the answers supported by each reference as the final answers, this strategy neglects the potential relations between retrieved passages. A more straightforward way is to directly compress the retrieved documents into short input tokens or even dense vectors [191–194].

To compress the retrieved references, an intuitive idea is to extract the most useful K sentences from the retrieved documents. LeanContext [191] applies this method and trains a small model by reinforcement learning (RL) to select the top K similar sentences to the queries. The researchers also augment this strategy by using a free open-sourced text reduction method for the rest sentences as a supplement. Instead of using RL-based methods, RECOMP [192] directly uses the probability or the match ratio of the generated answers to the golden answers as signals to build training datasets and tune the compressor model. For example, the sentence corresponding to the highest generating proba-

bility is the positive one while others are negative ones. Furthermore, FILCO [193] applies the “hindsight” methods, which directly align the prior distribution (the predicted importance probability distribution of sentences without knowing the gold answer) to the posterior distribution (the same distribution of sentences within knowing the gold answer) to tune language models to select sentences.

However, these extractive methods may lose potential intent among all references. Therefore, abstractive methods are proposed to summarize retrieved documents into short but concise summaries for downstream generation. These methods [192, 194] usually distill the summarizing abilities of LLMs to small models. For example, TCRA [194] leverages GPT-3.5-turbo to build abstractive compression datasets for MT5 model.

6.4 Analysis

With the rapid development of the above reader approaches, many researchers have begun to analyze the characteristics of retrieval-augmented LLMs:

- Liu et al. [195] find that the position of the relevant/golden reference has significant influences on the final generation performance. The performance is always better when the relevant reference is at the beginning or the end, which indicates the necessity of introducing a ranking module to order the retrieved knowledge.
- Ren et al. [196] observe that by applying retrieval augmentation generation strategy, LLMs can have a better awareness of their knowledge boundaries.
- Liu et al. [197] analyze different strategies of integrating retrieval systems and LLMs such as concatenate (*i.e.*, concatenating all references for answer generation) and post fusion (*i.e.*, aggregating the answers corresponding to each reference). They also explore several ways of combining these two strategies.
- Aksitov et al. [198] demonstrate that there exists an attribution and fluency tradeoff for retrieval-augmented LLMs: with more received references, the attribution of generated answers increases while the fluency decreases.
- Mallen et al. [199] argue that always retrieving references to support LLMs to generate answers hurts the question-answering performance. The reason is that LLMs themselves may have adequate knowledge while answering questions about popular entities and the retrieved noisy passages may interfere and bias the answering process. To overcome this challenge, they devise a simple strategy that only retrieves references while the popularity of entities in the query is quite low. By this means, the efficacy and efficiency of retrieval-augmented generation both improve.

6.5 Applications

Recently, researchers [200–205] have applied the retrieval-augmented generation strategy to areas such as clinical QA, medical QA, and financial QA to enhance LLMs with external knowledge and to develop domain-specific applications. For example, ATLANTIC [201] adapts Atlas to the scientific domain to derive a science QA system. Besides, some approaches [206] also apply techniques in federated learning such as multi-party computation to perform personal retrieval-augmented generation with privacy protection.

Furthermore, to better facilitate the deployment of these retrieval-augmented generation systems, some tools or frameworks are proposed [178, 207, 208]. For example, RETA-LLM [178] breaks down the whole complex generation task into several simple modules in the reader pipeline. These modules include a query rewriting module for refining query intents, a passage extraction module for aligning reference lengths with LLM limitations, and a fact verification module for confirming the absence of fabricated information in the generated answers.

6.6 Limitations

Several IR systems applying the retrieval-augmented generation strategy, such as New Bing and Langchain, have already entered commercial use. However, there are also some challenges in this novel retrieval-augmented content generation system. These include challenges such as effective query reformulation, optimal retrieval frequency, correct document comprehension, accurate passage extraction, and effective content summarization. It is crucial to address these challenges to effectively realize the potential of LLMs in this paradigm.

7 SEARCH AGENT

With the development of LLMs, IR systems are also facing new changes. Among them, developing LLMs as intelligent agents has attracted more and more attention. This conceptual shift aims to mimic human browsing patterns, thereby enhancing the capability of these models to handle complex retrieval tasks. Empowered by the advanced natural language understanding and generation capabilities of LLMs, these agents can autonomously search, interpret, and synthesize information from a wide range of sources.

One way to achieve this ability is to design a pipeline that combines a series of modules and assigns different roles to them. Such a pre-defined pipeline mimics users’ behaviors on the web by breaking it into several sub-tasks which are performed by different modules. However, this kind of static agent cannot deal with the complex nature of users’ behavior sequences on the web and may face challenges when interacting with real-world environments. An alternative solution is to allow LLMs to freely explore the web and make interactions themselves, namely letting the LLM itself decide what action it will take next based on the feedback from the environment (or humans). These agents have more flexibility and act more like human beings.

7.1 Static Agent

To mimic human search patterns, a straightforward approach is to design a static system to browse the web and synthesize information step by step [209–214]. By breaking the information-seeking process into multiple subtasks, they design a pipeline that contains various LLM-based modules in advance and assigns different subtasks to them.

LaMDA [209] serves as an early work of the static agent. It consists of a family of Transformer-based neural language models specialized for dialog, with up to 137B parameters, pre-trained on 1.56T tokens from public dialogue data and web text. The study emphasizes the model’s development

through a static pipeline, encompassing large-scale pre-training, followed by strategic fine-tuning stages aimed at enhancing three critical aspects: dialogue quality, safety, and groundedness. It can integrate external IR systems for factual grounding. This integration allows LaMDA to access and use external and authoritative sources when generating responses. SeeKeR [210] also incorporates the Internet search into its modular architecture for generating more factual responses. It performs three sequential tasks: generating a search query, generating knowledge from search results, and generating a final response. GopherCite [213] uses a search engine like Google Search to find relevant sources. It then synthesizes a response that includes verbatim quotes from these sources as evidence, aligning the Gopher’s output with verified information. WebAgent [212] develops a series of tasks, including instruction decomposition and planning, action programming, and HTML summarization. It can navigate the web, understand and synthesize information from multiple sources, and execute web-based tasks, effectively functioning as an advanced search and interaction agent. WebGLM [211] designs an LLM-augmented retriever, a bootstrapped generator, and a human preference-aware scorer. These components work together to provide accurate web-enhanced question-answering capabilities that are sensitive to human preferences. Shi et al. [214] focus on enhancing the relevance, responsibility, and trustworthiness of LLMs in web search applications via an intent-aware generator, an evidence-sensitive validator, and a multi-strategy supported optimizer.

7.2 Dynamic Agent

Instead of statically arranging LLMs in a pipeline, WebGPT [24] takes an alternate approach by training LLMs to use search engines automatically. This is achieved through the application of a reinforcement learning framework, within which a simulated environment is constructed for GPT-3 models. Specifically, the WebGPT model employs special tokens to execute actions such as querying, scrolling through rankings, and quoting references on search engines. This innovative approach allows the GPT-3 model to use search engines for text generation, enhancing the reliability and real-time capability of the generated texts. A following study [215] has extended this paradigm to the domain of Chinese question answering. Besides, some works develop important benchmarks for interactive web-based agents [216–218]. For example, WebShop [217] aims to provide a scalable, interactive web-based environment for language understanding and decision-making, focusing on the task of online shopping. ASH (Actor-Summarizer-Hierarchical) prompting [219] significantly enhances the ability of LLMs on WebShop benchmark. It first takes a raw observation from the environment and produces a new, more meaningful representation that aligns with the specific goal. Then, it dynamically predicts the next action based on the summarized observation and the interaction history.

7.3 Limitations

Though the aspect of static search agents has been thoroughly studied, the literature on dynamic search agents remains limited. Some agents may lack mechanisms for

real-time fact-checking or verification against authoritative sources, leading to the potential dissemination of misinformation. Moreover, since LLMs are trained on data from the Internet, they may inadvertently perpetuate biases present in the training data. This can lead to biased or offensive outputs and may collect unethical content from the web. Finally, as LLMs process user queries, there are concerns regarding user privacy and data security, especially if sensitive or personal information is involved in the queries.

8 FUTURE DIRECTION

In this survey, we comprehensively reviewed recent advancements in LLM-enhanced IR systems and discussed their limitations. Since the integration of LLMs into IR systems is still in its early stages, there are still many opportunities and challenges. In this section, we summarize the potential future directions in terms of the four modules in an IR system we just discussed, namely query rewriter, retriever, reranker, and reader. In addition, as evaluation has also emerged as an important aspect, we will also introduce the corresponding research problems that need to be addressed in the future. Another discussion about important research topics on applying LLMs to IR can be found in a recent perspective paper [53].

8.1 Query Rewriter

LLMs have enhanced query rewriting for both ad-hoc and conversational search scenarios. Most of the existing methods rely on prompting LLMs to generate new queries. While yielding remarkable results, the refinement of rewriting quality and the exploration of potential application scenarios require further investigation.

- *Rewriting queries according to ranking performance.* A typical paradigm of prompting-based methods is providing LLMs with several ground-truth rewriting cases (optional) and the task description of query rewriting. Despite LLMs being capable of identifying potential user intents of the query [220], they lack awareness of the resulting retrieval quality of the rewritten query. The absence of this connection can result in rewritten queries that seem correct yet produce unsatisfactory ranking results. Although some existing studies have used reinforcement learning to adjust the query rewriting process according to generation results [100], a substantial realm of research remains unexplored concerning the integration of ranking results.

- *Improving query rewriting in conversational search.* As yet, primary efforts have been made to improve query rewriting in ad-hoc search. In contrast, conversational search presents a more developed landscape with a broader scope for LLMs to contribute to query understanding. By incorporating historical interactive information, LLMs can adapt system responses based on user preferences, providing a more effective conversational experience. However, this potential has not been explored in depth. In addition, LLMs could also be used to simulate user behavior in conversational search scenarios, providing more training data, which are urgently needed in current research.

- *Achieving personalized query rewriting.* LLMs offer valuable contributions to personalized search through their capacity to analyze user-specific data. In terms of query rewriting, with the excellent language comprehension ability of

LLMs, it is possible to leverage them to build user profiles based on users' search histories (e.g., issued queries, click-through behaviors, and dwell time). This empowers the achievement of personalized query rewriting for enhanced IR and finally benefits personalized search or personalized recommendation.

8.2 Retriever

Leveraging LLMs to improve retrieval models has received considerable attention, promising an enhanced understanding of queries and documents for improved ranking performance. However, despite strides in this field, several challenges and limitations still need to be investigated in the future:

- *Reducing the latency of LLM-based retrievers.* LLMs, with their massive parameters and world knowledge, often entail high latency during the inferring process. This delay poses a significant challenge for practical applications of LLM-based retrievers, as search engines require in-time responses. To address this issue, promising research directions include transferring the capabilities of LLMs to smaller models, exploring quantization techniques for LLMs in IR tasks, and so on.

- *Simulating realistic queries for data augmentation.* Since the high latency of LLMs usually blocks their online application for retrieval tasks, many existing studies have leveraged LLMs to augment training data, which is insensitive to inference latency. Existing methods that leverage LLMs for data augmentation often generate queries without aligning them with real user queries, leading to noise in the training data and limiting the effectiveness of retrievers. As a consequence, exploring techniques such as reinforcement learning to enable LLMs to simulate the way that real queries are issued holds the potential for improving retrieval tasks.

- *Incremental indexing for generative retrieval.* As elaborated in Section 4.2.2, the emergence of LLMs has paved the way for generative retrievers to generate document identifiers for retrieval tasks. This approach encodes document indexes and knowledge into the LLM parameters. However, the static nature of LLM parameters, coupled with the expensive fine-tuning costs, poses challenges for updating document indexes in generative retrievers when new documents are added. Therefore, it is crucial to explore methods for constructing an incremental index that allows for efficient updates in LLM-based generative retrievers.

- *Supporting multi-modal search.* Web pages usually contain multi-modal information, including texts, images, audios, and videos. However, existing LLM-enhanced IR systems mainly support retrieval for text-based content. A straightforward solution is to replace the backbone with multi-modal large models, such as GPT-4 [80]. However, this undoubtedly increases the cost of deployment. A promising yet challenging direction is to combine the language understanding capability of LLMs with existing multi-modal retrieval models. By this means, LLMs can contribute their language skills in handling different types of content.

8.3 Reranker

In Section 5, we have discussed the recent advanced techniques of utilizing LLMs for the reranking task. Some potential future directions in reranking are discussed as follows.

- *Enhancing the online availability of LLMs.* Though effective, many LLMs have a massive number of parameters, making it challenging to deploy them in online applications. Besides, many reranking methods [152, 153] rely on calling LLM APIs, incurring considerable costs. Consequently, devising effective approaches (such as distilling to small models) to enhance the online applicability of LLMs emerges as a research direction worth exploring.

- *Improving personalized search.* Many existing LLM-based reranking methods mainly focus on the ad-hoc reranking task. However, by incorporating user-specific information, LLMs can also improve the effectiveness of the personalized reranking task. For example, by analyzing users' search history, LLMs can construct accurate user profiles and rerank the search results accordingly, providing personalized results with higher user satisfaction.

- *Adapting to diverse ranking tasks.* In addition to document reranking, there are also other ranking tasks, such as response ranking, evidence ranking, entity ranking and etc., which also belong to the universal information access system. Navigating LLMs towards adeptness in these diverse ranking tasks can be achieved through specialized methodologies, such as instruction tuning. Exploring this avenue holds promise as an intriguing and valuable research trajectory.

8.4 Reader

With the increasing capabilities of LLMs, the future interaction between users and IR systems will be significantly changed. Due to the powerful natural language processing and understanding capabilities of LLMs, the traditional search paradigm of providing ranking results is expected to be progressively replaced by synthesizing conclusive answering passages for user queries using the reader module. Although such strategies have already been investigated by academia and facilitated by industry as we stated in Section 6, there still exists much room for exploration.

- *Improving the reference quality for LLMs.* To support answer generation, existing approaches usually directly feed the retrieved documents to the LLMs as references. However, since a document usually covers many topics, some passages in it may be irrelevant to the user queries and can introduce noise during LLMs' generation. Therefore, it is necessary to explore techniques for extracting relevant snippets from retrieved documents, enhancing the performance of retrieval-augmented generation.

- *Improving the answer reliability of LLMs.* Incorporating the retrieved references has significantly alleviated the "hallucination" problem of LLMs. However, it remains uncertain whether the LLMs refer to these supported materials during answering queries. Some studies [196] have revealed that LLMs can still provide unfaithful answers even with additional references. Therefore, the reliability of the conclusive answers might be lower compared to the ranking results provided by traditional IR systems. It is essential to investigate the influence of these references on the generation process, thereby improving the credibility of reader-based novel IR systems.

8.5 Search Agent

With the outstanding performance of LLMs, the patterns of searching may completely change from traditional IR systems to autonomous search agents. In Section 7, we have discussed many existing works that utilize a static or dynamic pipeline to autonomously browse the web. These works are believed to be the pioneering works of the new searching paradigm. However, there is still plenty of room for further improvements.

- *Enhancing the Trustworthiness of LLMs.* When LLMs are enabled to browse the web, it is important to ensure the validity of retrieved documents. Otherwise, the unfaithful information may increase the LLMs’ “hallucination” problem. Besides, even if the gathered information has high quality, it remains unclear whether they are really used for synthesizing responses. A potential strategy to address this issue is enabling LLMs to autonomously validate the documents they scrape. This self-validation process could incorporate mechanisms for assessing the credibility and accuracy of the information within these documents.

- *Mitigating Bias and Offensive Content in LLMs.* The presence of biases and offensive content within LLM outputs is a pressing concern. This issue primarily stems from biases inherent in the training data and will be amplified by the low-quality information gathered from the web. Achieving this requires a multi-faceted approach, including improvements in training data, algorithmic adjustments, and continuous monitoring for bias and inappropriate content that LLMs collect and generate.

8.6 Evaluation

LLMs have attracted significant attention in the field of IR due to their strong ability in context understanding and text generation. To validate the effectiveness of LLM-enhanced IR approaches, it is crucial to develop appropriate evaluation metrics. Given the growing significance of readers as integral components of IR systems, the evaluation should consider two aspects: assessing ranking performance and evaluating generation performance.

- *Generation-oriented ranking evaluation.* Traditional evaluation metrics for ranking primarily focus on comparing the retrieval results of IR models with ground-truth (relevance) labels. Typical metrics include precision, recall, mean reciprocal rank (MRR) [221], mean average precision (MAP), and normalized discounted cumulative gain (nDCG) [222]. These metrics measure the alignment between ranking results and human preference on using these results. Nevertheless, these metrics may fall short in capturing a document’s role in the generation of passages or answers, as their relevance to the query alone might not adequately reflect this aspect. This effect could be leveraged as a means to evaluate the usefulness of documents more comprehensively. A formal and rigorous evaluation metric for ranking that centers on generation quality has yet to be defined.

- *Text generation evaluation.* The wide application of LLMs in IR has led to a notable enhancement in their generation capability. Consequently, there is an imperative demand for novel evaluation strategies to effectively evaluate the performance of passage or answer generation. Previous evaluation metrics for text generation have several limitations,

including: (1) Dependency on lexical matching: methods such as BLEU [223] or ROUGE [224] primarily evaluate the quality of generated outputs based on n -gram matching. This approach cannot account for lexical diversity and contextual semantics. As a result, models may favor generating common phrases or sentence structures rather than producing creative and novel content. (2) Insensitivity to subtle differences: existing evaluation methods may be insensitive to subtle differences in generated outputs. For example, if a generated output has minor semantic differences from the reference answer but is otherwise similar, traditional methods might overlook these nuanced distinctions. (3) Lack of ability to evaluate factuality: LLMs are prone to generating “hallucination” problems [225–228]. The hallucinated texts can closely resemble the oracle texts in terms of vocabulary usage, sentence structures, and patterns, while having non-factual content. Existing methods are hard to identify such problems, while the incorporation of additional knowledge sources such as knowledge bases or reference texts could potentially aid in addressing this challenge.

8.7 Bias

Since ChatGPT was released, LLMs have drawn much attention from both academia and industry. The wide applications of LLMs have led to a notable increase in content on the Internet that is not authored by humans but rather generated by these language models. However, as LLMs may hallucinate and generate non-factual texts, the increasing number of LLM-generated contents also brings worries that these contents may provide fictitious information for users across IR systems. More severely, researchers [229, 230] show that some modules in IR systems such as retriever and reranker, especially those based on neural models, may prefer LLM-generated documents, since their topics are more consistent and the perplexity of them are lower compared with human-written documents. The authors refer to this phenomenon as the “source bias” towards LLM-generated text. It is challenging but necessary to consider how to build IR systems free from this category of bias.

9 CONCLUSION

In this survey, we have conducted a thorough exploration of the transformative impact of LLMs on IR across various dimensions. We have organized existing approaches into distinct categories based on their functions: query rewriting, retrieval, reranking, and reader modules. In the domain of query rewriting, LLMs have demonstrated their effectiveness in understanding ambiguous or multi-faceted queries, enhancing the accuracy of intent identification. In the context of retrieval, LLMs have improved retrieval accuracy by enabling more nuanced matching between queries and documents, considering context as well. Within the reranking realm, LLM-enhanced models consider more fine-grained linguistic nuances when re-ordering results. The incorporation of reader modules in IR systems represents a significant step towards generating comprehensive responses instead of mere document lists. The integration of LLMs into IR systems has brought about a fundamental change in how users engage with information and knowledge. From query rewriting to retrieval, reranking, and

reader modules, LLMs have enriched each aspect of the IR process with advanced linguistic comprehension, semantic representation, and context-sensitive handling. As this field continues to progress, the journey of LLMs in IR portends a future characterized by more personalized, precise, and user-centric search encounters.

This survey focuses on reviewing recent studies of applying LLMs to different IR components and using LLMs as search agents. Beyond this, a more significant problem brought by the appearance of LLMs is: is the conventional IR framework necessary in the era of LLMs? For example, traditional IR aims to return a ranking list of documents that are relevant to issued queries. However, the development of generative language models has introduced a novel paradigm: the direct generation of answers to input questions. Furthermore, according to a recent perspective paper [53], IR might evolve into a fundamental service for diverse systems. For example, in a multi-agent simulation system [231], an IR component can be used for memory recall. This implies that there will be many new challenges in future IR.

REFERENCES

- [1] Y. Wu, W. Wu, C. Xing, M. Zhou, and Z. Li, "Sequential matching network: A new architecture for multi-turn response selection in retrieval-based chatbots," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, R. Barzilay and M. Kan, Eds. Association for Computational Linguistics, 2017, pp. 496–505.
- [2] H. Shum, X. He, and D. Li, "From eliza to xiaoice: challenges and opportunities with social chatbots," *Frontiers Inf. Technol. Electron. Eng.*, vol. 19, no. 1, pp. 10–26, 2018.
- [3] V. Karpukhin, B. Oguz, S. Min, P. S. H. Lewis, L. Wu, S. Edunov, D. Chen, and W. Yih, "Dense passage retrieval for open-domain question answering," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, B. Webber, T. Cohn, Y. He, and Y. Liu, Eds. Association for Computational Linguistics, 2020, pp. 6769–6781.
- [4] R. Datta, D. Joshi, J. Li, and J. Z. Wang, "Image retrieval: Ideas, influences, and trends of the new age," *ACM Comput. Surv.*, vol. 40, no. 2, pp. 5:1–5:60, 2008.
- [5] C. Yuan, W. Zhou, M. Li, S. Lv, F. Zhu, J. Han, and S. Hu, "Multi-hop selector network for multi-turn response selection in retrieval-based chatbots," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, K. Inui, J. Jiang, V. Ng, and X. Wan, Eds. Association for Computational Linguistics, 2019, pp. 111–120.
- [6] Y. Zhu, J. Nie, K. Zhou, P. Du, and Z. Dou, "Content selection network for document-grounded retrieval-based chatbots," in *Advances in Information Retrieval - 43rd European Conference on IR Research, ECIR 2021, Virtual Event, March 28 - April 1, 2021, Proceedings, Part I*, ser. Lecture Notes in Computer Science, D. Hiemstra, M. Moens, J. Mothe, R. Perego, M. Potthast, and F. Sebastiani, Eds., vol. 12656. Springer, 2021, pp. 755–769.
- [7] Y. Zhu, J. Nie, K. Zhou, P. Du, H. Jiang, and Z. Dou, "Proactive retrieval-based chatbots based on relevant knowledge and goals," in *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, F. Diaz, C. Shah, T. Suel, P. Castells, R. Jones, and T. Sakai, Eds. ACM, 2021, pp. 2000–2004.
- [8] H. Qian, Z. Dou, Y. Zhu, Y. Ma, and J. Wen, "Learning implicit user profiles for personalized retrieval-based chatbot," *CoRR*, vol. abs/2108.07935, 2021.
- [9] Y. Qu, Y. Ding, J. Liu, K. Liu, R. Ren, W. X. Zhao, D. Dong, H. Wu, and H. Wang, "Rocketqa: An optimized training approach to dense passage retrieval for open-domain question answering," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, K. Toutanova, A. Rumshisky, L. Zettlemoyer, D. Hakkani-Tür, I. Beltagy, S. Bethard, R. Cotterell, T. Chakraborty, and Y. Zhou, Eds. Association for Computational Linguistics, 2021, pp. 5835–5847.
- [10] Y. Arens, C. A. Knoblock, and W. Shen, "Query reformulation for dynamic information integration," *J. Intell. Inf. Syst.*, vol. 6, no. 2/3, pp. 99–130, 1996.
- [11] J. Huang and E. N. Efthimiadis, "Analyzing and evaluating query reformulation strategies in web search logs," in *Proceedings of the 18th ACM Conference on Information and Knowledge Management, CIKM 2009, Hong Kong, China, November 2-6, 2009*, D. W. Cheung, I. Song, W. W. Chu, X. Hu, and J. Lin, Eds. ACM, 2009, pp. 77–86.
- [12] R. F. Nogueira, W. Yang, K. Cho, and J. Lin, "Multi-stage document ranking with BERT," *CoRR*, vol. abs/1910.14424, 2019.
- [13] R. F. Nogueira, Z. Jiang, R. Pradeep, and J. Lin, "Document ranking with a pretrained sequence-to-sequence model," in *EMNLP (Findings)*, ser. Findings of ACL, vol. EMNLP 2020. Association for Computational Linguistics, 2020, pp. 708–718.
- [14] Y. Zhu, J. Nie, Z. Dou, Z. Ma, X. Zhang, P. Du, X. Zuo, and H. Jiang, "Contrastive learning of user behavior sequence for context-aware document ranking," in *CIKM '21: The 30th ACM International Conference on Information and Knowledge Management, Virtual Event, Queensland, Australia, November 1 - 5, 2021*, G. Demartini, G. Zuccon, J. S. Culpepper, Z. Huang, and H. Tong, Eds. ACM, 2021, pp. 2780–2791.
- [15] J. Teevan, S. T. Dumais, and E. Horvitz, "Personalizing search via automated analysis of interests and activities," in *SIGIR 2005: Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Salvador, Brazil, August 15-19, 2005*, R. A. Baeza-Yates, N. Ziviani, G. Marchionini, A. Moffat, and J. Tait, Eds. ACM, 2005, pp. 449–456.

- [16] P. N. Bennett, R. W. White, W. Chu, S. T. Dumais, P. Bailey, F. Borisyuk, and X. Cui, "Modeling the impact of short- and long-term behavior on search personalization," in *The 35th International ACM SIGIR conference on research and development in Information Retrieval, SIGIR '12, Portland, OR, USA, August 12-16, 2012*, W. R. Hersch, J. Callan, Y. Maarek, and M. Sanderson, Eds. ACM, 2012, pp. 185–194.
- [17] S. Ge, Z. Dou, Z. Jiang, J. Nie, and J. Wen, "Personalizing search results using hierarchical RNN with query-aware attention," in *Proceedings of the 27th ACM International Conference on Information and Knowledge Management, CIKM 2018, Torino, Italy, October 22-26, 2018*, A. Cuzzocrea, J. Allan, N. W. Paton, D. Srivastava, R. Agrawal, A. Z. Broder, M. J. Zaki, K. S. Candan, A. Labrinidis, A. Schuster, and H. Wang, Eds. ACM, 2018, pp. 347–356.
- [18] Y. Zhou, Z. Dou, Y. Zhu, and J. Wen, "PSSL: self-supervised learning for personalized search with contrastive sampling," in *CIKM '21: The 30th ACM International Conference on Information and Knowledge Management, Virtual Event, Queensland, Australia, November 1 - 5, 2021*, G. Demartini, G. Zuccon, J. S. Culpepper, Z. Huang, and H. Tong, Eds. ACM, 2021, pp. 2749–2758.
- [19] J. G. Carbonell and J. Goldstein, "The use of mmr, diversity-based reranking for reordering documents and producing summaries," in *SIGIR '98: Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, August 24-28 1998, Melbourne, Australia*, W. B. Croft, A. Moffat, C. J. van Rijsbergen, R. Wilkinson, and J. Zobel, Eds. ACM, 1998, pp. 335–336.
- [20] R. Agrawal, S. Gollapudi, A. Halverson, and S. Jeong, "Diversifying search results," in *Proceedings of the Second International Conference on Web Search and Web Data Mining, WSDM 2009, Barcelona, Spain, February 9-11, 2009*, R. Baeza-Yates, P. Boldi, B. A. Ribeiro-Neto, and B. B. Cambazoglu, Eds. ACM, 2009, pp. 5–14.
- [21] J. Liu, Z. Dou, X. Wang, S. Lu, and J. Wen, "DVGAN: A minimax game for search result diversification combining explicit and implicit features," in *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020*, J. X. Huang, Y. Chang, X. Cheng, J. Kamps, V. Murdock, J. Wen, and Y. Liu, Eds. ACM, 2020, pp. 479–488.
- [22] Z. Su, Z. Dou, Y. Zhu, X. Qin, and J. Wen, "Modeling intent graph for search result diversification," in *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, F. Diaz, C. Shah, T. Suel, P. Castells, R. Jones, and T. Sakai, Eds. ACM, 2021, pp. 736–746.
- [23] S. Borgeaud, A. Mensch, J. Hoffmann, T. Cai, E. Rutherford, K. Millican, G. van den Driessche, J. Lespiau, B. Damoc, A. Clark, D. de Las Casas, A. Guy, J. Menick, R. Ring, T. Hennigan, S. Huang, L. Maggiore, C. Jones, A. Cassirer, A. Brock, M. Paganini, G. Irving, O. Vinyals, S. Osindero, K. Simonyan, J. W. Rae, E. Elsen, and L. Sifre, "Improving language models by retrieving from trillions of tokens," in *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, ser. Proceedings of Machine Learning Research, K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvári, G. Niu, and S. Sabato, Eds., vol. 162. PMLR, 2022, pp. 2206–2240.
- [24] R. Nakano, J. Hilton, S. Balaji, J. Wu, L. Ouyang, C. Kim, C. Hesse, S. Jain, V. Kosaraju, W. Saunders, X. Jiang, K. Cobbe, T. Eloundou, G. Krueger, K. Button, M. Knight, B. Chess, and J. Schulman, "Webgpt: Browser-assisted question-answering with human feedback," *CoRR*, vol. abs/2112.09332, 2021.
- [25] G. Salton and M. McGill, *Introduction to Modern Information Retrieval*. McGraw-Hill Book Company, 1984.
- [26] G. Salton, A. Wong, and C. Yang, "A vector space model for automatic indexing," *Commun. ACM*, vol. 18, no. 11, pp. 613–620, 1975.
- [27] F. Song and W. B. Croft, "A general language model for information retrieval," in *Proceedings of the 1999 ACM CIKM International Conference on Information and Knowledge Management, Kansas City, Missouri, USA, November 2-6, 1999*. ACM, 1999, pp. 316–321.
- [28] J. Martineau and T. Finin, "Delta TFIDF: an improved feature space for sentiment analysis," in *Proceedings of the Third International Conference on Weblogs and Social Media, ICWSM 2009, San Jose, California, USA, May 17-20, 2009*, E. Adar, M. Hurst, T. Finin, N. S. Glance, N. Nicolov, and B. L. Tseng, Eds. The AAAI Press, 2009.
- [29] S. E. Robertson, S. Walker, S. Jones, M. Hancock-Beaulieu, and M. Gatford, "Okapi at TREC-3," in *Proceedings of The Third Text REtrieval Conference, TREC 1994, Gaithersburg, Maryland, USA, November 2-4, 1994*, ser. NIST Special Publication, D. K. Harman, Ed., vol. 500-225. National Institute of Standards and Technology (NIST), 1994, pp. 109–126.
- [30] J. Guo, Y. Fan, Q. Ai, and W. B. Croft, "A deep relevance matching model for ad-hoc retrieval," in *Proceedings of the 25th ACM International Conference on Information and Knowledge Management, CIKM 2016, Indianapolis, IN, USA, October 24-28, 2016*, S. Mukhopadhyay, C. Zhai, E. Bertino, F. Crestani, J. Mostafa, J. Tang, L. Si, X. Zhou, Y. Chang, Y. Li, and P. Sondhi, Eds. ACM, 2016, pp. 55–64.
- [31] L. Xiong, C. Xiong, Y. Li, K. Tang, J. Liu, P. N. Bennett, J. Ahmed, and A. Overwijk, "Approximate nearest neighbor negative contrastive learning for dense text retrieval," in *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [32] J. Lin, R. F. Nogueira, and A. Yates, *Pretrained Transformers for Text Ranking: BERT and Beyond*, ser. Synthesis Lectures on Human Language Technologies. Morgan & Claypool Publishers, 2021.
- [33] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language models are unsupervised multitask learners," 2019.
- [34] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger,

- T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., 2020.
- [35] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample, "Llama: Open and efficient foundation language models," *CoRR*, vol. abs/2302.13971, 2023.
- [36] J. Zhang, R. Xie, Y. Hou, W. X. Zhao, L. Lin, and J. Wen, "Recommendation as instruction following: A large language model empowered recommendation approach," *CoRR*, vol. abs/2305.07001, 2023.
- [37] Y. Hou, J. Zhang, Z. Lin, H. Lu, R. Xie, J. J. McAuley, and W. X. Zhao, "Large language models are zero-shot rankers for recommender systems," *CoRR*, vol. abs/2305.08845, 2023.
- [38] Y. Xi, W. Liu, J. Lin, J. Zhu, B. Chen, R. Tang, W. Zhang, R. Zhang, and Y. Yu, "Towards open-world recommendation with knowledge augmentation from large language models," *CoRR*, vol. abs/2306.10933, 2023.
- [39] W. Fan, Z. Zhao, J. Li, Y. Liu, X. Mei, Y. Wang, J. Tang, and Q. Li, "Recommender systems in the era of large language models (llms)," *CoRR*, vol. abs/2307.02046, 2023.
- [40] S. Wu, O. Irsoy, S. Lu, V. Dabrovolski, M. Dredze, S. Gehrmann, P. Kambadur, D. S. Rosenberg, and G. Mann, "Bloomberggpt: A large language model for finance," *CoRR*, vol. abs/2303.17564, 2023.
- [41] J. Li, Y. Liu, W. Fan, X. Wei, H. Liu, J. Tang, and Q. Li, "Empowering molecule discovery for molecule-caption translation with large language models: A chatgpt perspective," *CoRR*, vol. abs/2306.06615, 2023.
- [42] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, E. H. Chi, T. Hashimoto, O. Vinyals, P. Liang, J. Dean, and W. Fedus, "Emergent abilities of large language models," *Trans. Mach. Learn. Res.*, vol. 2022, 2022.
- [43] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. F. Christiano, J. Leike, and R. Lowe, "Training language models to follow instructions with human feedback," in *NeurIPS*, 2022.
- [44] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le, "Fine-tuned language models are zero-shot learners," in *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022.
- [45] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. H. Chi, Q. V. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," in *NeurIPS*, 2022.
- [46] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, "Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing," *ACM Comput. Surv.*, vol. 55, no. 9, pp. 195:1–195:35, 2023.
- [47] X. Qiu, T. Sun, Y. Xu, Y. Shao, N. Dai, and X. Huang, "Pre-trained models for natural language processing: A survey," *CoRR*, vol. abs/2003.08271, 2020.
- [48] Y. Cao, S. Li, Y. Liu, Z. Yan, Y. Dai, P. S. Yu, and L. Sun, "A comprehensive survey of ai-generated content (AIGC): A history of generative AI from GAN to chatgpt," *CoRR*, vol. abs/2303.04226, 2023.
- [49] J. Li, T. Tang, W. X. Zhao, and J. Wen, "Pretrained language model for text generation: A survey," in *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021*, Z. Zhou, Ed. ijcai.org, 2021, pp. 4492–4499.
- [50] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, L. Li, and Z. Sui, "A survey for in-context learning," *CoRR*, vol. abs/2301.00234, 2023.
- [51] J. Huang and K. C. Chang, "Towards reasoning in large language models: A survey," in *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, A. Rogers, J. L. Boyd-Graber, and N. Okazaki, Eds. Association for Computational Linguistics, 2023, pp. 1049–1065.
- [52] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, Y. Du, C. Yang, Y. Chen, Z. Chen, J. Jiang, R. Ren, Y. Li, X. Tang, Z. Liu, P. Liu, J. Nie, and J. Wen, "A survey of large language models," *CoRR*, vol. abs/2303.18223, 2023.
- [53] Q. Ai, T. Bai, Z. Cao, Y. Chang, J. Chen, Z. Chen, Z. Cheng, S. Dong, Z. Dou, F. Feng, S. Gao, J. Guo, X. He, Y. Lan, C. Li, Y. Liu, Z. Lyu, W. Ma, J. Ma, Z. Ren, P. Ren, Z. Wang, M. Wang, J. Wen, L. Wu, X. Xin, J. Xu, D. Yin, P. Zhang, F. Zhang, W. Zhang, M. Zhang, and X. Zhu, "Information retrieval meets large language models: A strategic report from chinese IR community," *CoRR*, vol. abs/2307.09751, 2023.
- [54] X. Liu and W. B. Croft, "Statistical language modeling for information retrieval," *Annu. Rev. Inf. Sci. Technol.*, vol. 39, no. 1, pp. 1–31, 2005.
- [55] B. Mitra and N. Craswell, "Neural models for information retrieval," *CoRR*, vol. abs/1705.01509, 2017.
- [56] W. X. Zhao, J. Liu, R. Ren, and J. Wen, "Dense text retrieval based on pretrained language models: A survey," *CoRR*, vol. abs/2211.14876, 2022.
- [57] 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," *J. Mach. Learn. Res.*, vol. 21, pp. 140:1–140:67, 2020.
- [58] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, "Deep contextualized word representations," in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana,*

- USA, June 1-6, 2018, Volume 1 (Long Papers), M. A. Walker, H. Ji, and A. Stent, Eds. Association for Computational Linguistics, 2018, pp. 2227–2237.
- [59] J. Devlin, M. Chang, K. Lee, and K. Toutanova, “BERT: pre-training of deep bidirectional transformers for language understanding,” in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, J. Burstein, C. Doran, and T. Solorio, Eds. Association for Computational Linguistics, 2019, pp. 4171–4186.
- [60] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, Eds., 2017, pp. 5998–6008.
- [61] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer, “BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, D. Jurafsky, J. Chai, N. Schluter, and J. R. Tetreault, Eds. Association for Computational Linguistics, 2020, pp. 7871–7880.
- [62] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei, “Scaling laws for neural language models,” *CoRR*, vol. abs/2001.08361, 2020.
- [63] A. Clark, D. de Las Casas, A. Guy, A. Mensch, M. Paganini, J. Hoffmann, B. Damoc, B. A. Hechtman, T. Cai, S. Borgeaud, G. van den Driessche, E. Rutherford, T. Hennigan, M. J. Johnson, A. Cassirer, C. Jones, E. Buchatskaya, D. Budden, L. Sifre, S. Osindero, O. Vinyals, M. Ranzato, J. W. Rae, E. Elsen, K. Kavukcuoglu, and K. Simonyan, “Unified scaling laws for routed language models,” in *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, ser. Proceedings of Machine Learning Research, K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvári, G. Niu, and S. Sabato, Eds., vol. 162. PMLR, 2022, pp. 4057–4086.
- [64] L. Dong, N. Yang, W. Wang, F. Wei, X. Liu, Y. Wang, J. Gao, M. Zhou, and H. Hon, “Unified language model pre-training for natural language understanding and generation,” in *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, and R. Garnett, Eds., 2019, pp. 13042–13054.
- [65] L. Xue, N. Constant, A. Roberts, M. Kale, R. Al-Rfou, A. Siddhant, A. Barua, and C. Raffel, “mt5: A massively multilingual pre-trained text-to-text transformer,” in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, K. Toutanova, A. Rumshisky, L. Zettlemoyer, D. Hakkani-Tür, I. Beltagy, S. Bethard, R. Cotterell, T. Chakraborty, and Y. Zhou, Eds. Association for Computational Linguistics, 2021, pp. 483–498.
- [66] V. Sanh, A. Webson, C. Raffel, S. H. Bach, L. Sutawika, Z. Alyafeai, A. Chaffin, A. Stiegler, A. Raja, M. Dey, M. S. Bari, C. Xu, U. Thakker, S. S. Sharma, E. Szczechla, T. Kim, G. Chhablani, N. V. Nayak, D. Datta, J. Chang, M. T. Jiang, H. Wang, M. Manica, S. Shen, Z. X. Yong, H. Pandey, R. Bawden, T. Wang, T. Neeraj, J. Rozen, A. Sharma, A. Santilli, T. Févry, J. A. Fries, R. Teehan, T. L. Scao, S. Biderman, L. Gao, T. Wolf, and A. M. Rush, “Multitask prompted training enables zero-shot task generalization,” in *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022.
- [67] H. Bao, L. Dong, F. Wei, W. Wang, N. Yang, X. Liu, Y. Wang, J. Gao, S. Piao, M. Zhou, and H. Hon, “Unilmv2: Pseudo-masked language models for unified language model pre-training,” in *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 642–652.
- [68] A. Zeng, X. Liu, Z. Du, Z. Wang, H. Lai, M. Ding, Z. Yang, Y. Xu, W. Zheng, X. Xia, W. L. Tam, Z. Ma, Y. Xue, J. Zhai, W. Chen, Z. Liu, P. Zhang, Y. Dong, and J. Tang, “GLM-130B: an open bilingual pre-trained model,” in *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- [69] W. Fedus, B. Zoph, and N. Shazeer, “Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity,” *J. Mach. Learn. Res.*, vol. 23, pp. 120:1–120:39, 2022.
- [70] Z. Yang, Z. Dai, Y. Yang, J. G. Carbonell, R. Salakhutdinov, and Q. V. Le, “Xlnet: Generalized autoregressive pretraining for language understanding,” in *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, and R. Garnett, Eds., 2019, pp. 5754–5764.
- [71] S. Black, S. Biderman, E. Hallahan, Q. Anthony, L. Gao, L. Golding, H. He, C. Leahy, K. McDonnell, J. Phang, M. Pieler, U. S. Prashanth, S. Purohit, L. Reynolds, J. Tow, B. Wang, and S. Weinbach, “Gpt-neox-20b: An open-source autoregressive language model,” *CoRR*, vol. abs/2204.06745, 2022.
- [72] J. W. Rae, S. Borgeaud, T. Cai, K. Millican, J. Hoffmann, H. F. Song, J. Aslanides, S. Henderson, R. Ring, S. Young, E. Rutherford, T. Hennigan, J. Menick, A. Cassirer, R. Powell, G. van den Driessche, L. A. Hendricks, M. Rauh, P. Huang, A. Glaese, J. Welbl, S. Dathathri, S. Huang, J. Uesato, J. Mellor, I. Higgins, A. Creswell, N. McAleese, A. Wu, E. Elsen, S. M.

- Jayakumar, E. Buchatskaya, D. Budden, E. Sutherland, K. Simonyan, M. Paganini, L. Sifre, L. Martens, X. L. Li, A. Kuncoro, A. Nematzadeh, E. Gribovskaya, D. Donato, A. Lazaridou, A. Mensch, J. Lespiau, M. Tsimpoukelli, N. Grigorev, D. Fritz, T. Sottiaux, M. Pajarskas, T. Pohlen, Z. Gong, D. Toyama, C. de Masson d’Autume, Y. Li, T. Terzi, V. Mikulik, I. Babuschkin, A. Clark, D. de Las Casas, A. Guy, C. Jones, J. Bradbury, M. J. Johnson, B. A. Hechtman, L. Weidinger, I. Gabriel, W. Isaac, E. Lockhart, S. Osindero, L. Rimell, C. Dyer, O. Vinyals, K. Ayoub, J. Stanway, L. Bennett, D. Hassabis, K. Kavukcuoglu, and G. Irving, “Scaling language models: Methods, analysis & insights from training gopher,” *CoRR*, vol. abs/2112.11446, 2021.
- [73] N. Du, Y. Huang, A. M. Dai, S. Tong, D. Lepikhin, Y. Xu, M. Krikun, Y. Zhou, A. W. Yu, O. Firat, B. Zoph, L. Fedus, M. P. Bosma, Z. Zhou, T. Wang, Y. E. Wang, K. Webster, M. Pellat, K. Robinson, K. S. Meier-Hellstern, T. Duke, L. Dixon, K. Zhang, Q. V. Le, Y. Wu, Z. Chen, and C. Cui, “Glam: Efficient scaling of language models with mixture-of-experts,” in *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, ser. Proceedings of Machine Learning Research, K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvári, G. Niu, and S. Sabato, Eds., vol. 162. PMLR, 2022, pp. 5547–5569.
- [74] Y. Sun, S. Wang, S. Feng, S. Ding, C. Pang, J. Shang, J. Liu, X. Chen, Y. Zhao, Y. Lu, W. Liu, Z. Wu, W. Gong, J. Liang, Z. Shang, P. Sun, W. Liu, X. Ouyang, D. Yu, H. Tian, H. Wu, and H. Wang, “ERNIE 3.0: Large-scale knowledge enhanced pre-training for language understanding and generation,” *CoRR*, vol. abs/2107.02137, 2021.
- [75] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. T. Diab, X. Li, X. V. Lin, T. Mihaylov, M. Ott, S. Shleifer, K. Shuster, D. Simig, P. S. Koura, A. Sridhar, T. Wang, and L. Zettlemoyer, “OPT: open pre-trained transformer language models,” *CoRR*, vol. abs/2205.01068, 2022.
- [76] R. Thoppilan, D. D. Freitas, J. Hall, N. Shazeer, A. Kulshreshtha, H. Cheng, A. Jin, T. Bos, L. Baker, Y. Du, Y. Li, H. Lee, H. S. Zheng, A. Ghafouri, M. Menegali, Y. Huang, M. Krikun, D. Lepikhin, J. Qin, D. Chen, Y. Xu, Z. Chen, A. Roberts, M. Bosma, Y. Zhou, C. Chang, I. Krivokon, W. Rusch, M. Pickett, K. S. Meier-Hellstern, M. R. Morris, T. Doshi, R. D. Santos, T. Duke, J. Soraker, B. Zevenbergen, V. Prabhakaran, M. Diaz, B. Hutchinson, K. Olson, A. Molina, E. Hoffman-John, J. Lee, L. Aroyo, R. Rajakumar, A. Butryna, M. Lamm, V. Kuzmina, J. Fenton, A. Cohen, R. Bernstein, R. Kurzweil, B. A. y Arcas, C. Cui, M. Croak, E. H. Chi, and Q. Le, “Lamda: Language models for dialog applications,” *CoRR*, vol. abs/2201.08239, 2022.
- [77] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, P. Schuh, K. Shi, S. Tsvyashchenko, J. Maynez, A. Rao, P. Barnes, Y. Tay, N. Shazeer, V. Prabhakaran, E. Reif, N. Du, B. Hutchinson, R. Pope, J. Bradbury, J. Austin, M. Isard, G. Gur-Ari, P. Yin, T. Duke, A. Levskaya, S. Ghemawat, S. Dev, H. Michalewski, X. Garcia, V. Misra, K. Robinson, L. Fedus, D. Zhou, D. Ippolito, D. Luan, H. Lim, B. Zoph, A. Spiridonov, R. Sepassi, D. Dohan, S. Agrawal, M. Omernick, A. M. Dai, T. S. Pillelai, M. Pellat, A. Lewkowycz, E. Moreira, R. Child, O. Polozov, K. Lee, Z. Zhou, X. Wang, B. Saeta, M. Diaz, O. Firat, M. Catasta, J. Wei, K. Meier-Hellstern, D. Eck, J. Dean, S. Petrov, and N. Fiedel, “Palm: Scaling language modeling with pathways,” *CoRR*, vol. abs/2204.02311, 2022.
- [78] T. L. Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilic, D. Hesslow, R. Castagné, A. S. Luccioni, F. Yvon, M. Gallé, J. Tow, A. M. Rush, S. Biderman, A. Webson, P. S. Ammanamanchi, T. Wang, B. Sagot, N. Muennighoff, A. V. del Moral, O. Ruwase, R. Bawden, S. Bekman, A. McMillan-Major, I. Beltagy, H. Nguyen, L. Saulnier, S. Tan, P. O. Suarez, V. Sanh, H. Laurençon, Y. Jernite, J. Launay, M. Mitchell, C. Raffel, A. Gokaslan, A. Simhi, A. Soroa, A. F. Aji, A. Alfassy, A. Rogers, A. K. Nitzav, C. Xu, C. Mou, C. Emezue, C. Klamm, C. Leong, D. van Strien, D. I. Adelan, and et al., “BLOOM: A 176b-parameter open-access multilingual language model,” *CoRR*, vol. abs/2211.05100, 2022.
- [79] A. Lewkowycz, A. Andreassen, D. Dohan, E. Dyer, H. Michalewski, V. V. Ramasesh, A. Slone, C. Anil, I. Schlag, T. Gutman-Solo, Y. Wu, B. Neyshabur, G. Gur-Ari, and V. Misra, “Solving quantitative reasoning problems with language models,” in *NeurIPS*, 2022.
- [80] OpenAI, “GPT-4 technical report,” *CoRR*, vol. abs/2303.08774, 2023.
- [81] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. de Las Casas, L. A. Hendricks, J. Welbl, A. Clark, T. Hennigan, E. Noland, K. Millican, G. van den Driessche, B. Damoc, A. Guy, S. Osindero, K. Simonyan, E. Elsen, J. W. Rae, O. Vinyals, and L. Sifre, “Training compute-optimal large language models,” *CoRR*, vol. abs/2203.15556, 2022.
- [82] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, “Lora: Low-rank adaptation of large language models,” in *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022.
- [83] X. L. Li and P. Liang, “Prefix-tuning: Optimizing continuous prompts for generation,” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, C. Zong, F. Xia, W. Li, and R. Navigli, Eds. Association for Computational Linguistics, 2021, pp. 4582–4597.
- [84] B. Lester, R. Al-Rfou, and N. Constant, “The power of scale for parameter-efficient prompt tuning,” in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, M. Moens, X. Huang, L. Specia, and S. W. Yih, Eds. Association for Computational Linguistics, 2021,

- pp. 3045–3059.
- [85] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, “Qlora: Efficient finetuning of quantized llms,” *CoRR*, vol. abs/2305.14314, 2023.
- [86] L. Wang, N. Yang, and F. Wei, “Query2doc: Query expansion with large language models,” pp. 9414–9423, 2023.
- [87] N. A. Jaleel, J. Allan, W. B. Croft, F. Diaz, L. S. Larkey, X. Li, M. D. Smucker, and C. Wade, “Umass at TREC 2004: Novelty and HARD,” in *Proceedings of the Thirteenth Text REtrieval Conference, TREC 2004, Gaithersburg, Maryland, USA, November 16-19, 2004*, ser. NIST Special Publication, E. M. Voorhees and L. P. Buckland, Eds., vol. 500-261. National Institute of Standards and Technology (NIST), 2004.
- [88] D. Metzler and W. B. Croft, “Latent concept expansion using markov random fields,” in *SIGIR 2007: Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Amsterdam, The Netherlands, July 23-27, 2007*, W. Kraaij, A. P. de Vries, C. L. A. Clarke, N. Fuhr, and N. Kando, Eds. ACM, 2007, pp. 311–318.
- [89] C. Zhai and J. D. Lafferty, “Model-based feedback in the language modeling approach to information retrieval,” in *Proceedings of the 2001 ACM CIKM International Conference on Information and Knowledge Management, Atlanta, Georgia, USA, November 5-10, 2001*. ACM, 2001, pp. 403–410.
- [90] D. Metzler and W. B. Croft, “A markov random field model for term dependencies,” in *SIGIR 2005: Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Salvador, Brazil, August 15-19, 2005*, R. A. Baeza-Yates, N. Ziviani, G. Marchionini, A. Moffat, and J. Tait, Eds. ACM, 2005, pp. 472–479.
- [91] X. Wang, C. Macdonald, N. Tonello, and I. Ounis, “Pseudo-relevance feedback for multiple representation dense retrieval,” in *ICTIR ’21: The 2021 ACM SIGIR International Conference on the Theory of Information Retrieval, Virtual Event, Canada, July 11, 2021*, F. Hasibi, Y. Fang, and A. Aizawa, Eds. ACM, 2021, pp. 297–306.
- [92] Z. Zheng, K. Hui, B. He, X. Han, L. Sun, and A. Yates, “BERT-QE: contextualized query expansion for document re-ranking,” in *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, ser. Findings of ACL, T. Cohn, Y. He, and Y. Liu, Eds., vol. EMNLP 2020. Association for Computational Linguistics, 2020, pp. 4718–4728.
- [93] F. Diaz, B. Mitra, and N. Craswell, “Query expansion with locally-trained word embeddings,” in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics, 2016.
- [94] S. Kuzi, A. Shtok, and O. Kurland, “Query expansion using word embeddings,” in *Proceedings of the 25th ACM International Conference on Information and Knowledge Management, CIKM 2016, Indianapolis, IN, USA, October 24-28, 2016*, S. Mukhopadhyay, C. Zhai, E. Bertino, F. Crestani, J. Mostafa, J. Tang, L. Si, X. Zhou, Y. Chang, Y. Li, and P. Sondhi, Eds. ACM, 2016, pp. 1929–1932.
- [95] K. Mao, Z. Dou, F. Mo, J. Hou, H. Chen, and H. Qian, “Large language models know your contextual search intent: A prompting framework for conversational search,” pp. 1211–1225, 2023.
- [96] I. Mackie, I. Sekulic, S. Chatterjee, J. Dalton, and F. Crestani, “GRM: generative relevance modeling using relevance-aware sample estimation for document retrieval,” *CoRR*, vol. abs/2306.09938, 2023.
- [97] K. Srinivasan, K. Raman, A. Samanta, L. Liao, L. Bertelli, and M. Bendersky, “QUILL: query intent with large language models using retrieval augmentation and multi-stage distillation,” in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: EMNLP 2022 - Industry Track, Abu Dhabi, UAE, December 7 - 11, 2022*, Y. Li and A. Lazaridou, Eds. Association for Computational Linguistics, 2022, pp. 492–501.
- [98] J. Feng, C. Tao, X. Geng, T. Shen, C. Xu, G. Long, D. Zhao, and D. Jiang, “Knowledge refinement via interaction between search engines and large language models,” *CoRR*, vol. abs/2305.07402, 2023.
- [99] I. Mackie, S. Chatterjee, and J. Dalton, “Generative and pseudo-relevant feedback for sparse, dense and learned sparse retrieval,” *CoRR*, vol. abs/2305.07477, 2023.
- [100] X. Ma, Y. Gong, P. He, H. Zhao, and N. Duan, “Query rewriting for retrieval-augmented large language models,” *CoRR*, vol. abs/2305.14283, 2023.
- [101] L. Gao, X. Ma, J. Lin, and J. Callan, “Precise zero-shot dense retrieval without relevance labels,” *CoRR*, vol. abs/2212.10496, 2022.
- [102] R. Jagerman, H. Zhuang, Z. Qin, X. Wang, and M. Bendersky, “Query expansion by prompting large language models,” *CoRR*, vol. abs/2305.03653, 2023.
- [103] Y. Tang, R. Qiu, and X. Li, “Prompt-based effective input reformulation for legal case retrieval,” in *Databases Theory and Applications - 34th Australasian Database Conference, ADC 2023, Melbourne, VIC, Australia, November 1-3, 2023, Proceedings*, ser. Lecture Notes in Computer Science, Z. Bao, R. Borovica-Gajic, R. Qiu, F. M. Choudhury, and Z. Yang, Eds., vol. 14386. Springer, 2023, pp. 87–100.
- [104] F. Ye, M. Fang, S. Li, and E. Yilmaz, “Enhancing conversational search: Large language model-aided informative query rewriting,” *arXiv preprint arXiv:2310.09716*, 2023.
- [105] C. Huang, C. Hsu, T. Hsu, C. Li, and Y. Chen, “CONVERSER: few-shot conversational dense retrieval with synthetic data generation,” in *Proceedings of the 24th Meeting of the Special Interest Group on Discourse and Dialogue, SIGDIAL 2023, Prague, Czechia, September 11 - 15, 2023*, D. Schlangen, S. Stoyanchev, S. Joty, O. Dusek, C. Kennington, and M. Alikhani, Eds. Association for Computational Linguistics, 2023, pp. 381–387.
- [106] M. Li, H. Zhuang, K. Hui, Z. Qin, J. Lin, R. Jagerman, X. Wang, and M. Bendersky, “Generate, filter, and fuse: Query expansion via multi-step keyword generation for zero-shot neural rankers,” *CoRR*, vol.

- abs/2311.09175, 2023.
- [107] A. Anand, V. V. V. Setty, and A. Anand, "Context aware query rewriting for text rankers using LLM," *CoRR*, vol. abs/2308.16753, 2023.
- [108] T. Shen, G. Long, X. Geng, C. Tao, T. Zhou, and D. Jiang, "Large language models are strong zero-shot retriever," *CoRR*, vol. abs/2304.14233, 2023.
- [109] M. Alaofi, L. Gallagher, M. Sanderson, F. Scholer, and P. Thomas, "Can generative llms create query variants for test collections? an exploratory study," in *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2023, Taipei, Taiwan, July 23-27, 2023*, H. Chen, W. E. Duh, H. Huang, M. P. Kato, J. Mothe, and B. Poblete, Eds. ACM, 2023, pp. 1869–1873.
- [110] W. Yu, D. Iter, S. Wang, Y. Xu, M. Ju, S. Sanyal, C. Zhu, M. Zeng, and M. Jiang, "Generate rather than retrieve: Large language models are strong context generators," in *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- [111] T. Nguyen, M. Rosenberg, X. Song, J. Gao, S. Tiwary, R. Majumder, and L. Deng, "MS MARCO: A human generated machine reading comprehension dataset," in *CoCo@NIPS*, ser. CEUR Workshop Proceedings, vol. 1773. CEUR-WS.org, 2016.
- [112] T. Kwiatkowski, J. Palomaki, O. Redfield, M. Collins, A. P. Parikh, C. Alberti, D. Epstein, I. Polosukhin, J. Devlin, K. Lee, K. Toutanova, L. Jones, M. Kelcey, M. Chang, A. M. Dai, J. Uszkoreit, Q. Le, and S. Petrov, "Natural questions: a benchmark for question answering research," *Trans. Assoc. Comput. Linguistics*, vol. 7, pp. 452–466, 2019.
- [113] W. Peng, G. Li, Y. Jiang, Z. Wang, D. Ou, X. Zeng, D. Xu, T. Xu, and E. Chen, "Large language model based long-tail query rewriting in taobao search," *CoRR*, vol. abs/2311.03758, 2023.
- [114] Z. Du, Y. Qian, X. Liu, M. Ding, J. Qiu, Z. Yang, and J. Tang, "GLM: general language model pretraining with autoregressive blank infilling," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, S. Muresan, P. Nakov, and A. Villavicencio, Eds. Association for Computational Linguistics, 2022, pp. 320–335.
- [115] A. Yang, B. Xiao, B. Wang, B. Zhang, C. Bian, C. Yin, C. Lv, D. Pan, D. Wang, D. Yan, F. Yang, F. Deng, F. Wang, F. Liu, G. Ai, G. Dong, H. Zhao, H. Xu, H. Sun, H. Zhang, H. Liu, J. Ji, J. Xie, J. Dai, K. Fang, L. Su, L. Song, L. Liu, L. Ru, L. Ma, M. Wang, M. Liu, M. Lin, N. Nie, P. Guo, R. Sun, T. Zhang, T. Li, T. Li, W. Cheng, W. Chen, X. Zeng, X. Wang, X. Chen, X. Men, X. Yu, X. Pan, Y. Shen, Y. Wang, Y. Li, Y. Jiang, Y. Gao, Y. Zhang, Z. Zhou, and Z. Wu, "Baichuan 2: Open large-scale language models," *CoRR*, vol. abs/2309.10305, 2023.
- [116] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang, B. Hui, L. Ji, M. Li, J. Lin, R. Lin, D. Liu, G. Liu, C. Lu, K. Lu, J. Ma, R. Men, X. Ren, X. Ren, C. Tan, S. Tan, J. Tu, P. Wang, S. Wang, W. Wang, S. Wu, B. Xu, J. Xu, A. Yang, H. Yang, J. Yang, S. Yang, Y. Yao, B. Yu, H. Yuan, Z. Yuan, J. Zhang, X. Zhang, Y. Zhang, Z. Zhang, C. Zhou, J. Zhou, X. Zhou, and T. Zhu, "Qwen technical report," *CoRR*, vol. abs/2309.16609, 2023.
- [117] D. Alexander, W. Kusa, and A. P. de Vries, "ORCAS-I: queries annotated with intent using weak supervision," in *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Madrid, Spain, July 11 - 15, 2022*, E. Amigó, P. Castells, J. Gonzalo, B. Carterette, J. S. Culpepper, and G. Kazai, Eds. ACM, 2022, pp. 3057–3066.
- [118] K. D. Dhole, R. Chandradevan, and E. Agichtein, "An interactive query generation assistant using llm-based prompt modification and user feedback," *CoRR*, vol. abs/2311.11226, 2023.
- [119] O. Weller, K. Lo, D. Wadden, D. J. Lawrie, B. V. Durme, A. Cohan, and L. Soldaini, "When do generative query and document expansions fail? A comprehensive study across methods, retrievers, and datasets," *CoRR*, vol. abs/2309.08541, 2023.
- [120] L. H. Bonifacio, H. Abonizio, M. Fadaee, and R. F. Nogueira, "Inpars: Data augmentation for information retrieval using large language models," *CoRR*, vol. abs/2202.05144, 2022.
- [121] G. Ma, X. Wu, P. Wang, Z. Lin, and S. Hu, "Pre-training with large language model-based document expansion for dense passage retrieval," *CoRR*, vol. abs/2308.08285, 2023.
- [122] V. Jeronymo, L. H. Bonifacio, H. Abonizio, M. Fadaee, R. de Alencar Lotufo, J. Zavrel, and R. F. Nogueira, "Inpars-v2: Large language models as efficient dataset generators for information retrieval," *CoRR*, vol. abs/2301.01820, 2023.
- [123] Z. Dai, V. Y. Zhao, J. Ma, Y. Luan, J. Ni, J. Lu, A. Bakalov, K. Guu, K. B. Hall, and M. Chang, "Promptagator: Few-shot dense retrieval from 8 examples," in *ICLR*. OpenReview.net, 2023.
- [124] R. Meng, Y. Liu, S. Yavuz, D. Agarwal, L. Tu, N. Yu, J. Zhang, M. Bhat, and Y. Zhou, "Augtriever: Unsupervised dense retrieval by scalable data augmentation," 2023.
- [125] J. Saad-Falcon, O. Khattab, K. Santhanam, R. Florian, M. Franz, S. Roukos, A. Sil, M. A. Sultan, and C. Potts, "UDAPDR: unsupervised domain adaptation via LLM prompting and distillation of rerankers," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 11 265–11 279.
- [126] Z. Peng, X. Wu, and Y. Fang, "Soft prompt tuning for augmenting dense retrieval with large language models," 2023.
- [127] D. S. Sachan, M. Lewis, D. Yogatama, L. Zettlemoyer, J. Pineau, and M. Zaheer, "Questions are all you need to train a dense passage retriever," *Transactions of the Association for Computational Linguistics*, vol. 11, pp. 600–616, 2023.
- [128] N. Thakur, N. Reimers, A. Rücklé, A. Srivastava, and I. Gurevych, "BEIR: A heterogeneous benchmark for zero-shot evaluation of information retrieval models,"

- in *NeurIPS Datasets and Benchmarks*, 2021.
- [129] N. Thakur, J. Ni, G. H. Ábrego, J. Wieting, J. Lin, and D. Cer, “Leveraging llms for synthesizing training data across many languages in multilingual dense retrieval,” *CoRR*, vol. abs/2311.05800, 2023.
- [130] A. Neelakantan, T. Xu, R. Puri, A. Radford, J. M. Han, J. Tworek, Q. Yuan, N. Tezak, J. W. Kim, C. Hallacy, J. Heidecke, P. Shyam, B. Power, T. E. Niekoul, G. Sastry, G. Krueger, D. Schnurr, F. P. Such, K. Hsu, M. Thompson, T. Khan, T. Sherbakov, J. Jang, P. Welinder, and L. Weng, “Text and code embeddings by contrastive pre-training,” *CoRR*, vol. abs/2201.10005, 2022.
- [131] X. Ma, L. Wang, N. Yang, F. Wei, and J. Lin, “Fine-tuning llama for multi-stage text retrieval,” *CoRR*, vol. abs/2310.08319, 2023.
- [132] A. Asai, T. Schick, P. S. H. Lewis, X. Chen, G. Izacard, S. Riedel, H. Hajishirzi, and W. Yih, “Task-aware retrieval with instructions,” in *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, A. Rogers, J. L. Boyd-Graber, and N. Okazaki, Eds. Association for Computational Linguistics, 2023, pp. 3650–3675.
- [133] J. Ni, C. Qu, J. Lu, Z. Dai, G. H. Ábrego, J. Ma, V. Y. Zhao, Y. Luan, K. B. Hall, M. Chang, and Y. Yang, “Large dual encoders are generalizable retrievers,” in *EMNLP*. Association for Computational Linguistics, 2022, pp. 9844–9855.
- [134] G. Izacard, M. Caron, L. Hosseini, S. Riedel, P. Bojanowski, A. Joulin, and E. Grave, “Unsupervised dense information retrieval with contrastive learning,” *Trans. Mach. Learn. Res.*, vol. 2022, 2022.
- [135] D. Metzler, Y. Tay, D. Bahri, and M. Najork, “Rethinking search: making domain experts out of dilettantes,” *SIGIR Forum*, vol. 55, no. 1, pp. 13:1–13:27, 2021.
- [136] Y. Zhou, J. Yao, Z. Dou, L. Wu, and J. Wen, “Dynamicretriever: A pre-trained model-based IR system without an explicit index,” *Mach. Intell. Res.*, vol. 20, no. 2, pp. 276–288, 2023.
- [137] J. Chen, R. Zhang, J. Guo, Y. Liu, Y. Fan, and X. Cheng, “Corpusbrain: Pre-train a generative retrieval model for knowledge-intensive language tasks,” in *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17-21, 2022*, M. A. Hasan and L. Xiong, Eds. ACM, 2022, pp. 191–200.
- [138] Y. Tay, V. Tran, M. Dehghani, J. Ni, D. Bahri, H. Mehta, Z. Qin, K. Hui, Z. Zhao, J. P. Gupta, T. Schuster, W. W. Cohen, and D. Metzler, “Transformer memory as a differentiable search index,” in *NeurIPS*, 2022.
- [139] N. Ziemis, W. Yu, Z. Zhang, and M. Jiang, “Large language models are built-in autoregressive search engines,” in *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, A. Rogers, J. L. Boyd-Graber, and N. Okazaki, Eds. Association for Computational Linguistics, 2023, pp. 2666–2678.
- [140] R. F. Nogueira, W. Yang, K. Cho, and J. Lin, “Multi-stage document ranking with BERT,” *CoRR*, vol. abs/1910.14424, 2019.
- [141] J. Ju, J. Yang, and C. Wang, “Text-to-text multi-view learning for passage re-ranking,” in *SIGIR*. ACM, 2021, pp. 1803–1807.
- [142] R. Pradeep, R. F. Nogueira, and J. Lin, “The expando-mono-duo design pattern for text ranking with pre-trained sequence-to-sequence models,” *CoRR*, vol. abs/2101.05667, 2021.
- [143] H. Zhuang, Z. Qin, R. Jagerman, K. Hui, J. Ma, J. Lu, J. Ni, X. Wang, and M. Bendersky, “Rankt5: Fine-tuning T5 for text ranking with ranking losses,” in *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2023, Taipei, Taiwan, July 23-27, 2023*, H. Chen, W. E. Duh, H. Huang, M. P. Kato, J. Mothe, and B. Poblete, Eds. ACM, 2023, pp. 2308–2313.
- [144] L. Zhang, Y. Zhang, D. Long, P. Xie, M. Zhang, and M. Zhang, “Rankinggpt: Empowering large language models in text ranking with progressive enhancement,” *CoRR*, vol. abs/2311.16720, 2023.
- [145] X. Zhang, S. Hofstätter, P. Lewis, R. Tang, and J. Lin, “Rank-without-gpt: Building gpt-independent list-wise rerankers on open-source large language models,” *arXiv preprint arXiv:2312.02969*, 2023.
- [146] P. Liang, R. Bommasani, T. Lee, D. Tsipras, D. Soylu, M. Yasunaga, Y. Zhang, D. Narayanan, Y. Wu, A. Kumar, B. Newman, B. Yuan, B. Yan, C. Zhang, C. Cosgrove, C. D. Manning, C. Ré, D. Acosta-Navas, D. A. Hudson, E. Zelikman, E. Durmus, F. Ladhak, F. Rong, H. Ren, H. Yao, J. Wang, K. Santhanam, L. J. Orr, L. Zheng, M. Yüsekçönlü, M. Suzgun, N. Kim, N. Guha, N. S. Chatterji, O. Khattab, P. Henderson, Q. Huang, R. Chi, S. M. Xie, S. Santurkar, S. Ganguli, T. Hashimoto, T. Icard, T. Zhang, V. Chaudhary, W. Wang, X. Li, Y. Mai, Y. Zhang, and Y. Koreeda, “Holistic evaluation of language models,” *CoRR*, vol. abs/2211.09110, 2022.
- [147] H. Zhuang, Z. Qin, K. Hui, J. Wu, L. Yan, X. Wang, and M. Bendersky, “Beyond yes and no: Improving zero-shot LLM rankers via scoring fine-grained relevance labels,” *CoRR*, vol. abs/2310.14122, 2023.
- [148] D. S. Sachan, M. Lewis, M. Joshi, A. Aghajanyan, W. Yih, J. Pineau, and L. Zettlemoyer, “Improving passage retrieval with zero-shot question generation,” in *EMNLP*. Association for Computational Linguistics, 2022, pp. 3781–3797.
- [149] S. Zhuang, B. Liu, B. Koopman, and G. Zuccon, “Open-source large language models are strong zero-shot query likelihood models for document ranking,” in *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 8807–8817.
- [150] S. Cho, S. Jeong, J. Seo, and J. C. Park, “Discrete prompt optimization via constrained generation for zero-shot re-ranker,” in *ACL (Findings)*. Association for Computational Linguistics, 2023, pp. 960–971.
- [151] 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 LLMs,” in *Findings of the Association for Computational Linguistics: EMNLP 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Singapore: Association

- for Computational Linguistics, Dec. 2023, pp. 14242–14252.
- [152] 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,” in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 14918–14937.
- [153] X. Ma, X. Zhang, R. Pradeep, and J. Lin, “Zero-shot listwise document reranking with a large language model,” *CoRR*, vol. abs/2305.02156, 2023.
- [154] R. Tang, X. Zhang, X. Ma, J. Lin, and F. Ture, “Found in the middle: Permutation self-consistency improves listwise ranking in large language models,” *CoRR*, vol. abs/2310.07712, 2023.
- [155] Z. Qin, R. Jagerman, K. Hui, H. Zhuang, J. Wu, J. Shen, T. Liu, J. Liu, D. Metzler, X. Wang *et al.*, “Large language models are effective text rankers with pairwise ranking prompting,” *arXiv preprint arXiv:2306.17563*, 2023.
- [156] S. Zhuang, H. Zhuang, B. Koopman, and G. Zuccon, “A setwise approach for effective and highly efficient zero-shot ranking with large language models,” *CoRR*, vol. abs/2310.09497, 2023.
- [157] F. Ferraretto, T. Laitz, R. de Alencar Lotufo, and R. F. Nogueira, “Exaranker: Synthetic explanations improve neural rankers,” in *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2023, Taipei, Taiwan, July 23-27, 2023*, H. Chen, W. E. Duh, H. Huang, M. P. Kato, J. Mothe, and B. Poblete, Eds. ACM, 2023, pp. 2409–2414.
- [158] L. Boytsov, P. Patel, V. Sourabh, R. Nisar, S. Kundu, R. Ramanathan, and E. Nyberg, “Inpars-light: Cost-effective unsupervised training of efficient rankers,” *CoRR*, vol. abs/2301.02998, 2023.
- [159] A. Askari, M. Aliannejadi, E. Kanoulas, and S. Verberne, “Generating synthetic documents for cross-encoder re-rankers: A comparative study of chatgpt and human experts,” *CoRR*, vol. abs/2305.02320, 2023.
- [160] R. Pradeep, S. Sharifmoghaddam, and J. Lin, “Rankvicuna: Zero-shot listwise document reranking with open-source large language models,” *CoRR*, vol. abs/2309.15088, 2023.
- [161] —, “Rankzephyr: Effective and robust zero-shot listwise reranking is a breeze!” *CoRR*, vol. abs/2312.02724, 2023.
- [162] W. Sun, Z. Chen, X. Ma, L. Yan, S. Wang, P. Ren, Z. Chen, D. Yin, and Z. Ren, “Instruction distillation makes large language models efficient zero-shot rankers,” *arXiv preprint arXiv:2311.01555*, 2023.
- [163] C. J. C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, and G. N. Hullender, “Learning to rank using gradient descent,” in *ICML*, ser. ACM International Conference Proceeding Series, vol. 119. ACM, 2005, pp. 89–96.
- [164] J. A. Baktash and M. Dawodi, “Gpt-4: A review on advancements and opportunities in natural language processing,” *arXiv preprint arXiv:2305.03195*, 2023.
- [165] H. Wachsmuth, S. Syed, and B. Stein, “Retrieval of the best counterargument without prior topic knowledge,” in *ACL (1)*. Association for Computational Linguistics, 2018, pp. 241–251.
- [166] K. Guu, K. Lee, Z. Tung, P. Pasupat, and M. Chang, “Retrieval augmented language model pre-training,” in *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 3929–3938.
- [167] P. S. H. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. Yih, T. Rocktäschel, S. Riedel, and D. Kiela, “Retrieval-augmented generation for knowledge-intensive NLP tasks,” in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., 2020.
- [168] W. Shi, S. Min, M. Yasunaga, M. Seo, R. James, M. Lewis, L. Zettlemoyer, and W. Yih, “REPLUG: retrieval-augmented black-box language models,” *CoRR*, vol. abs/2301.12652, 2023.
- [169] G. Izacard, P. S. H. Lewis, M. Lomeli, L. Hosseini, F. Petroni, T. Schick, J. Dwivedi-Yu, A. Joulin, S. Riedel, and E. Grave, “Atlas: Few-shot learning with retrieval augmented language models,” *J. Mach. Learn. Res.*, vol. 24, pp. 251:1–251:43, 2023.
- [170] A. Lazaridou, E. Gribovskaya, W. Stokowiec, and N. Grigorev, “Internet-augmented language models through few-shot prompting for open-domain question answering,” *CoRR*, vol. abs/2203.05115, 2022.
- [171] H. He, H. Zhang, and D. Roth, “Rethinking with retrieval: Faithful large language model inference,” *CoRR*, vol. abs/2301.00303, 2023.
- [172] W. Yu, H. Zhang, X. Pan, K. Ma, H. Wang, and D. Yu, “Chain-of-note: Enhancing robustness in retrieval-augmented language models,” *CoRR*, vol. abs/2311.09210, 2023.
- [173] O. Ram, Y. Levine, I. Dalmedigos, D. Muhlgay, A. Shashua, K. Leyton-Brown, and Y. Shoham, “In-context retrieval-augmented language models,” *CoRR*, vol. abs/2302.00083, 2023.
- [174] Z. Shao, Y. Gong, Y. Shen, M. Huang, N. Duan, and W. Chen, “Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy,” in *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 9248–9274.
- [175] H. Trivedi, N. Balasubramanian, T. Khot, and A. Sabharwal, “Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, A. Rogers, J. L. Boyd-Graber, and N. Okazaki, Eds. Association for Computational Linguistics, 2023, pp. 10014–10037.
- [176] Z. Jiang, F. F. Xu, L. Gao, Z. Sun, Q. Liu, J. Dwivedi-

- Yu, Y. Yang, J. Callan, and G. Neubig, "Active retrieval augmented generation," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 7969–7992.
- [177] A. Asai, Z. Wu, Y. Wang, A. Sil, and H. Hajishirzi, "Self-rag: Learning to retrieve, generate, and critique through self-reflection," *CoRR*, vol. abs/2310.11511, 2023.
- [178] J. Liu, J. Jin, Z. Wang, J. Cheng, Z. Dou, and J. Wen, "RETA-LLM: A retrieval-augmented large language model toolkit," *CoRR*, vol. abs/2306.05212, 2023.
- [179] T. Vu, M. Iyyer, X. Wang, N. Constant, J. W. Wei, J. Wei, C. Tar, Y. Sung, D. Zhou, Q. V. Le, and T. Luong, "Freshllms: Refreshing large language models with search engine augmentation," *CoRR*, vol. abs/2310.03214, 2023.
- [180] X. Lyu, S. Grafberger, S. Biegel, S. Wei, M. Cao, S. Schelter, and C. Zhang, "Improving retrieval-augmented large language models via data importance learning," *CoRR*, vol. abs/2307.03027, 2023.
- [181] T. Gao, H. Yen, J. Yu, and D. Chen, "Enabling large language models to generate text with citations," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 6465–6488.
- [182] H. Luo, T. Zhang, Y. Chuang, Y. Gong, Y. Kim, X. Wu, H. Meng, and J. R. Glass, "Search augmented instruction learning," in *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 3717–3729.
- [183] X. V. Lin, X. Chen, M. Chen, W. Shi, M. Lomeli, R. James, P. Rodriguez, J. Kahn, G. Szilvasy, M. Lewis, L. Zettlemoyer, and S. Yih, "RA-DIT: retrieval-augmented dual instruction tuning," *CoRR*, vol. abs/2310.01352, 2023.
- [184] W. Yu, Z. Zhang, Z. Liang, M. Jiang, and A. Sabharwal, "Improving language models via plug-and-play retrieval feedback," *CoRR*, vol. abs/2305.14002, 2023.
- [185] Z. Feng, X. Feng, D. Zhao, M. Yang, and B. Qin, "Retrieval-generation synergy augmented large language models," *CoRR*, vol. abs/2310.05149, 2023.
- [186] S. Kadavath, T. Conerly, A. Askell, T. Henighan, D. Drain, E. Perez, N. Schiefer, Z. Hatfield-Dodds, N. DasSarma, E. Tran-Johnson, S. Johnston, S. E. Showk, A. Jones, N. Elhage, T. Hume, A. Chen, Y. Bai, S. Bowman, S. Fort, D. Ganguli, D. Hernandez, J. Jacobson, J. Kernion, S. Kravec, L. Lovitt, K. Ndousse, C. Olsson, S. Ringer, D. Amodei, T. Brown, J. Clark, N. Joseph, B. Mann, S. McCandlish, C. Olah, and J. Kaplan, "Language models (mostly) know what they know," *CoRR*, vol. abs/2207.05221, 2022.
- [187] Z. Jiang, J. Araki, H. Ding, and G. Neubig, "How can we know *When* language models know? on the calibration of language models for question answering," *Trans. Assoc. Comput. Linguistics*, vol. 9, pp. 962–977, 2021.
- [188] O. Press, M. Zhang, S. Min, L. Schmidt, N. A. Smith, and M. Lewis, "Measuring and narrowing the compositionality gap in language models," in *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 5687–5711.
- [189] O. Khattab, K. Santhanam, X. L. Li, D. Hall, P. Liang, C. Potts, and M. Zaharia, "Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive NLP," *CoRR*, vol. abs/2212.14024, 2022.
- [190] O. Yoran, T. Wolfson, B. Bogin, U. Katz, D. Deutch, and J. Berant, "Answering questions by meta-reasoning over multiple chains of thought," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 5942–5966.
- [191] M. A. Arefeen, B. Debnath, and S. Chakradhar, "Lean-context: Cost-efficient domain-specific question answering using llms," *CoRR*, vol. abs/2309.00841, 2023.
- [192] F. Xu, W. Shi, and E. Choi, "RECOMP: improving retrieval-augmented lms with compression and selective augmentation," *CoRR*, vol. abs/2310.04408, 2023.
- [193] Z. Wang, J. Araki, Z. Jiang, M. R. Parvez, and G. Neubig, "Learning to filter context for retrieval-augmented generation," *CoRR*, vol. abs/2311.08377, 2023.
- [194] J. Liu, L. Li, T. Xiang, B. Wang, and Y. Qian, "TCRA-LLM: token compression retrieval augmented large language model for inference cost reduction," in *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 9796–9810.
- [195] 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," *CoRR*, vol. abs/2307.03172, 2023.
- [196] R. Ren, Y. Wang, Y. Qu, W. X. Zhao, J. Liu, H. Tian, H. Wu, J. Wen, and H. Wang, "Investigating the factual knowledge boundary of large language models with retrieval augmentation," *CoRR*, vol. abs/2307.11019, 2023.
- [197] Y. Liu, S. Yavuz, R. Meng, M. Moorthy, S. Joty, C. Xiong, and Y. Zhou, "Exploring the integration strategies of retriever and large language models," *CoRR*, vol. abs/2308.12574, 2023.
- [198] R. Aksitov, C. Chang, D. Reitter, S. Shakeri, and Y. Sung, "Characterizing attribution and fluency tradeoffs for retrieval-augmented large language models," *CoRR*, vol. abs/2302.05578, 2023.
- [199] A. Mallen, A. Asai, V. Zhong, R. Das, D. Khashabi, and H. Hajishirzi, "When not to trust language models: Investigating effectiveness of parametric and non-parametric memories," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2023, Toronto,

- Canada, July 9-14, 2023, A. Rogers, J. L. Boyd-Graber, and N. Okazaki, Eds. Association for Computational Linguistics, 2023, pp. 9802–9822.
- [200] Y. Wang, X. Ma, and W. Chen, “Augmenting black-box llms with medical textbooks for clinical question answering,” *CoRR*, vol. abs/2309.02233, 2023.
- [201] S. Munikoti, A. Acharya, S. Wagle, and S. Horawalavithana, “ATLANTIC: structure-aware retrieval-augmented language model for interdisciplinary science,” *CoRR*, vol. abs/2311.12289, 2023.
- [202] X. Li, E. Nie, and S. Liang, “Crosslingual retrieval augmented in-context learning for bangla,” *CoRR*, vol. abs/2311.00587, 2023.
- [203] A. Lozano, S. L. Fleming, C. Chiang, and N. Shah, “Clinfo.ai: An open-source retrieval-augmented large language model system for answering medical questions using scientific literature,” *CoRR*, vol. abs/2310.16146, 2023.
- [204] B. Zhang, H. Yang, T. Zhou, A. Babar, and X. Liu, “Enhancing financial sentiment analysis via retrieval augmented large language models,” in *4th ACM International Conference on AI in Finance, ICAIF 2023, Brooklyn, NY, USA, November 27-29, 2023*. ACM, 2023, pp. 349–356.
- [205] A. Louis, G. van Dijck, and G. Spanakis, “Interpretable long-form legal question answering with retrieval-augmented large language models,” *CoRR*, vol. abs/2309.17050, 2023.
- [206] G. Zyskind, T. South, and A. Pentland, “Don’t forget private retrieval: distributed private similarity search for large language models,” *CoRR*, vol. abs/2311.12955, 2023.
- [207] W. Jiang, M. Zeller, R. Waleffe, T. Hoefler, and G. Alonso, “Chameleon: a heterogeneous and disaggregated accelerator system for retrieval-augmented language models,” *CoRR*, vol. abs/2310.09949, 2023.
- [208] Y. Hoshi, D. Miyashita, Y. Ng, K. Tatsuno, Y. Morioka, O. Torii, and J. Deguchi, “Ralle: A framework for developing and evaluating retrieval-augmented large language models,” in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023 - System Demonstrations, Singapore, December 6-10, 2023*, Y. Feng and E. Lefever, Eds. Association for Computational Linguistics, 2023, pp. 52–69.
- [209] R. Thoppilan, D. D. Freitas, J. Hall, N. Shazeer, A. Kulshreshtha, H. Cheng, A. Jin, T. Bos, L. Baker, Y. Du, Y. Li, H. Lee, H. S. Zheng, A. Ghafouri, M. Menegali, Y. Huang, M. Krikun, D. Lepikhin, J. Qin, D. Chen, Y. Xu, Z. Chen, A. Roberts, M. Bosma, Y. Zhou, C. Chang, I. Krivokon, W. Rusch, M. Pickett, K. S. Meier-Hellstern, M. R. Morris, T. Doshi, R. D. Santos, T. Duke, J. Soraker, B. Zevenbergen, V. Prabhakaran, M. Diaz, B. Hutchinson, K. Olson, A. Molina, E. Hoffman-John, J. Lee, L. Aroyo, R. Rajakumar, A. Butryna, M. Lamm, V. Kuzmina, J. Fenton, A. Cohen, R. Bernstein, R. Kurzweil, B. A. y Arcas, C. Cui, M. Croak, E. H. Chi, and Q. Le, “Lamda: Language models for dialog applications,” *CoRR*, vol. abs/2201.08239, 2022.
- [210] K. Shuster, M. Komeili, L. Adolphs, S. Roller, A. Szlam, and J. Weston, “Language models that seek for knowledge: Modular search & generation for dialogue and prompt completion,” in *Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, Y. Goldberg, Z. Kozareva, and Y. Zhang, Eds. Association for Computational Linguistics, 2022, pp. 373–393.
- [211] X. Liu, H. Lai, H. Yu, Y. Xu, A. Zeng, Z. Du, P. Zhang, Y. Dong, and J. Tang, “Webglm: Towards an efficient web-enhanced question answering system with human preferences,” in *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2023, Long Beach, CA, USA, August 6-10, 2023*, A. K. Singh, Y. Sun, L. Akoglu, D. Gunopulos, X. Yan, R. Kumar, F. Ozcan, and J. Ye, Eds. ACM, 2023, pp. 4549–4560.
- [212] I. Gur, H. Furuta, A. Huang, M. Safdari, Y. Matsuo, D. Eck, and A. Faust, “A real-world webagent with planning, long context understanding, and program synthesis,” *CoRR*, vol. abs/2307.12856, 2023.
- [213] J. Menick, M. Trebacz, V. Mikulik, J. Aslanides, H. F. Song, M. J. Chadwick, M. Glaese, S. Young, L. Campbell-Gillingham, G. Irving, and N. McAleese, “Teaching language models to support answers with verified quotes,” *CoRR*, vol. abs/2203.11147, 2022.
- [214] X. Shi, J. Liu, Y. Liu, Q. Cheng, and W. Lu, “Know where to go: Make LLM a relevant, responsible, and trustworthy searcher,” *CoRR*, vol. abs/2310.12443, 2023.
- [215] Y. Qin, Z. Cai, D. Jin, L. Yan, S. Liang, K. Zhu, Y. Lin, X. Han, N. Ding, H. Wang, R. Xie, F. Qi, Z. Liu, M. Sun, and J. Zhou, “Webcpm: Interactive web search for chinese long-form question answering,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, A. Rogers, J. L. Boyd-Graber, and N. Okazaki, Eds. Association for Computational Linguistics, 2023, pp. 8968–8988.
- [216] X. Deng, Y. Gu, B. Zheng, S. Chen, S. Stevens, B. Wang, H. Sun, and Y. Su, “Mind2web: Towards a generalist agent for the web,” *CoRR*, vol. abs/2306.06070, 2023.
- [217] S. Yao, H. Chen, J. Yang, and K. Narasimhan, “Webshop: Towards scalable real-world web interaction with grounded language agents,” in *NeurIPS*, 2022.
- [218] S. Zhou, F. F. Xu, H. Zhu, X. Zhou, R. Lo, A. Sridhar, X. Cheng, Y. Bisk, D. Fried, U. Alon, and G. Neubig, “Webarena: A realistic web environment for building autonomous agents,” *CoRR*, vol. abs/2307.13854, 2023.
- [219] R. Lo, A. Sridhar, F. F. Xu, H. Zhu, and S. Zhou, “Hierarchical prompting assists large language model on web navigation,” in *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 10217–10244.
- [220] S. MacAvaney, C. Macdonald, R. Murray-Smith, and I. Ounis, “Intent5: Search result diversification using causal language models,” *CoRR*, vol. abs/2108.04026, 2021.

- [221] N. Craswell, "Mean reciprocal rank," in *Encyclopedia of Database Systems*, L. Liu and M. T. Özsu, Eds. Springer US, 2009, p. 1703.
- [222] K. Järvelin and J. Kekäläinen, "Cumulated gain-based evaluation of IR techniques," *ACM Trans. Inf. Syst.*, vol. 20, no. 4, pp. 422–446, 2002.
- [223] K. Papineni, S. Roukos, T. Ward, and W. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA*. ACL, 2002, pp. 311–318.
- [224] C.-Y. Lin, "ROUGE: A package for automatic evaluation of summaries," in *Text Summarization Branches Out*. Barcelona, Spain: Association for Computational Linguistics, Jul. 2004, pp. 74–81.
- [225] P. Manakul, A. Liusie, and M. J. F. Gales, "Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models," *CoRR*, vol. abs/2303.08896, 2023.
- [226] H. Qian, Y. Zhu, Z. Dou, H. Gu, X. Zhang, Z. Liu, R. Lai, Z. Cao, J. Nie, and J. Wen, "Webbrain: Learning to generate factually correct articles for queries by grounding on large web corpus," *CoRR*, vol. abs/2304.04358, 2023.
- [227] J. Li, X. Cheng, W. X. Zhao, J. Nie, and J. Wen, "Halueval: A large-scale hallucination evaluation benchmark for large language models," *CoRR*, vol. abs/2305.11747, 2023.
- [228] L. Chen, Y. Deng, Y. Bian, Z. Qin, B. Wu, T. Chua, and K. Wong, "Beyond factuality: A comprehensive evaluation of large language models as knowledge generators," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Association for Computational Linguistics, 2023, pp. 6325–6341.
- [229] S. Xu, D. Hou, L. Pang, J. Deng, J. Xu, H. Shen, and X. Cheng, "Ai-generated images introduce invisible relevance bias to text-image retrieval," *CoRR*, vol. abs/2311.14084, 2023.
- [230] S. Dai, Y. Zhou, L. Pang, W. Liu, X. Hu, Y. Liu, X. Zhang, and J. Xu, "Llms may dominate information access: Neural retrievers are biased towards llm-generated texts," *CoRR*, vol. abs/2310.20501, 2023.
- [231] J. S. Park, J. C. O'Brien, C. J. Cai, M. R. Morris, P. Liang, and M. S. Bernstein, "Generative agents: Interactive simulacra of human behavior," *CoRR*, vol. abs/2304.03442, 2023.