



# Trends in Integration of Knowledge and Large Language Models: A Survey and Taxonomy of Methods, Benchmarks, and Applications

Zhangyin Feng<sup>1\*</sup>, Weitao Ma<sup>1\*</sup>, Weijiang Yu<sup>2✉</sup>, Lei Huang<sup>1</sup>, Haotian Wang<sup>1</sup>,  
Qianglong Chen<sup>2</sup>, Weihua Peng<sup>2</sup>, Xiaocheng Feng<sup>1✉</sup>, Bing Qin<sup>1</sup>, Ting Liu<sup>1</sup>

<sup>1</sup>Harbin Institute of Technology, Harbin, China

<sup>2</sup>Huawei Inc., Shenzhen, China

{zyfeng, wtma, lhuang, xcfeng, qinb, tliu}@ir.hit.edu.cn

{weijiangyu8, wanght1998, chenqianglong.ai, pengwh.hit}@gmail.com

## Abstract

Large language models (LLMs) exhibit superior performance on various natural language tasks, but they are susceptible to issues stemming from outdated data and domain-specific limitations. In order to address these challenges, researchers have pursued two primary strategies, knowledge editing and retrieval augmentation, to enhance LLMs by incorporating external information from different aspects. Nevertheless, there is still a notable absence of a comprehensive survey. In this paper, we propose a review to discuss the trends in integration of knowledge and large language models, including taxonomy of methods, benchmarks, and applications. In addition, we conduct an in-depth analysis of different methods and point out potential research directions in the future. We hope this survey offers the community quick access and a comprehensive overview of this research area, with the intention of inspiring future research endeavors.

## 1 Introduction

Large language models (LLMs) have demonstrated an impressive ability to encode real-world knowledge in their parameters and a remarkable capacity for solving various natural language processing tasks (Brown et al., 2020; Hoffmann et al., 2022; Zeng et al., 2022; Chowdhery et al., 2022; Touvron et al., 2023; Zhao et al., 2023b). However, they still suffer from serious challenges in knowledge-intensive tasks (Petroni et al., 2021), which require a substantial volume of real-world knowledge.

Recent works show that LLMs struggle to learn long-tail knowledge (Kandpal et al., 2023; Mallen et al., 2023), are not able to update their parameters in time to capture the changing world (De Cao et al., 2021; Kasai et al., 2022) (i.e., the parameters of ChatGPT<sup>1</sup> only contain information be-

fore September 2021, and are completely unaware of the latest world knowledge.), and suffer from hallucinations (Zhang et al., 2023a; Rawte et al., 2023; Huang et al., 2023a). To alleviate these problems, there has been growing interest in integrating knowledge and large language models through knowledge editing or retrieval augmentation. Knowledge editing (De Cao et al., 2021; Sinitin et al., 2020) aims to modify obsolete knowledge in LLMs using an efficient method that only updates partial model parameters. Retrieval augmentation (Mallen et al., 2023; Shi et al., 2023; Trivedi et al., 2023) employs an off-the-shelf retrieval model to fetch relevant documents from an external corpus to aid large language models and maintains their parameters unchanged. Numerous works have been proposed to integrate knowledge and large language models, focusing on the aforementioned two aspects. Nevertheless, these endeavors remain rather fragmented, lacking a comprehensive and systematic review.

To fill the gap, in this paper, we present a concrete organization of our survey, focusing on knowledge editing and retrieval augmentation, as depicted in Figure 1. We begin by systematically introducing the knowledge editing methods according to the processed structure of the model (§2), including input editing (§2.1), model editing (§2.2) and assess knowledge editing (§2.3) which covers representative methods and general benchmarks. Furthermore, we provide a detailed discussion of retrieval augmentation (§3), including retrieval judgement (§3.1), document retrieval (§3.2), document utilization (§3.3), knowledge conflict (§3.4) and benchmark (§3.5). Then, we summarize some cutting-edge applications of integration of knowledge and large language models (§4), such as New Bing<sup>2</sup>. Finally, to stimulate further research in this field, we offer insights into prospective directions for future investigations (§5).

\* means Equal Contribution, ✉ means Corresponding Author

<sup>1</sup><https://chat.openai.com>

<sup>2</sup><https://www.bing.com/new>

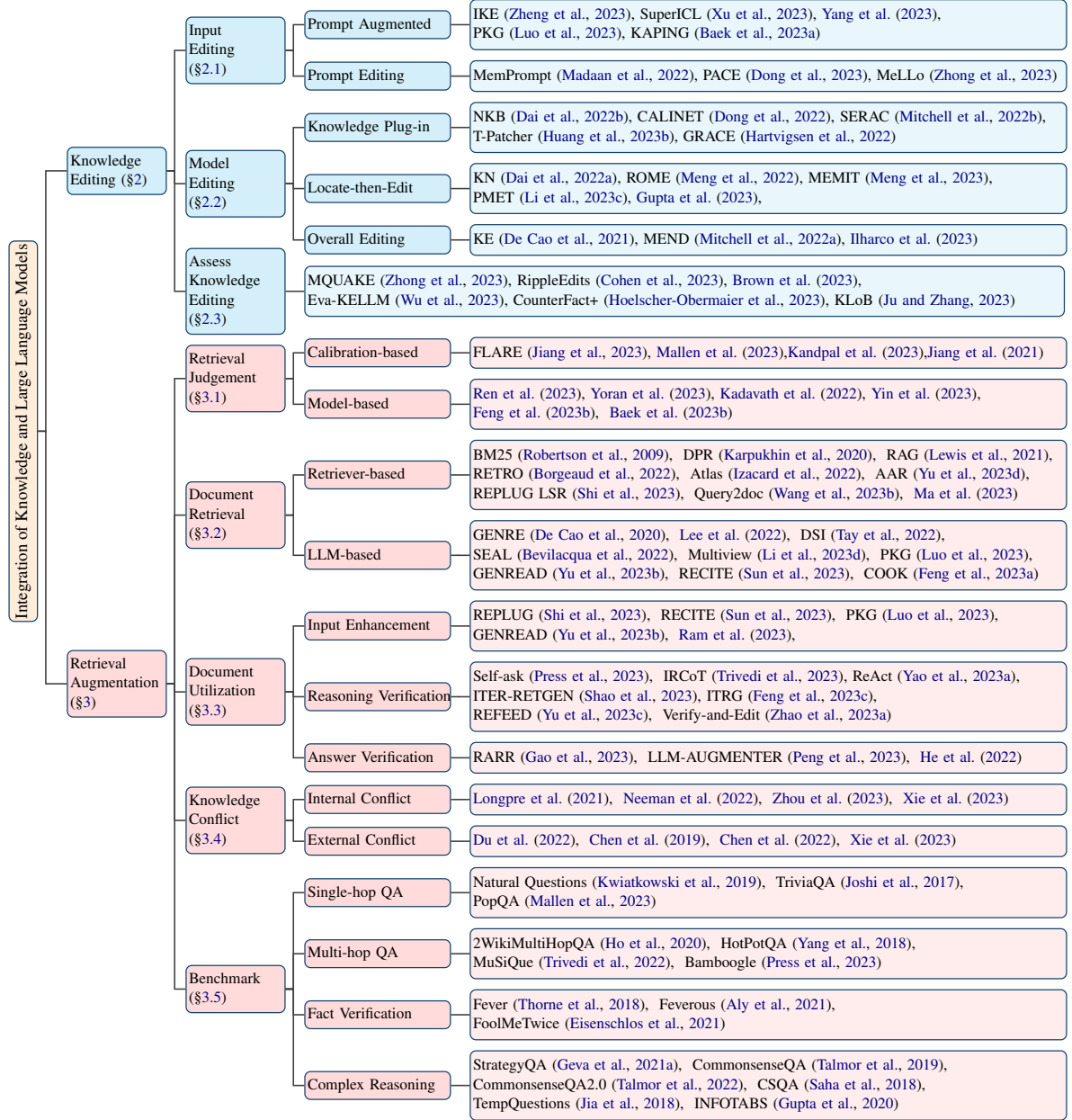


Figure 1: Taxonomy of trends in integration of knowledge and large language models.

**Related Work** Some previous reviews also discussed the interaction between language models and external knowledge. Hu et al. (2023) deliberates on the various forms of knowledge and methods employed to augment language models in prior research, which only focuses on pre-trained language models with smaller parameter sizes. Yao et al. (2023b) specializes in a comprehensive and empirical discussion on existing knowledge editing methods which does not involve other retrieval-related knowledge application methods. Zhang et al. (2023b) offers a broad and comprehensive overview of the strategies for adapting LLMs to dynamically updated world knowledge, while we pro-

vide a deep and detailed analysis of current methods and benchmarks. Other existing knowledge-enhanced LMs surveys (Yin et al., 2022; Yu et al., 2022; Yang et al., 2023) center on pre-trained language models and frequently employ re-training methods to incorporate additional knowledge that is not suitable for the current LLMs. We will comprehensively introduce the recent advancements in integration of knowledge and large language models, with a specific focus on two effective methods: knowledge editing and retrieval augmentation through providing a detailed analysis and offering insights into future developments.

## 2 Knowledge Editing

Knowledge editing is an emerging method for rectifying inaccuracies and updating outdated information in LLMs through the incorporation of new knowledge. In this section, we delve into the current works about knowledge editing, with a specific focus on the processed structure of LLMs across various methods. As shown in Figure 2, we organize them into three categories: input editing (§2.1), model editing (§2.2), and assess knowledge editing (§2.3).

### 2.1 Input Editing

The extensive parameter scale and "black-box" form of numerous large models often impede them from undergoing fine-tuning commonly for the acquisition of new knowledge, such as ChatGPT, Bard<sup>3</sup>. Therefore, the most direct approach to infusing knowledge into LLMs involves editing the input (Zheng et al., 2023; Luo et al., 2023), which incurs minimal costs and resource requirements. There are two aspects of input editing: the inclusion of external information to enhance prompts, and editing prompts based on feedback. Adjusting input not only offers an intuitive and comprehensible depiction of the process of new knowledge but also guarantees the preservation of the original model’s knowledge.

**Prompt Augmented.** In-context learning (ICL) has been proven a useful paradigm for LLMs based on a few demonstrations that are tailored for diverse tasks. Based on the goal of knowledge editing, IKE (Zheng et al., 2023) designs three different types of prompts: copy, update, and retain to enhance the generalization of newly injected knowledge and keep original information in LLMs. SuperICL (Xu et al., 2023) introduces a smaller model fine-tuned on task-specific data as a plug-in for LLMs to augment the conventional ICL method. The smaller model produces the predicted label and confidence score for each original demonstration which assists the LLM in grasping the complexity of the given examples. During the inference phase, the LLM generates a final prediction for the test input which comprises the reconstructed context, input text, and plug-in model’s prediction. Luo et al. (2023) similarly train a smaller model as an auxiliary parametric knowledge guiding (PKG) framework to generate background documents for domain-specific

tasks enhancing the prompt. In addition, leveraging knowledge graphs (KGs) to augment pre-trained language models has become a prevalent practice (Yang et al., 2023; Moiseev et al., 2022), some methods expand the input using knowledge graphs as an external resource without training. For example, KAPING (Baek et al., 2023a) retrieves the relevant facts from KG and integrates them with the original question as a new prompt to generate the answer. Andrus et al. (2022) propose an architecture aiming at enhancing story comprehension for LLMs. Initially, it extracts entities from the story document to construct KGs. Subsequently, question-related triples are also retrieved and transformed into natural language sentences to create an informative prompt.

*Highlight:* These methods mainly concentrate on incorporating new facts into prompts using diverse approaches, which are relatively straightforward and highly achievable. However, they exhibit a limited capacity to correct errors in LLMs.

**Prompt Editing.** In addition to introducing background information into the prompt, the process of deconstructing and refining the prompt can also contribute to enhancing the accuracy of the LLMs’ responses. MemPrompt (Madaan et al., 2022) is designed by incorporating user feedback to mitigate instances where the LLM misunderstood prompts. It establishes an expanding memory to retain users’ past assessments of the model’s understanding of the task. Subsequently, the architecture produces an improved prompt by combining it with user feedback retrieved from this memory resource. In contrast to the utilization of user feedback, PACE (Dong et al., 2023) leverages LLMs’ own judgment to criticize the response and refine the prompt, which successfully boosts the performance of medium/low-quality human-written prompts. To enhance the generality of injected knowledge for multi-hop questions, MeLLO (Zhong et al., 2023) decomposes the original question into sub-questions. This approach based on a self-checking mechanism, empowers the LLM to adapt the output of the sub-questions according to the retrieved facts from an explicit memory.

*Highlight:* It is noteworthy that editing prompts can facilitate LLMs in accurately tackling complicated reasoning tasks and rectifying errors, making it particularly suitable in "black-box" scenarios. Nevertheless, these methods introduce facts in a disposable manner, without fundamentally modify-

<sup>3</sup><https://bard.google.com/>

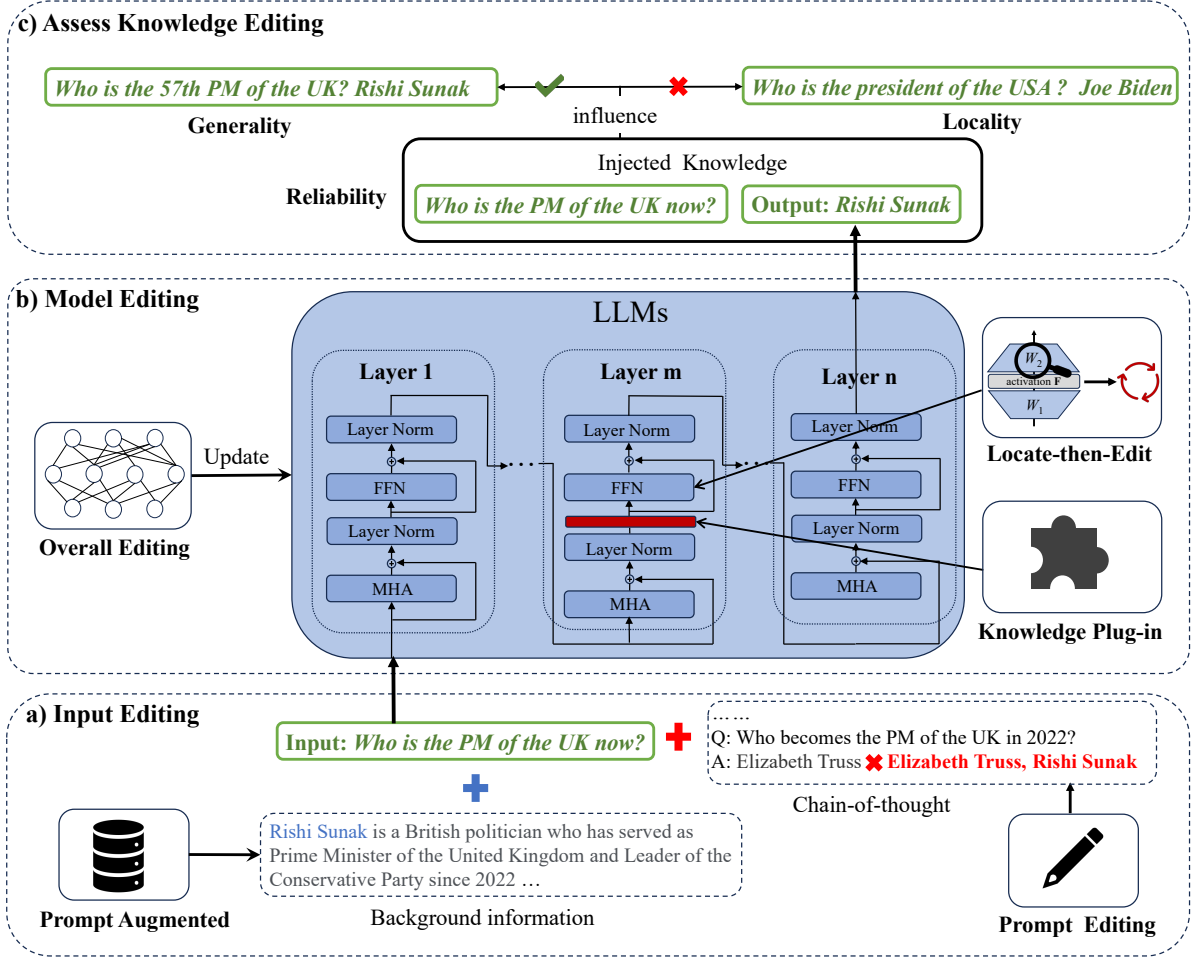


Figure 2: An overview of knowledge editing from three perspectives, which encompasses input (input editing), processing (model editing), and output (assess knowledge editing) aspects of the model’s processing structure.

ing the intrinsic knowledge within LLMs.

## 2.2 Model Editing

Instead of editing input, numerous works are dedicated to fine-grained model editing in a parametric fashion, which can ensure the persistence of injected knowledge. According to different operations targeting LLMs’ parameters, we categorize them into three classes, namely knowledge plug-in, locate-then-edit, and overall editing.

**Knowledge Plug-in.** This paradigm introduces external plug-in parametric knowledge to edit LLMs, without altering the original weights. Following the perspective that knowledge resides within MLP layers (Geva et al., 2021b), NKB (Dai et al., 2022b) and CALINET (Dong et al., 2022) both adjust the output of the feed-forward-network (FFN) by adding additional memory slots. SERAC (Mitchell et al., 2022b) maintains a scope classifier that determines whether the inputs fall within the scope of each edit stored in an external edit mem-

ory. Inputs deemed to be in scope will be passed to a counterfactual model along with the associated edit examples, while out-of-scope inputs will be processed by the original model. In addition to the traditional one-step editing scenario, plug-in knowledge is also applied to sequential model editing tasks suitable for practical situations. T-Patcher (Huang et al., 2023b) incorporates an extra trainable patch into the last FFN layer for each mistake. This mechanism allows the model to train these patches to correct errors while keeping the original parameters frozen. GRACE (Hartvigsen et al., 2022) implements sequential editing by adding an adapter at a selected layer. This adapter contains a discrete codebook designed to map activations to the corresponding values, along with a deferral mechanism that achieves codebook updates to encourage similar edits with the same values. The frame enables efficient model editing with low impact on unrelated inputs.

*Highlight:* These methods introduce plug-in knowl-



edge to quickly complete model editing in lower resource scenes, but they often cannot achieve precise editing and the added objects are mostly in the last FFN layer of the model.

**Locate-then-Edit.** This paradigm employs a dual-stage framework to edit the model in a fine-grained manner. It first locates the specific parameters about new facts and makes targeted modifications to inject this knowledge. Dai et al. (2022a) adapt a method on integrated gradients to evaluate the contribution of each neuron in the second linear layer of FFN to knowledge predictions and then selects the knowledge neurons (KN) to update. ROME (Meng et al., 2022) focuses on GPT-like autoregressive models, which introduces corruption to the embeddings of the subject token and successively restore internal activations to their clean value to acquire the causal importance for MLP layers. This approach will compute the new key-value pair vector which is inserted into the original matrix in specific FFN for each fact to edit the model precisely. However, most model editing methods can only inject a single piece of knowledge at a time. MEMIT (Meng et al., 2023) based on ROME defines a new expansion objective solved by the normal equations for updating thousands of edits at once. Instead of modifying a single layer, MEMIT spreads updates from a set of key-value pairs over the identified layers to improve the robustness. As finding multi-head self-attention (MHSA) works as a knowledge extractor in transformer component, PMET (Li et al., 2023c) concurrently optimizes the hidden states of MHSA and FFN, surpassing the performance of MEMIT. Moreover, in accordance with the specific attributes of distinct knowledge domains, the model editing approach is tailored accordingly. Gupta et al. (2023) adapt the MEMIT method to edit commonsense knowledge by broadening the scope of corruption tokens to encompass subjects, objects, and verbs.

*Highlight:* These methods require an additional step to identify the parameters that store knowledge about edits. However, Hase et al. (2023) demonstrates that these knowledge-storage parameters are not exactly the same as the neurons that can be efficiently adjusted to edit the model. Therefore, the selection of appropriate parameters plays a significant role in model editing.

**Overall Editing.** This paradigm involves editing the model by directly modifying it, bypassing the

need to locate the specific parameters where knowledge is stored. A common approach is to train a hyper-network based on meta-learning to learn the model weights on new facts. Knowledge editor (KE) (De Cao et al., 2021) trains a bidirectional-LSTM as a hyper-network with constrained optimizations, to predict the updates required for injecting each atomic fact. However, this method cannot be applied to large models and computationally infeasible when editing knowledge. MEND (Mitchell et al., 2022a) introduces a low-rank decomposition of the gradients to optimize the process of hyper-network learning the parameter update from standard fine-tuning. This approach is capable of training small auxiliary networks with limited resources which effectively edits the behavior of LLMs. In addition to meta-learning, Ilharco et al. (2023) propose a method to simply edit models through task arithmetic. The framework illustrates that by combining task vectors derived from the element-wise difference between parameters after pre-training and fine-tuning, becomes possible to control the performance of models on various tasks.

*Highlight:* These methods are careless about the internals of the model, and obtain the update of model editing in a data-driven manner, without being constrained to a particular model. Nevertheless, in the case of LLMs, these approaches often require many additional cost, particularly in the context of hyper-network methods, which demand their own parameter size reduction in careful design.

### 2.3 Assess Knowledge Editing

After editing the input and model, assessing the extent of knowledge integration can be accomplished by scrutinizing the output. This subsection will primarily introduce the characteristics of model evaluation and provide an overview of general benchmarks for knowledge editing in Table 1.

The current methods for editing knowledge mainly aim at incorporating triple-fact knowledge, which concentrate on question-answer (QA) tasks, i.e., ZsRE (Levy et al., 2017). In addition, CounterFact is an evaluation dataset specially constructed for knowledge editing tasks to measure the effectiveness of significant changes in comparison to merely superficial alterations of target words (Meng et al., 2022). There are three main properties of assessing knowledge editing accounting for *reliability*, *generality*, *locality* (Yao et al., 2023b; Huang et al., 2023b).

DataSet	Language	Size	Evaluation
ZsRE (Levy et al., 2017)	En	182,282	Reliability/Generality/Locality
Counterful (Meng et al., 2022)	En	21,919	Reliability/Generality/Locality
Counterful+ (Hoelscher-Obermaier et al., 2023)	En	21,919	Reliability/Generality/Locality
Bi-ZsRE (Wang et al., 2023a)	En&Zh	14037&14037	Reliability/Generality/Locality/Portability/Cross-Lingual
MQUAKE (Zhong et al., 2023)	En	11,043	Generality
RippleEdits (Cohen et al., 2023)	En	4,000	Reliability/Generality/Locality
Eva-KELLM (Wu et al., 2023)	En&Zh	8,882&6,930	Reliability/Generality/Locality/Cross-Lingual

Table 1: Comparison of knowledge editing evaluation benchmarks, including language, size and evaluation.

**Reliability.** The primary objective of the post-edited model is to generate the desired predictions for edited input (De Cao et al., 2021). The reliability which is directly evaluated using the datasets mentioned, is assessed based on the success rate in knowledge editing (Levy et al., 2017; Meng et al., 2022).

**Generality.** The generality is reflected in the fact that the post-edited model can successfully update the relevant facts of the edits, which includes numerous aspects. In addition to fundamental semantic transcription like sentence rephrasing and back translation, there are other more comprehensive assessment methods. Yao et al. (2023b) introduce a new metric called *portability* which assesses the robustness of generalization in order to ascertain whether the post-edit model merely memorizes the superficial alterations in wording. MQUAKE serves as a challenging benchmark to evaluate whether the edited model can correctly and synchronously update the answers to multi-hop questions related to the edits (Zhong et al., 2023). RippleEdits also focuses on the ripple effects of knowledge editing based on six evaluation criteria, most of which assess the generality of injected knowledge within a 2-hop distance from edited facts (Cohen et al., 2023).

**Locality.** Determining whether the edited model retains unrelated knowledge of the edited fact is also a crucial property of the evaluation process. There are also two criteria to evaluate the locality of ripple effects of knowledge editing in RippleEdits (Cohen et al., 2023). CounterFact+ is a more sensitive benchmark extending the original Counterfact to detect the unwanted side effects brought

by knowledge editing, which influences unrelated facts and the next token probability distribution (Hoelscher-Obermaier et al., 2023).

Besides the aforementioned evaluation methods, several studies have raised attention towards examining a special kind of generalization: the cross-lingual capabilities of knowledge editing methods. Bi-ZsRE (Wang et al., 2023a) constructs samples in both the Chinese and English languages, while also assessing traditional editing methods. The findings indicate that current methods struggle to effectively transfer edited knowledge from one language to another. In addition, Eva-KELLM (Wu et al., 2023) also evaluates the cross-language capabilities of knowledge editing, which pioneers the direct editing of raw documents for greater convenience.

*Highlight:* Aside from assessing the based effects, there are profound analyses to comprehensively evaluate knowledge editing methods. KLoB (Ju and Zhang, 2023) meticulously emphasizes knowledge-locating methods based on three crucial assessing criteria: consistency, relevance, and unbiasedness. Brown et al. (2023) explore the robustness of the post-edited model, which experimentally identifies the degrades of general robustness from knowledge editing. These thorough evaluations facilitate the selection of appropriate editing methods in real-world deployment scenarios.

### 3 Retrieval Augmentation

As discussed in section §2, knowledge editing (De Cao et al., 2021) is a productive approach to updating outdated knowledge by modifying the parameters of a specific part of large language models. However, knowledge editing faces some other

problems. Firstly, it is not entirely clear how and where knowledge is stored in large language models. Secondly, the mapping relationship between knowledge and parameters is very complicated, and modifying the parameters corresponding to some knowledge may affect other knowledge. In this section, we introduce retrieval augmentation, an alternative route to integrate knowledge and large language models while keeping the parameters unchanged.

Unlike knowledge editing, which mainly parameterizes external knowledge to update the large language models, retrieval augmentation utilizes external knowledge in a non-parameterized form in the inference stage. Retrieval augmentation typically consists of a retriever and a large language model. Given an input context, the retriever first fetches relevant documents from an external corpus. Then, we can use relevant documents at different stages to improve the performance of large language models. In this section, we focus on the following key questions of retrieval augmentation:

- When do large language models need to be enhanced by retrieval? (§3.1)
- How to retrieve relevant documents? (§3.2)
- How do large language models utilize retrieved documents? (§3.3)
- How to resolve knowledge conflicts from different documents? (§3.4)

### 3.1 Retrieval Judgement

A very important problem for retrieval-augmented large language models is to know the knowledge boundaries (Yin et al., 2023) of LLMs and determine when to retrieve supplementary knowledge. The current retrieval judgment methods are mainly divided into two categories: calibration-based judgment and model-based judgment.

**Calibration-based.** A simple and intuitive idea is to set a metric and a threshold. When the metric is above or below the threshold, we trigger the retriever to fetch relevant documents. Kandpal et al. (2023) study the relationship between the knowledge memorized by large language models and the information in pre-training datasets. Their results demonstrate the strong correlational relationship between accuracy and relevant document count for numerous question-answering datasets. In order to deeply analyze the relationship between

the parameterized knowledge of LLMs and the data popularity, Mallen et al. (2023) build an open domain question answering dataset PopQA, which contains entity popularity from Wikipedia. Then, they devise an adaptive retrieval method to only use retrieval for questions whose popularity is lower than the popularity threshold. In addition to popularity, Jiang et al. (2021) show that LLMs tend to be well-calibrated and low probability or confidence often indicates a lack of relevant knowledge. Si et al. (2022) and Manakul et al. (2023) utilize the token probability to indicate the uncertainty of their output. Following this idea, Jiang et al. (2023) propose a confidence-based active retrieval approach, named FLARE. If the confidence of each word in the generated sentence is above the threshold, they accept the sentence without retrieving additional information. Otherwise, they actively trigger retrieval and utilize the retrieved relevant information to regenerate the current sentence.

**Model-based.** There are two kinds of model-based judgment methods: *normal setting* and *retrieval setting*. The *normal setting* is to directly determine whether to trigger retrieval based on the question. Considering that large language models have very powerful capabilities, some researchers directly employ large language models to determine whether retrieval is needed. Yin et al. (2023) investigate the self-knowledge of LLMs by assessing their ability to identify unanswerable or unknowable questions. Kadavath et al. (2022) prompt LLMs to predict the probabilities of whether their responses are reliable. Those unreliable responses indicate that LLMs require additional information to answer the corresponding questions. Feng et al. (2023b) ask LLMs “Do you need more information? (Yes or No)” to determine whether external knowledge is needed for the given question through in-context learning. Ren et al. (2023) adopt priori and posteriori judgment instructions to investigate whether LLMs are capable of perceiving their own factual knowledge boundary for both normal setting and retrieval setting. Priori judgment asks LLMs whether they can provide an answer to the question. Posteriori judgment asks LLMs to evaluate the correctness of the answer to the question. They observe that LLMs perceive their factual knowledge boundary inaccurately and have a tendency to be overconfident in normal setting.

The *retrieval setting* is to first retrieve the relevant documents for all questions, and then judge

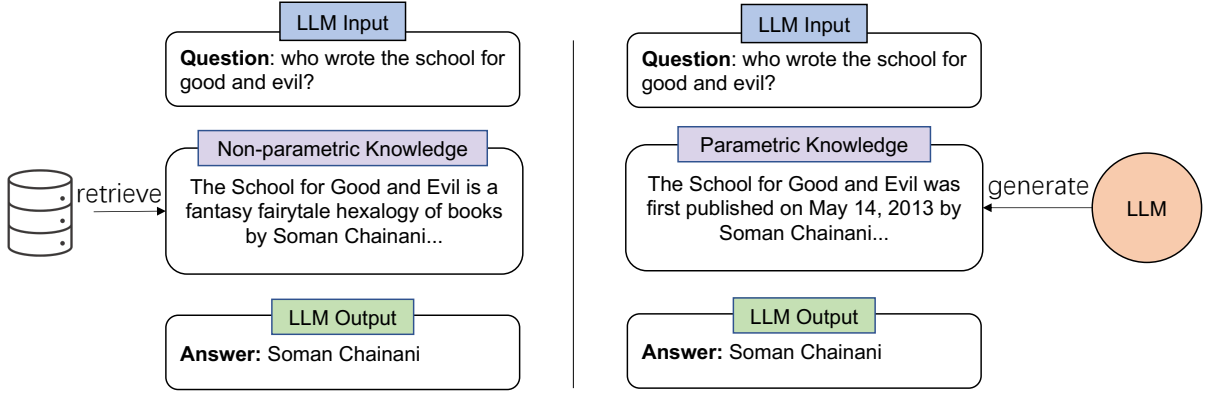


Figure 3: Two kinds of document retrieval methods: the document on the left is fetch from an external corpus with a retriever, and the document on the right is generated by a large language model.

whether the relevant documents can answer the question. If relevant, large language models utilize the retrieved documents to generate the answer. If irrelevant, large language models directly generate the answer. Yoran et al. (2023) regard the judgment of the relevance of retrieved documents and questions as a natural language inference (NLI) problem (Dagan et al., 2005; Bowman et al., 2015) and use a well-trained BART-Large (Lewis et al., 2019) NLI model to identify irrelevant retrieved documents. The retrieved documents serve as the premise, while the question and generated answer are concatenated and serve as hypothesis. Baek et al. (2023b) propose to fine-tune LLMs with instruction data to identify the relatedness between the input question and the retrieved documents, and ensemble results from various instructions to further improve accuracy. Ren et al. (2023) find that the accuracy of LLMs’ self-assessment improves after incorporating relevant documents and it is effective to dynamically introduce retrieved documents for LLMs.

*Highlight:* Calibration-based methods are simple and effective. However, these methods are mostly ad-hoc, while the score may not be available for commercial LLMs. Additionally, it may be challenging to find appropriate thresholds for different data and models due to the sensitivity of thresholds. Judging whether to trigger retrieval with a well-trained model is a promising direction. At present, this direction is in an early stage, and there is still a large space for exploration.

### 3.2 Document Retrieval

As shown in Figure 3, there are two ways to get the relevant documentation. One approach is to use a

retriever to fetch relevant documents from an external corpus (e.g. Wikipedia). Another approach is to use a large language model to generate relevant documents.

**Retriever-based.** Given an input context  $x$ , a retriever aims to retrieve a small set of documents from a corpus  $\mathcal{D} = \{d_1 \dots d_m\}$  that are relevant to  $x$ . There are different types of searchers, including term-based sparse retriever, embedding-based dense retriever and commercial search engines. Sparse retriever is usually implemented using TF-IDF or BM25 (Robertson et al., 2009), which matches keywords efficiently with an inverted index. However, term-matching methods are sensitive to highly selective keywords and phrases. Dense retriever (Karpukhin et al., 2020) encodes text into a continuous dense semantic space, where synonyms or paraphrases that consist of completely different tokens may still be mapped to vectors close to each other. Commercial search engines, such as Google and Baidu, are complex systems that are able to retrieve the latest world knowledge. These three methods have different advantages and application scenarios, and we will focus on dense retriever next.

Given a collection of text passages, the goal of the dense retriever is to index all the passages in a low-dimensional and continuous space, such that it can retrieve efficiently the top  $k$  passages relevant to the input question for the reader at runtime. Dense retriever uses a dense encoder  $E(\cdot)$  which maps any text passage to a  $d$ -dimensional real-valued vectors. Specifically, the encoder maps each document  $d \in \mathcal{D}$  to an embedding  $E(d)$  by taking the mean pooling of the last hidden represen-



tation over the tokens in  $d$ . At query time, the same encoder is applied to the input context  $q$  to obtain a query embedding  $E(q)$ . The similarity between the query embedding and the document embedding is computed by their cosine similarity:

$$s(d, q) = \cos(E(d), E(q)) \quad (1)$$

The top- $k$  documents that have the highest similarity scores when compared with the input  $q$  are retrieved in this step.

Prior works explore different ways to train the whole retriever-LM system in an end-to-end fashion, using retrieval augmented sequence log-likelihood (Lewis et al., 2021; Borgeaud et al., 2022), fusion-in-decoder attention distills (Izacard and Grave, 2022; Izacard et al., 2022), or knowledge graph (Ju et al., 2022). This kind of fine-tuning can be expensive when more and more unique demands emerge (Maronikolakis and Schütze, 2021). More importantly, many toptier LMs can only be accessed through black-box APIs (Ouyang et al., 2022; OpenAI, 2023). These APIs allow users to submit queries and receive responses but typically do not support fine-tuning. In contrast to prior work that adapts language models to the retriever, recent work attempts to adapt the retriever to language models. REPLUG LSR (Shi et al., 2023) further improves the initial retrieval model in REPLUG with supervision signals from a black-box language model, i.e. GPT-3 Curie (Brown et al., 2020). AAR (Yu et al., 2023d) proposes to leverage a small source LM to provide LM-preferred signals for the retriever’s training. The retriever after training can be directly utilized to assist a large target LM by plugging in the retrieved documents.

Unlike prior studies focusing on adapting either the retriever or the language models, another research line focuses on bridging the semantic gap between the input text and the knowledge that is really needed to query. Query2doc (Wang et al., 2023b) prompts the LLMs to generate a pseudo-document by employing a few-shot prompting paradigm. Subsequently, the original query is expanded by incorporating the pseudo-document. The retriever module uses this new query to retrieve a list of relevant documents. Ma et al. (2023) introduce a Rewrite-Retrieve-Read framework for retrieval augmentation, which can be further tuned for adapting to LLMs. They also add a query rewriting step before the retriever. Different with Query2doc,

they adopt a trainable language model to perform the rewriting step. The rewriting language model is trained by reinforcement learning to using the LLM performance as a reward.

**LLM-based.** Generative retrieval is a new paradigm of retrieval method that mainly includes two schemes: generating identifier strings of documents and generating completed documents. The former uses identifiers to reduce the amount of useless information and make it easier for the model to memorize and learn (Li et al., 2023d). De Cao et al. (2020) propose GENRE, which retrieves an entity by generating the entity text itself. GENRE also could be applied in page-level retrieval, where each document contains a unique title as the identifier. Lee et al. (2022) introduce generative retrieval to the multi-hop setting, and the retrieved items are short sentences. Tay et al. (2022) propose the DSI method, which takes numeric IDs as identifiers for documents. Wang et al. (2022) improve the DSI by generating more queries as extra training data. However, the numeric IDs-based methods usually are evaluated on small datasets, partially because they suffer from the large scaling problem. Bevilacqua et al. (2022) propose SEAL, which takes substrings as identifiers. The retrieval process is effectively completed upon the FM-Index structure. Li et al. (2023d) propose multiview identifiers that represented a passage from different perspectives to enhance generative retrieval and achieve state-of-the-art performance.

Instead of generating identifiers, the latter aims to use large language models to directly generate complete documents. Generate-then-read (Yu et al., 2023b) shows that generated contextual documents contain the correct answer more often than the top retrieved documents and significantly outperform directly generating answers from large language models despite not incorporating any new external information. RECITE (Sun et al., 2023) takes a similar approach, which tackles knowledge-intensive NLP tasks by first reciting relevant information and then generating the outputs. PKG (Luo et al., 2023) equips LLMs with a background knowledge generation module to access relevant knowledge. The parametric knowledge module is based on open-source small language models and can be fine-tuned efficiently offline to store any knowledge. Feng et al. (2023a) propose to empower general large language models with modular and collaboratively sourced knowledge through the integration

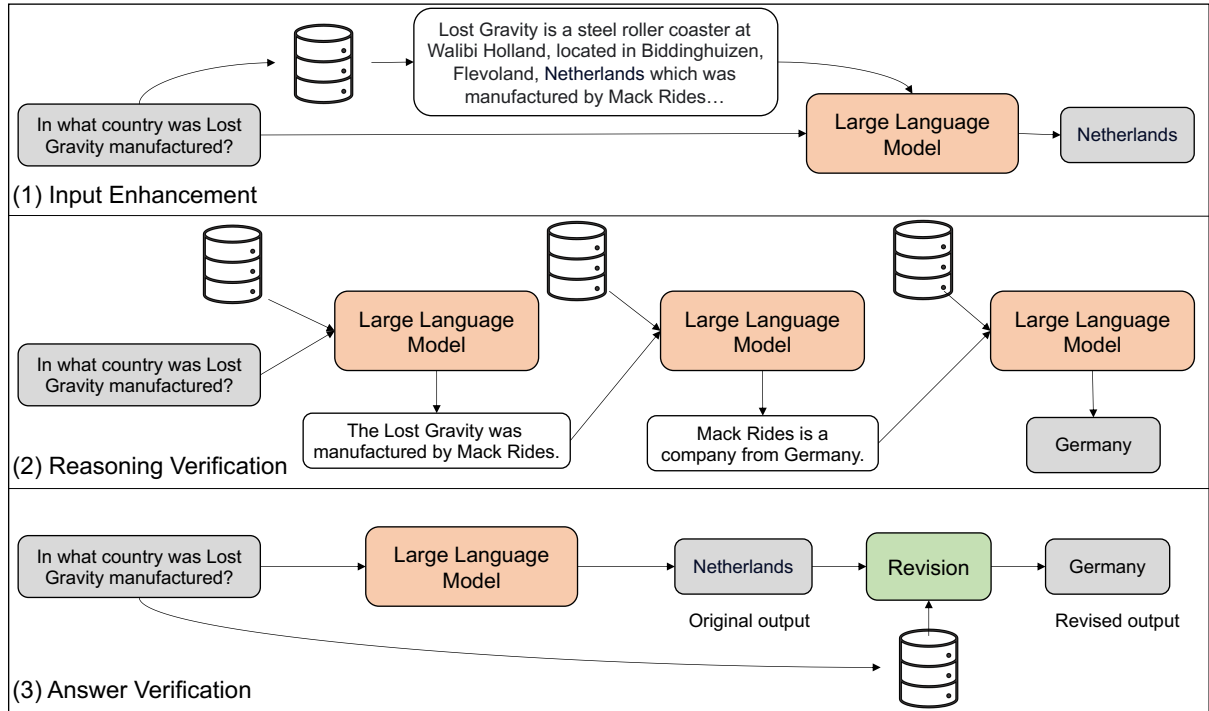


Figure 4: Three kinds of document utilization: (1) input enhancement uses relevant documents as part of the input prompt; (2) reasoning verification uses relevant documents to ensure that the reasoning process is correct; (3) answer verification modifies the original answers of large language models using relevant documents.

of specialized language models. The specialized language models are trained on corpora from diverse sources and domains. They also propose three levels of knowledge filters to dynamically select and refine generated documents and control for topic relevance, document brevity, and knowledge factuality.

*Highlight:* Retriever-based retrieval uses a retriever to fetch relevant documents from an external corpus. However, the retrieved documents might contain noisy information that is irrelevant to the question. Another option is to use large language models to directly generate relevant documents. However, the latter method cannot obtain real-time information. So we should make an appropriate choice according to the actual scenario.

### 3.3 Document Utilization

Once we have the relevant documents, how can we use them to improve the capability of the large language models? As shown in Figure 4, we divide the different ways of using documents into three categories: input enhancement, reasoning verification, and answer verification.

**Input Enhancement.** Language language models define probability distributions over sequences

of tokens. The retrieved top- $k$  documents provide rich information about the original input context and can potentially help the LLMs to make a better prediction. A common approach is to prepend the retrieved documents to the input and feed them into large language models to make the final prediction (Khattab et al., 2023; Yu et al., 2023b; Luo et al., 2023; Feng et al., 2023b).

However, this simple scheme is fundamentally restricted by the number of documents we can include, given the language model’s context window size. REPLUG (Shi et al., 2023) introduces a new ensemble scheme that encodes the retrieved documents in parallel with the same black-box LM. They prepend each document to input separately and then ensemble output probabilities from all  $k$  passes.

**Reasoning Verification.** Large language models are capable of answering complex questions by generating step-by-step natural language reasoning steps — called chain of thought (CoT) (Wei et al., 2022). However, for many open-domain questions, all required knowledge is not always available or up-to-date in models’ parameters and it’s beneficial to retrieve knowledge from external sources (Lazaridou et al., 2022). IRCOT (Trivedi

et al., 2023) proposes an interleaving approach to use retrieval to guide the chain-of-thought reasoning steps and use CoT reasoning to guide the retrieval. Self-ask (Press et al., 2023) builds on chain of thought prompting instead of outputting a continuous undemarcated chain-of-thought. Self-ask clearly demarcates the beginning and end of every sub-question and uses a search engine to answer the sub-questions instead of the LLMs. ReAct (Yao et al., 2023a) prompts LLMs to generate both verbal reasoning traces and actions in an interleaved manner, which allows the model to perform dynamic reasoning to create, maintain, and adjust high-level plans for acting, while also interacting with the external environments to incorporate additional information into reasoning. Verify-and-Edit (Zhao et al., 2023a) seeks to increase prediction factuality by post-editing reasoning chains according to external knowledge. In addition, several recent works (Shao et al., 2023; Feng et al., 2023c; Yu et al., 2023c) first generate initial outputs, then utilize a retrieval model to acquire relevant information from large document collections, and finally incorporate the retrieved information into the in-context demonstration for output refinement.

**Answer Verification.** He et al. (2022) present a post-processing approach called rethinking with retrieval (RR) for utilizing external knowledge in LLMs. They begin by using the chain-of-thought (CoT) prompting method (Wei et al., 2022) to generate a diverse set of reasoning paths. They then use each reasoning step in those paths to retrieve relevant external knowledge, which enables RR to provide more faithful explanations and more accurate predictions. Instead of constraining LMs to generate attributed text, RARR (Gao et al., 2023) proposes a model-agnostic approach to improve the attribution of any existing LM. After generating text, RARR fetches relevant evidence, and then revises the text to make it consistent with the evidence while preserving qualities like style or structure, enabling the revised text to be seamlessly used in place of the original. RARR can be viewed as a retrieval augmented model where retrieval happens after generation rather than before. Peng et al. (2023) present LLM-AUGMENTER to improve LLMs with external knowledge and automated feedback. Given a user query, LLM-AUGMENTER first retrieves evidence from external knowledge and further consolidates evidence by linking retrieved raw evidence with related context

and performing reasoning to form evidence chains. LLM-AUGMENTER then verifies the candidate’s response by checking whether it hallucinates evidence.

*Highlight:* Input enhancement, reasoning verification, and answer verification are three common ways to use relevant documents. They use the relevant documents in three different stages. Input enhancement uses relevant documents as part of the input prompt. Reasoning verification uses relevant documents to ensure that the reasoning process is correct. Answer verification modifies the original answers of large language models using relevant documents. It is worth exploring the effects of using related documents simultaneously at different stages.

### 3.4 Knowledge Conflict

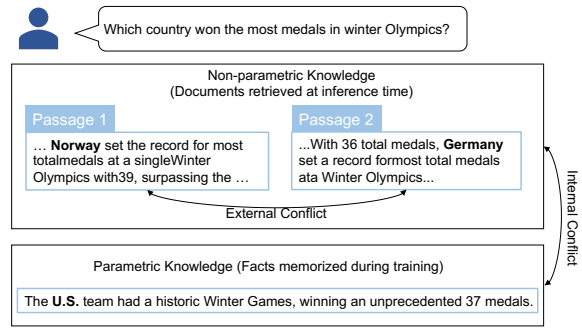


Figure 5: Two kinds of knowledge conflict: internal conflict: inconsistency between the knowledge in large language models and the knowledge in the retrieved documents, and external conflict: inconsistency between the retrieved multiple documents.

In retrieval-augmented LLMs, there are two sources of knowledge contributing to model inference with an ambiguous and opaque division of labor. The first is the implicit parametric knowledge (i.e., their learned weights) instilled by pre-training and fine-tuning. The second is contextual knowledge, usually sourced as passages of text from the retriever. Knowledge conflict means that the information contained is inconsistent and contradictory. As shown in Figure 5, there are two types of knowledge conflicts: internal conflict and external conflict. Internal conflict refers to the inconsistency between the knowledge in large language models and the knowledge in the retrieved documents. External conflict refers to the inconsistency between the retrieved multiple documents.

**Internal Conflict.** As the world is constantly evolving, memorized facts may become outdated (Liška et al., 2022; Kasai et al., 2022). Augment LLM prompting with external context containing relevant knowledge is a promising direction. However, such methods face the challenge that LLMs may persist with the memorized facts and ignore the provided context (Longpre et al., 2021). To tackle this challenge, recent works (Neeman et al., 2022; Li et al., 2022) finetune LLMs on counterfactual contexts, where the original facts are replaced with counterfactual ones. They find that such finetuning processes can effectively improve the LLMs’ utilization of contexts instead of relying solely on their parametric knowledge. Zhou et al. (2023) propose an approach using prompting to improve context faithfulness in LLMs without additional finetuning, which offers a more general and cost-effective method for LLMs. They present various prompting strategies to improve the faithfulness of LLMs, including designing effective prompts and choosing appropriate in-context demonstrations. Xie et al. (2023) present comprehensive and controlled investigation into the behavior of LLMs when encountering counter-memory. They find that with both supportive and contradictory evidence to their parametric memory, LLMs show a strong confirmation bias and tend to cling to their parametric memory.

**External Conflict.** Such scenarios are common in settings where some passages are updated with new information, while other passages remain outdated (Zhang and Choi, 2021). Such conflicts can also occur when passages are adversarially edited to contain false information (Du et al., 2022), or when passages are authored by multiple people who have differing opinions about an answer (Chen et al., 2019). Chen et al. (2022) simulate a setting where a subset of evidence passages are perturbed to suggest a different answer to reflect the realistic scenario where retrieval returns a mixed bag of information. They find that when different passages suggest multiple conflicting answers, models prefer the answer that matches their parametric knowledge. In addition to analyzing simple internal and external conflicts, Xie et al. (2023) also experiment on more complicated knowledge conflict scenarios. With both relevant and irrelevant evidence provided, LLMs can filter out the irrelevant ones to a certain extent. However, as the quantity of irrelevant evidence increases, such an

ability diminishes.

*Highlight:* Knowledge conflict is a very important problem. However, current research focuses on the analysis of knowledge conflicts. The next step should be to resolve knowledge conflicts from different aspects, such as data filtering, model finetuning.

### 3.5 Benchmark

Knowledge-sensitive tasks are ideal for evaluating retrieval-enhanced large language models because solving knowledge-sensitive tasks requires access to a large amount of information. We investigate the commonly used datasets in detail and divide them into the following categories by task type.

**Single-hop QA.** Single-hop questions have relatively simple structures and can be answered using information contained in the paragraph. There are several commonly used datasets, including Natural Questions (NQ), TriviaQA and PopQA. Natural Questions (Kwiatkowski et al., 2019) consists of questions aggregated from the Google search engine, and the answers are annotated by human experts. TriviaQA (Joshi et al., 2017) consists of questions authored by trivia enthusiasts, and evidence documents are collected retrospectively from Wikipedia and the Web. PopQA (Mallen et al., 2023) is a large-scale entity-centric QA dataset, which is constructed to sample more heavily from the tail and has significantly more low-popularity entities.

**Multi-hop QA.** Multi-hop question answering is a challenging subfield of QA that involves answering questions that cannot be resolved with a direct answer from a single source or passage. Models need to perform multiple steps of reasoning in order to answer a question. There are several commonly used datasets, including HotPotQA, 2WikiMultiHopQA, MuSiQue and Bamboogle. HotPotQA (Yang et al., 2018) is collected by explicitly composing questions requiring reasoning about multiple supporting context documents. 2WikiMultihopQA (Ho et al., 2020) is also constructed via composition, but they use a limited set of hand-authored compositional rules. MuSiQue (Trivedi et al., 2022) is constructed with a bottom-up process by carefully selecting and composing single-hop questions. With six composition structures, MuSiQue is more challenging and less cheatable than HotPotQA and 2WikiMultihopQA. Bamboogle (Press et al., 2023) is a small dataset



Task Type	DataSet	Metric
Single-hop QA	Natural Questions (Kwiatkowski et al., 2019)	EM / F <sub>1</sub>
	TriviaQA (Joshi et al., 2017)	EM / F <sub>1</sub>
	PopQA (Mallen et al., 2023)	Accuracy
Multi-hop QA	2WikiMultiHopQA (Ho et al., 2020)	EM / F <sub>1</sub>
	HotPotQA (Yang et al., 2018)	EM / F <sub>1</sub>
	MuSiQue (Trivedi et al., 2022)	EM / F <sub>1</sub>
	Bamboogle (Press et al., 2023)	EM / F <sub>1</sub>
Fact Verification	Fever (Thorne et al., 2018)	Accuracy
	Feverous (Aly et al., 2021)	Accuracy
	FoolMeTwice (Eisenschlos et al., 2021)	Accuracy
Complex Reasoning	StrategyQA (Geva et al., 2021a)	Accuracy
	CommonsenseQA (Talmor et al., 2019)	Accuracy
	CommonsenseQA2.0 (Talmor et al., 2022)	Accuracy
	CSQA (Saha et al., 2018)	EM / ROUGE
	TempQuestions (Jia et al., 2018)	EM / F <sub>1</sub>
	INFOTABS (Gupta et al., 2020)	EM / F <sub>1</sub>

Table 2: Benchmarks for evaluating retrieval-augmented large language models.

with 2-hop questions written by the authors, where all questions are sufficiently difficult to be unanswerable by a popular internet search engine, but where both supporting pieces of evidence can be found in Wikipedia.

**Fact Verification.** Fact verification, also called fact checking, is a challenging task that requires retrieving relevant evidence from plain text and use the evidence to verify given claims. There are several commonly used datasets, including Fever, Feverous and FoolMeTwice (FM2). Fever (Thorne et al., 2018) is a large dataset for fact verification that requires retrieving sentence-level evidence to support if a claim is supported or refuted. In addition to unstructured text evidence, Feverous (Aly et al., 2021) also considers Wikipedia tables as a form of evidence. The evidence retrieval in Feverous considers the entirety of a Wikipedia article and thus the evidence can be located in any section of the article except the reference sections. Feverous is balanced, having almost an equal amount of instances containing, either exclusively text, tables, or both as evidence. FoolMeTwice (Eisenschlos et al., 2021) is collected through a fun multiplayer game, which encourages adversarial examples, drastically lowering the number of examples that can be solved using shortcuts compared to other datasets.

**Complex Reasoning.** Complex Reasoning includes different types of reasoning, such as commonsense reasoning, tabular reasoning, etc. Commonsense reasoning is the foundation of human understanding, rooted in the basic knowledge and life experiences accumulated through daily life and social practice, which outlines practical knowledge of how the world works (Sap et al., 2020). Commonsense reasoning tasks evaluate models’ reasoning skills in the physical world. StrategyQA, CommonsenseQA and CommonsenseQA2.0 are widely used commonsense reasoning datasets. StrategyQA (Geva et al., 2021a) is a question-answering benchmark focusing on open domain questions where the required reasoning steps are implicit in the question and should be inferred using a strategy. CommonsenseQA (Talmor et al., 2019) and CommonsenseQA2.0 (Talmor et al., 2022) is proposed to explore the commonsense understanding ability of large language models, which includes yes/no questions (or assertions) about everyday commonsense knowledge. CSQA (Saha et al., 2018) is a long-form QA, which aims to generate comprehensive answers to questions seeking complex information. TempQuestions (Jia et al., 2018) is built to investigate temporal reasoning. This dataset includes 1,271 temporal questions that are divided into four classes: explicit temporal, implicit temporal, temporal answer, and ordinal

constraints. INFOTABS (Gupta et al., 2020) consists of 23, 738 human-written textual hypotheses based on premises in the form of tables extracted from Wikipedia info-boxes.

## 4 Frontier Applications

### 4.1 Knowledge editing

Knowledge editing offers a cost-effective way of refreshing the outdated information in LLMs. Consequently, its primary purpose is to keep LLMs aligned with the continually evolving world. Towards the errors reported by users following deployment, sequential model editing (SME) also can stand out as an efficient method for rectifying a series of mistakes as a patching mechanism (Huang et al., 2023b; Hartvigsen et al., 2022). In addition to general applications, knowledge editing methods open up new avenues by focusing on information beyond factual knowledge. LLMs might output toxic text or leak personal information when subjected to adversarial prompts (Carlini et al., 2023; Qiu et al., 2023). To mitigate this concern, Patil et al. (2023) introduces an attack-and-defense framework based on knowledge editing methods to remove sensitive information from the model. In addition, compared to traditional controllable text generation techniques (Qian et al., 2022; Gu et al., 2022), knowledge editing can serve as a control methods for the generation of LLMs. From the social psychology, Mao et al. (2023) innovatively employs knowledge editing methods to modify the personality for LLMs which can control the LLM’s open view on specified topics and establish a benchmark dataset called PersonalityEdit. In summary, model editing will continue to play a significant role in the domains of model security and stance control.

### 4.2 Retrieval Augmentation

Augmenting language models with relevant information retrieved from various knowledge stores has been shown to be effective in improving performance on various knowledge-intensive tasks. In open-domain question answering and fact verification, the model can answer the question more accurately by searching relevant documents in a large corpus or on the web. In addition to classical natural language processing tasks, many new applications have emerged with the development of retrieval-augmented large language mod-

els. LangChain <sup>4</sup> is a powerful framework that provides a set of tools, components, and interfaces to simplify the process of creating applications powered by large language models and chat models. LangChain makes it easy to manage interactions with large language models, link multiple components together, and integrate additional resources. ChatPDF <sup>5</sup> is an AI tool that helps you understand and chat with PDF documents. It can identify key information, provide concise summaries, and answer your questions. ChatDoctor (Li et al., 2023e) is an advanced language model that is specifically designed for medical applications. Patients can interact with the ChatDoctor model through a chat interface, asking questions about their health, symptoms, or medical conditions. The model will then analyze the input and provide a response that is tailored to the patient’s unique situation. New Bing uses retrieval augmentation by combining ChatGPT with Microsoft’s search engine. New Bing Chat generates a search query from your prompt, retrieves relevant documents, and uses them as context for its results. New Bing also provides links to sources of information for the sentences it generates. In summary, a large language model has more powerful knowledge understanding, and reasoning ability by retrieving relevant documents, and will have more application scenarios. Baidu, iFlytek and Kunlun also offer similar services, such as ERNIE Bot <sup>6</sup>, Spark <sup>7</sup> and Skywork <sup>8</sup>.

## 5 Future Directions

The development of knowledge editing and retrieval augmentation are still in a rudimentary stage and thus leaves much room for improvement. In this section, we offer a succinct overview of future research.

**Multi-source Knowledge Augmentation.** Existing knowledge enhancement methods mainly exhibit limitations in terms of the formats and varieties of incorporated knowledge. Most knowledge editing methods primarily center around factual knowledge represented as triples, thereby constraining the extent of modifiable knowledge (Meng et al., 2022, 2023). Current retrieval augmentation methods mainly focus on unstructured text

<sup>4</sup><https://www.langchain.com/>

<sup>5</sup><https://www.chatpdf.com/>

<sup>6</sup><https://yiban.baidu.com/>

<sup>7</sup><https://xinghuo.xfyun.cn/>

<sup>8</sup><https://search.tiangong.cn>

retrieval from Wikipedia or web (Shi et al., 2023; Feng et al., 2023b; Vu et al., 2023). In real-world scenarios, a complex question may require fragmented evidence gathered from different sources to have the final answer. It is worth exploring the impact of different sources and different formats of evidence for large language models. In addition, it is extremely important to find a suitable way to integrate evidence from different sources.

**Knowledge Augmented Multi-modal Large Language Models.** Multi-modal learning has attracted increasing research attention due to its huge application potential, as a fundamental technique for vision-to-language reasoning (Li et al., 2023b; Yu et al., 2023a; Chen et al., 2023; Yu et al., 2021b). How to endow the large language model with the ability of multi-modal reasoning is becoming a hot research topic. Yu et al. (2021a) propose to use extra multi-modal knowledge to augment language generation processing for reasoning. Cheng et al. (2023) explore the application of current knowledge editing methods for refining multi-modal models and reveal that the effects remain further improved. RA-CM3 (Yasunaga et al., 2023) proposes a retrieval-augmented multi-modal model, which enables a base multi-modal model to refer to relevant text and images fetched by a retriever from external memory. Future research can further investigate the integration of knowledge and multi-modal large language models to address complex challenges in the real world.

**Large Language Model Based Agents.** Autonomous agents have long been a prominent research focus in both academic and industry communities (Padgham and Winikoff, 2005). Through the acquisition of vast amounts of web knowledge, LLMs have demonstrated remarkable potential in achieving human-level intelligence and brought a glimmer of hope for the further development of agents (OpenAI, 2023; Sumers et al., 2023). These LLM-based agents can exhibit reasoning and planning abilities and have been applied to various real-world scenarios (Qian et al., 2023; Li et al., 2023a). Due to the diversity of the real world, LLM-based agents need additional information to make decisions. It is very important for the development of LLM-based Agents to explore the integration of knowledge and large language model methods in actual and complex scenarios.

## Analysis of Knowledge Enhancement Methods.

Current knowledge enhancement methods mainly focus on the generation results of the model, and other aspects still need to be studied. Li et al. (2023f) introduce innovative evaluation metrics to analyze the side effects of knowledge editing. In addition, Pinter and Elhadad (2023) scrutinize the limitations of extensive models and the ramifications of knowledge editing, and raise concerns about the current methods, contending that they pose a potential risk to users of LLMs when compared to alternative interactive methods (e.g. retrieval-based architectures). Future research can establish exhaustive and empirically-driven analysis, which will enhance the comprehension of the viability of current methods and guide the application of knowledge enhancement in actual scenarios.

## 6 Conclusion

In this paper, we perform a survey on integration of knowledge and large language models and offer a broad view of its main directions, including knowledge editing and retrieval augmentation. Moreover, we summarize the commonly used benchmarks and frontier applications and point out some promising research directions. We hope this survey can offer readers a clear picture of the current progress and inspire more work.

## References

- Rami Aly, Zhijiang Guo, Michael Schlichtkrull, James Thorne, Andreas Vlachos, Christos Christodoulopoulos, Oana Cocarascu, and Arpit Mittal. 2021. Feverous: Fact extraction and verification over unstructured and structured information. *arXiv preprint arXiv:2106.05707*.
- Berkeley R Andrus, Yeganeh Nasiri, Shilong Cui, Benjamin Cullen, and Nancy Fulda. 2022. Enhanced story comprehension for large language models through dynamic document-based knowledge graphs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10436–10444.
- Jinheon Baek, Alham Fikri Aji, and Amir Saffari. 2023a. Knowledge-augmented language model prompting for zero-shot knowledge graph question answering. *arXiv preprint arXiv:2306.04136*.
- Jinheon Baek, Soyeong Jeong, Minki Kang, Jong C. Park, and Sung Ju Hwang. 2023b. [Knowledge-augmented language model verification](#).
- Michele Bevilacqua, Giuseppe Ottaviano, Patrick Lewis, Scott Yih, Sebastian Riedel, and Fabio Petroni. 2022. Autoregressive search engines: Generating substrings

- as document identifiers. *Advances in Neural Information Processing Systems*, 35:31668–31683.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. 2022. [Improving language models by retrieving from trillions of tokens](#).
- Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*.
- Davis Brown, Charles Godfrey, Cody Nizinski, Jonathan Tu, and Henry Kvinge. 2023. [Edit at your own risk: evaluating the robustness of edited models to distribution shifts](#).
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#).
- Nicholas Carlini, Milad Nasr, Christopher A. Choquette-Choo, Matthew Jagielski, Irena Gao, Anas Awadalla, Pang Wei Koh, Daphne Ippolito, Katherine Lee, Florian Tramèr, and Ludwig Schmidt. 2023. [Are aligned neural networks adversarially aligned?](#)
- Hung-Ting Chen, Michael JQ Zhang, and Eunsol Choi. 2022. Rich knowledge sources bring complex knowledge conflicts: Recalibrating models to reflect conflicting evidence. *arXiv preprint arXiv:2210.13701*.
- Sihao Chen, Daniel Khashabi, Wenpeng Yin, Chris Callison-Burch, and Dan Roth. 2019. [Seeing things from a different angle: discovering diverse perspectives about claims](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 542–557, Minneapolis, Minnesota. Association for Computational Linguistics.
- Yangyi Chen, Karan Sikka, Michael Cogswell, Heng Ji, and Ajay Divakaran. 2023. Measuring and improving chain-of-thought reasoning in vision-language models. *arXiv preprint arXiv:2309.04461*.
- Siyuan Cheng, Bozhong Tian, Qingbin Liu, Xi Chen, Yongheng Wang, Huajun Chen, and Ningyu Zhang. 2023. [Can we edit multimodal large language models?](#)
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.
- Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. 2023. [Evaluating the ripple effects of knowledge editing in language models](#).
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In *Machine learning challenges workshop*, pages 177–190. Springer.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022a. Knowledge neurons in pretrained transformers. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8493–8502.
- Damai Dai, Wenbin Jiang, Qingxiu Dong, Yajuan Lyu, Qiaoqiao She, and Zhifang Sui. 2022b. Neural knowledge bank for pretrained transformers. *arXiv preprint arXiv:2208.00399*.
- Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. [Editing factual knowledge in language models](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6491–6506, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. 2020. Autoregressive entity retrieval. *arXiv preprint arXiv:2010.00904*.
- Qingxiu Dong, Damai Dai, Yifan Song, Jingjing Xu, Zhifang Sui, and Lei Li. 2022. [Calibrating factual knowledge in pretrained language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5937–5947, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Yihong Dong, Kangcheng Luo, Xue Jiang, Zhi Jin, and Ge Li. 2023. Pace: Improving prompt with actor-critic editing for large language model. *arXiv preprint arXiv:2308.10088*.
- Yibing Du, Antoine Bosselut, and Christopher D Manning. 2022. Synthetic disinformation attacks on automated fact verification systems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10581–10589.
- Julian Eisenschlos, Bhuwan Dhingra, Jannis Bulian, Benjamin Börschinger, and Jordan Boyd-Graber. 2021. [Fool me twice: Entailment from Wikipedia gamification](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association*



- for *Computational Linguistics: Human Language Technologies*, pages 352–365, Online. Association for Computational Linguistics.
- Shangbin Feng, Weijia Shi, Yuyang Bai, Vidhisha Balachandran, Tianxing He, and Yulia Tsvetkov. 2023a. Cook: Empowering general-purpose language models with modular and collaborative knowledge. *arXiv preprint arXiv:2305.09955*.
- Shangbin Feng, Weijia Shi, Yuyang Bai, Vidhisha Balachandran, Tianxing He, and Yulia Tsvetkov. 2023b. Knowledge card: Filling llms’ knowledge gaps with plug-in specialized language models.
- Zhangyin Feng, Xiaocheng Feng, Dezhi Zhao, Maojin Yang, and Bing Qin. 2023c. Retrieval-generation synergy augmented large language models. *arXiv preprint arXiv:2310.05149*.
- Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and Kelvin Guu. 2023. RARR: Researching and revising what language models say, using language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16477–16508, Toronto, Canada. Association for Computational Linguistics.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021a. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361.
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021b. Transformer feed-forward layers are key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yuxuan Gu, Xiaocheng Feng, Sicheng Ma, Lingyuan Zhang, Heng Gong, and Bing Qin. 2022. A distributional lens for multi-aspect controllable text generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1023–1043, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Anshita Gupta, Debanjan Mondal, Akshay Krishna Sheshadri, Wenlong Zhao, Xiang Lorraine Li, Sarah Wiegrefe, and Niket Tandon. 2023. Editing commonsense knowledge in GPT. *CoRR*, abs/2305.14956.
- Vivek Gupta, Maitrey Mehta, Pegah Nokhiz, and Vivek Srikumar. 2020. INFOTABS: Inference on tables as semi-structured data. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2309–2324, Online. Association for Computational Linguistics.
- Thomas Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. 2022. Aging with grace: Lifelong model editing with discrete key-value adaptors. *arXiv preprint arXiv:2211.11031*.
- Peter Hase, Mohit Bansal, Been Kim, and Asma Ghan-deharioun. 2023. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models. *CoRR*, abs/2301.04213.
- Hangfeng He, Hongming Zhang, and Dan Roth. 2022. Rethinking with retrieval: Faithful large language model inference.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Jason Hoelscher-Obermaier, Julia Persson, Esben Kran, Ioannis Konstas, and Fazl Barez. 2023. Detecting edit failures in large language models: An improved specificity benchmark.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 2022. Training compute-optimal large language models.
- Linmei Hu, Zeyi Liu, Ziwang Zhao, Lei Hou, Liqiang Nie, and Juanzi Li. 2023. A survey of knowledge enhanced pre-trained language models. *IEEE Transactions on Knowledge and Data Engineering*.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023a. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions.
- Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. 2023b. Transformer-patcher: One mistake worth one neuron. *arXiv preprint arXiv:2301.09785*.
- Gabriel Ilharco, Marco Túlio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Gautier Izacard and Edouard Grave. 2022. Distilling knowledge from reader to retriever for question answering.

- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022. [Atlas: Few-shot learning with retrieval augmented language models](#).
- Zhen Jia, Abdalghani Abujabal, Rishiraj Saha Roy, Jan-nik Strötgen, and Gerhard Weikum. 2018. Tempques-tions: A benchmark for temporal question answering. In *Companion Proceedings of the The Web Conference 2018*, pages 1057–1062.
- Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. 2021. [How can we know when language models know? on the calibration of language models for question answering](#). *Transactions of the Association for Computational Linguistics*, 9:962–977.
- Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. [Active retrieval augmented generation](#).
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. [TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Mingxuan Ju, Wenhao Yu, Tong Zhao, Chuxu Zhang, and Yanfang Ye. 2022. [Grape: Knowledge graph enhanced passage reader for open-domain question answering](#).
- Yiming Ju and Zheng Zhang. 2023. [Klob: a benchmark for assessing knowledge locating methods in language models](#).
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*.
- Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In *International Conference on Machine Learning*, pages 15696–15707. PMLR.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick S. H. Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. [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*, pages 6769–6781. Association for Computational Linguistics.
- Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir Radev, Noah A. Smith, Yejin Choi, and Kentaro Inui. 2022. [Realtime qa: What’s the answer right now?](#)
- Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. 2023. [Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive nlp](#).
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natural questions: A benchmark for question answering research](#). *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Angeliki Lazaridou, Elena Gribovskaya, Wojciech Stokowiec, and Nikolai Grigorev. 2022. Internet-augmented language models through few-shot prompting for open-domain question answering. *arXiv preprint arXiv:2203.05115*.
- Hyunji Lee, Sohee Yang, Hanseok Oh, and Minjoon Seo. 2022. Generative multi-hop retrieval. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1417–1436.
- Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. [Zero-shot relation extraction via reading comprehension](#). In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, pages 333–342, Vancouver, Canada. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#).
- Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin Wang, Michal Lukasik, Andreas Veit, Felix Yu, and Sanjiv Kumar. 2022. [Large language models with controllable working memory](#).
- Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023a. [Camel: Communicative agents for "mind" exploration of large language model society](#).
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. [Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models](#). *arXiv preprint arXiv:2301.12597*.
- Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun Ma, and Jie Yu. 2023c. [PMET: precise model editing in a transformer](#). *CoRR*, abs/2308.08742.

- Yongqi Li, Nan Yang, Liang Wang, Furu Wei, and Wenjie Li. 2023d. [Multiview identifiers enhanced generative retrieval](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6636–6648, Toronto, Canada. Association for Computational Linguistics.
- Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. 2023e. Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge. *Cureus*, 15(6).
- Zhoubo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang, Xi Chen, and Huajun Chen. 2023f. [Unveiling the pitfalls of knowledge editing for large language models](#).
- Adam Liška, Tomáš Kočiský, Elena Gribovskaya, Tayfun Terzi, Eren Sezener, Devang Agrawal, Cyprien de Masson d’Autume, Tim Scholtes, Manzil Zaheer, Susannah Young, Ellen Gilsonan-McMahon, Sophia Austin, Phil Blunsom, and Angeliki Lazaridou. 2022. [Streamingqa: A benchmark for adaptation to new knowledge over time in question answering models](#).
- Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021. [Entity-based knowledge conflicts in question answering](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7052–7063, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ziyang Luo, Can Xu, Pu Zhao, Xiubo Geng, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. [Augmented large language models with parametric knowledge guiding](#). *arXiv preprint arXiv:2305.04757*.
- Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. 2023. [Query rewriting for retrieval-augmented large language models](#). *arXiv preprint arXiv:2305.14283*.
- Aman Madaan, Niket Tandon, Peter Clark, and Yiming Yang. 2022. [Memory-assisted prompt editing to improve GPT-3 after deployment](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2833–2861, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. [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)*, pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.
- Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. 2023. [Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models](#).
- Shengyu Mao, Ningyu Zhang, Xiaohan Wang, Mengru Wang, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Huajun Chen. 2023. [Editing personality for llms](#).
- Antonis Maronikolakis and Hinrich Schütze. 2021. [Multidomain pretrained language models for green NLP](#). In *Proceedings of the Second Workshop on Domain Adaptation for NLP*, pages 1–8, Kyiv, Ukraine. Association for Computational Linguistics.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. [Locating and editing factual associations in GPT](#). In *NeurIPS*.
- Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. 2023. [Mass-editing memory in a transformer](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022a. [Fast model editing at scale](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. 2022b. [Memory-based model editing at scale](#). In *International Conference on Machine Learning*, pages 15817–15831. PMLR.
- Fedor Moiseev, Zhe Dong, Enrique Alfonseca, and Martin Jaggi. 2022. [SKILL: Structured knowledge infusion for large language models](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1581–1588, Seattle, United States. Association for Computational Linguistics.
- Ella Neeman, Roei Aharoni, Or Honovich, Leshem Choshen, Idan Szpektor, and Omri Abend. 2022. [Disentqa: Disentangling parametric and contextual knowledge with counterfactual question answering](#).
- OpenAI. 2023. [Gpt-4 technical report](#).
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#).
- Lin Padgham and Michael Winikoff. 2005. *Developing intelligent agent systems: A practical guide*. John Wiley & Sons.
- Vaidehi Patil, Peter Hase, and Mohit Bansal. 2023. [Can sensitive information be deleted from llms? objectives for defending against extraction attacks](#).



- Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, and Jianfeng Gao. 2023. [Check your facts and try again: Improving large language models with external knowledge and automated feedback](#).
- Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim Rocktäschel, and Sebastian Riedel. 2021. [KILT: a benchmark for knowledge intensive language tasks](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2523–2544, Online. Association for Computational Linguistics.
- Yuval Pinter and Michael Elhadad. 2023. [Emptying the ocean with a spoon: Should we edit models?](#)
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. 2023. [Measuring and narrowing the compositionality gap in language models](#).
- Chen Qian, Xin Cong, Wei Liu, Cheng Yang, Weize Chen, Yusheng Su, Yufan Dang, Jiahao Li, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023. [Communicative agents for software development](#).
- Jing Qian, Li Dong, Yelong Shen, Furu Wei, and Weizhu Chen. 2022. [Controllable natural language generation with contrastive prefixes](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2912–2924, Dublin, Ireland. Association for Computational Linguistics.
- Huachuan Qiu, Shuai Zhang, Anqi Li, Hongliang He, and Zhenzhong Lan. 2023. [Latent jailbreak: A benchmark for evaluating text safety and output robustness of large language models](#).
- Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. [In-context retrieval-augmented language models](#).
- Vipula Rawte, Amit Sheth, and Amitava Das. 2023. [A survey of hallucination in large foundation models](#).
- Ruiyang Ren, Yuhao Wang, Yingqi Qu, Wayne Xin Zhao, Jing Liu, Hao Tian, Hua Wu, Ji-Rong Wen, and Haifeng Wang. 2023. Investigating the factual knowledge boundary of large language models with retrieval augmentation. *arXiv preprint arXiv:2307.11019*.
- Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.
- Amrita Saha, Vardaan Pahuja, Mitesh Khapra, Karthik Sankaranarayanan, and Sarath Chandar. 2018. Complex sequential question answering: Towards learning to converse over linked question answer pairs with a knowledge graph. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Maarten Sap, Vered Shwartz, Antoine Bosselut, Yejin Choi, and Dan Roth. 2020. [Commonsense reasoning for natural language processing](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*, pages 27–33, Online. Association for Computational Linguistics.
- Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. [Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy](#).
- Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen tau Yih. 2023. [Replug: Retrieval-augmented black-box language models](#).
- Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Boyd-Graber, and Li-juan Wang. 2022. Prompting gpt-3 to be reliable. *arXiv preprint arXiv:2210.09150*.
- Anton Sinitin, Vsevolod Plokhotnyuk, Dmitriy Pyrkun, Sergei Popov, and Artem Babenko. 2020. Editable neural networks. *arXiv preprint arXiv:2004.00345*.
- Theodore R. Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L. Griffiths. 2023. [Cognitive architectures for language agents](#).
- Zhiqing Sun, Xuezhi Wang, Yi Tay, Yiming Yang, and Denny Zhou. 2023. [Recitation-augmented language models](#).
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [CommonsenseQA: A question answering challenge targeting commonsense knowledge](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Alon Talmor, Ori Yoran, Ronan Le Bras, Chandra Bhagavatula, Yoav Goldberg, Yejin Choi, and Jonathan Berant. 2022. Commonsenseqa 2.0: Exposing the limits of ai through gamification. *arXiv preprint arXiv:2201.05320*.
- Yi Tay, Vinh Tran, Mostafa Dehghani, Jianmo Ni, Dara Bahri, Harsh Mehta, Zhen Qin, Kai Hui, Zhe Zhao, Jai Gupta, et al. 2022. Transformer memory as a differentiable search index. *Advances in Neural Information Processing Systems*, 35:21831–21843.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. [FEVER: a large-scale dataset for fact extraction and VERification](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics*:



- Human Language Technologies, Volume 1 (Long Papers)*, pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#).
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. [Musique: Multihop questions via single-hop question composition](#).
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. [Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions](#).
- Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and Thang Luong. 2023. [Freshllms: Refreshing large language models with search engine augmentation](#).
- Jiaan Wang, Yunlong Liang, Zengkui Sun, Yuxuan Cao, and Jiarong Xu. 2023a. [Cross-lingual knowledge editing in large language models](#).
- Liang Wang, Nan Yang, and Furu Wei. 2023b. [Query2doc: Query expansion with large language models](#).
- Yujing Wang, Yingyan Hou, Haonan Wang, Ziming Miao, Shibin Wu, Qi Chen, Yuqing Xia, Chengmin Chi, Guoshuai Zhao, Zheng Liu, et al. 2022. A neural corpus indexer for document retrieval. *Advances in Neural Information Processing Systems*, 35:25600–25614.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Suhang Wu, Minlong Peng, Yue Chen, Jinsong Su, and Mingming Sun. 2023. [Eva-kellm: A new benchmark for evaluating knowledge editing of llms](#). *CoRR*, abs/2308.09954.
- Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and Yu Su. 2023. [Adaptive chameleon or stubborn sloth: Unraveling the behavior of large language models in knowledge clashes](#).
- Canwen Xu, Yichong Xu, Shuohang Wang, Yang Liu, Chenguang Zhu, and Julian McAuley. 2023. Small models are valuable plug-ins for large language models. *arXiv preprint arXiv:2305.08848*.
- Linyao Yang, Hongyang Chen, Zhao Li, Xiao Ding, and Xindong Wu. 2023. Chatgpt is not enough: Enhancing large language models with knowledge graphs for fact-aware language modeling. *arXiv preprint arXiv:2306.11489*.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. [HotpotQA: A dataset for diverse, explainable multi-hop question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023a. [React: Synergizing reasoning and acting in language models](#).
- Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023b. [Editing large language models: Problems, methods, and opportunities](#).
- Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Rich James, Jure Leskovec, Percy Liang, Mike Lewis, Luke Zettlemoyer, and Wen tau Yih. 2023. [Retrieval-augmented multimodal language modeling](#).
- Da Yin, Li Dong, Hao Cheng, Xiaodong Liu, Kai-Wei Chang, Furu Wei, and Jianfeng Gao. 2022. [A survey of knowledge-intensive nlp with pre-trained language models](#).
- Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. 2023. [Do large language models know what they don’t know?](#) In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8653–8665, Toronto, Canada. Association for Computational Linguistics.
- Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. 2023. Making retrieval-augmented language models robust to irrelevant context.
- Weijiang Yu, Jian Liang, Lei Ji, Lu Li, Yuejian Fang, Nong Xiao, and Nan Duan. 2021a. Hybrid reasoning network for video-based commonsense captioning. In *Proceedings of the 29th ACM international conference on multimedia*, pages 5213–5221.
- Weijiang Yu, Haofan Wang, Guohao Li, Nong Xiao, and Bernard Ghanem. 2023a. Knowledge-aware global reasoning for situation recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Weijiang Yu, Haoteng Zheng, Mengfei Li, Lei Ji, Lijun Wu, Nong Xiao, and Nan Duan. 2021b. Learning from inside: Self-driven siamese sampling and reasoning for video question answering. *Advances in Neural Information Processing Systems*, 34:26462–26474.
- Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2023b. [Generate rather than retrieve: Large language models are strong context generators](#).
- Wenhao Yu, Zhihan Zhang, Zhenwen Liang, Meng Jiang, and Ashish Sabharwal. 2023c. [Improving language models via plug-and-play retrieval feedback](#).

- Wenhao Yu, Chenguang Zhu, Zaitang Li, Zhiting Hu, Qingyun Wang, Heng Ji, and Meng Jiang. 2022. [A survey of knowledge-enhanced text generation](#). *ACM Computing Surveys*, 54(11s):1–38.
- Zichun Yu, Chenyan Xiong, Shi Yu, and Zhiyuan Liu. 2023d. [Augmentation-adapted retriever improves generalization of language models as generic plug-in](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2421–2436, Toronto, Canada. Association for Computational Linguistics.
- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*.
- Michael Zhang and Eunsol Choi. 2021. [SituatingQA: Incorporating extra-linguistic contexts into QA](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7371–7387, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023a. [Siren’s song in the ai ocean: A survey on hallucination in large language models](#).
- Zihan Zhang, Meng Fang, Ling Chen, Mohammad-Reza Namazi-Rad, and Jun Wang. 2023b. [How do large language models capture the ever-changing world knowledge? a review of recent advances](#).
- Ruochen Zhao, Xingxuan Li, Shafiq Joty, Chengwei Qin, and Lidong Bing. 2023a. [Verify-and-edit: A knowledge-enhanced chain-of-thought framework](#).
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023b. [A survey of large language models](#).
- Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023. Can we edit factual knowledge by in-context learning? *arXiv preprint arXiv:2305.12740*.
- Zexuan Zhong, Zhengxuan Wu, Christopher D Manning, Christopher Potts, and Danqi Chen. 2023. Mquake: Assessing knowledge editing in language models via multi-hop questions. *arXiv preprint arXiv:2305.14795*.
- Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and Muhao Chen. 2023. [Context-faithful prompting for large language models](#).