



# A Survey on Knowledge Conflicts in the Era of LLMs

Anonymous ACL submission

## Abstract

This survey presents a comprehensive examination of knowledge conflicts in Large Language Models (LLMs). It explores the intricate challenges that arise when LLMs integrate contextual knowledge with their parametric knowledge. Our focus is on three primary types of knowledge conflicts: context-memory, inter-context, and intra-memory conflict. These conflicts can significantly impact the trustworthiness and accuracy of LLMs, especially in real-world applications where misinformation and noise are prevalent. The survey categorizes these conflicts, investigates their causes, and reviews potential mitigation strategies. It aims to provide insights into enhancing the robustness of LLMs, making it a valuable resource for advancing research in this evolving area.

## 1 Introduction

Large language models (LLMs) (Brown et al., 2020; OpenAI, 2023b; Touvron et al., 2023) are renowned for encapsulating a vast repository of world knowledge (Petroni et al., 2019; Roberts et al., 2020), often referred to as *parametric knowledge*. These models demonstrate exceptional proficiency in knowledge-intensive tasks including QA (Petroni et al., 2019), fact-checking (Gao et al., 2023a), knowledge generation (Chen et al., 2023c), *inter alia*. Concurrently, LLMs continue to engage with external *contextual knowledge* in their applications (Pan et al., 2022). This external knowledge may originate from various sources, including user prompts (Liu et al., 2023a), interactive dialogues (Zhang et al., 2020), or retrieved documents from the Web (Lewis et al., 2020; Shi et al., 2023c), and tools (Schick et al., 2023; Zhuang et al., 2023).

While integrating contextual information is intended to augment LLMs, enabling them to keep abreast of current events (Kasai et al., 2022) and generate more accurate responses (Shuster et al., 2021), it can also pose challenges. This integration

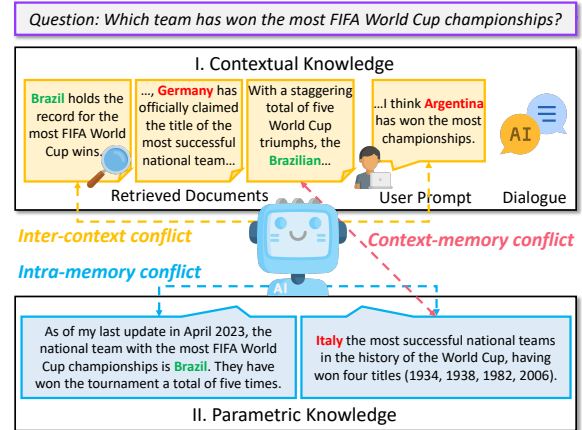


Figure 1: LLMs can encounter three distinct types of knowledge conflicts, which depend on the sources of knowledge — whether provided in the context (the yellow chatboxes) or parametric (the blue chatboxes) within the LLM itself. In the given example, the LLM is presented with a complex scenario involving rich conflicts about the knowledge related to a user’s question (the purple chatbox). This scenario requires the LLM to resolve the conflicts to provide accurate responses.

may lead to interference with the LLM’s parametric knowledge. Furthermore, in real-world scenarios, the context provided to these models might be fraught with noise (Zhang and Choi, 2021) or even deliberately crafted misinformation (Du et al., 2022b; Pan et al., 2023a), complicating their ability to process and respond accurately (Chen et al., 2022). **The discrepancies among the contexts and the model’s parametric knowledge are referred to as knowledge conflicts.**

Knowledge conflict is rooted in open-domain QA research. The concept gained attention in Longpre et al. (2021) that focused on the entity-based conflicts between parametric knowledge and external passages. Concurrently, discrepancies among multiple passages were also scrutinized in the same year (Chen et al., 2022). Knowledge conflicts attract significant attention with the recent advent of LLMs. For instance, recent studies find that LLMs

exhibit both adherence to parametric knowledge and susceptibility to contextual influences (Xie et al., 2023), which can be problematic when this external knowledge is factually incorrect (Pan et al., 2023b). Given the implications for the trustworthiness (Du et al., 2022b), real-time accuracy (Kasai et al., 2022), and robustness of LLMs (Ying et al., 2023), it is imperative to delve deeper into understanding and resolving knowledge conflicts (Xie et al., 2023; Wang et al., 2023g).

As of the time of writing, to the best of our knowledge, there is no systematic survey dedicated to the investigation of knowledge conflicts. Existing reviews (Zhang et al., 2023b; Wang et al., 2023a; Feng et al., 2023) touch upon knowledge conflicts as a subtopic within their broader contexts. To fill the gap, we aim to provide a comprehensive review encompassing the categorization, cause and behavior analysis, and mitigation strategies for addressing various forms of knowledge conflicts.

## 2 The Problem

### 2.1 Background

There are three kinds of knowledge conflicts, as illustrated in Figure 1. Considering a user’s question, the LLM’s provided context (which may include retrieved documents, user prompts, dialogue history, *inter alia*) and its parameterized information can lead to conflicting answers. The three types of conflicts can occur simultaneously, presenting significant challenges for LLMs in deriving factually accurate responses. The typical one is the discrepancies between contextual information and the models’ parameterized knowledge (Longpre et al., 2021; Li et al., 2022a; Xie et al., 2023). This phenomenon, dubbed as **context-memory conflict (CM)**, is detailed in § 3. Following Chen et al. (2022) and Feng et al. (2023), we also recognize the importance of conflicts within the contextual information itself. This is particularly relevant in the era of retrieval-augmented language models (RALMs) (Lewis et al., 2020), and is referred to as **inter-context conflict (IC)**, see § 4. Furthermore, we consider **intra-memory conflict (IM)**, which occurs when an LLM’s internal memorized knowledge is in contradiction (deferred to § A for space limitations).

### 2.2 Problem Formulation

A knowledge conflict can be formally represented using a pair of knowledge statements  $(x, x')$ , where

$x$  and  $x'$  are *distinct* pieces of information about the *same subject*. The conflict arises when these statements are contradictory or incompatible with each other, *i.e.*,

$$\text{KC}(x, x') = \begin{cases} \text{True}, & \text{if } x \wedge x' = \text{False} \\ \text{False}, & \text{otherwise,} \end{cases} \quad (1)$$

where  $x \wedge x'$  denotes two statements are true simultaneously. *Please note that knowledge conflicts can also be generalized to more than 2 statements.*

In the context of large language models, we use  $G \sim \text{LLM}(\cdot|C)$  to denote the generation process, where  $G$  is the generation and  $C$  is the given context. If  $\text{KC}(x, x') = \text{True}$ , then the LLM is encountered with a knowledge conflict. Depending on the origins of  $(x, x')$ , we categorize knowledge conflicts to three types:

- *Context-memory conflict*: Using a probe prompt  $p_{\text{probe}}$  to elicit the memorized (parametric) knowledge  $\text{Mem} \sim \text{LLM}(\cdot|p_{\text{probe}})$ , where  $\text{Mem} \vdash x^1$ . The LLM is provided with a context  $C$ , where  $C \vdash x'$ .
- *Inter-context conflict*: The LLM is provided with a context  $C$ , where  $C \vdash x \wedge C \vdash x'$ .
- *Intra-memory conflict*: Given two probes  $p_{\text{probe}}, p'_{\text{probe}}$ .  $\text{Mem} \sim \text{LLM}(\cdot|p_{\text{probe}})$  and  $\text{Mem}' \sim \text{LLM}(\cdot|p'_{\text{probe}})$ , where  $\text{Mem} \vdash x \wedge \text{Mem}' \vdash x'$ .

Note that the probe prompt query the LLM in a closed-book QA setting (Roberts et al., 2020), *i.e.*, it *does not involve* the direct knowledge about the subject, *i.e.*,  $p_{\text{probe}} \not\vdash x \wedge p_{\text{probe}} \not\vdash x'$ .

### 2.3 Our Philosophy

As illustrated in Figure 2, we conceptualize *lifecycle of knowledge conflicts* as both an *effect*, originating from various causes, and as a *cause* leading to various behaviors in the model. Knowledge conflicts serve as a crucial intermediary between causes and effects. For instance, they are a significant factor that can cause the model to produce factually incorrect information, *a.k.a.*, hallucinations (Ji et al., 2023; Zhang et al., 2023b). Our research, in a manner akin to Freudian psychoanalysis, underscores the significance of understanding the origins of these conflicts. Although some existing analyses (Chen et al., 2022; Xie et al., 2023;

<sup>1</sup> $x \vdash y$  denotes  $y$  is entailed by  $x$ .

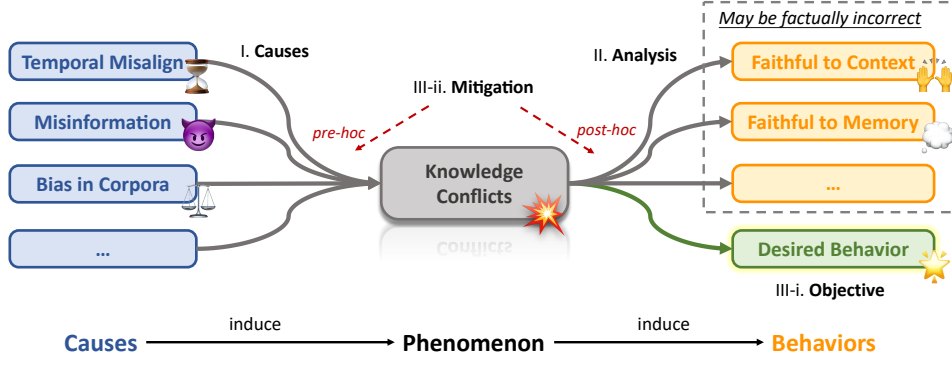


Figure 2: We view knowledge conflict not only as a standalone **phenomenon** but also as a nexus that connects various causal triggers (**causes**) with the **behaviors** of LLMs. While existing (research) literature mainly focus on *II. Analysis*, our survey involves systematically observing these conflicts, offering insights into their emergence and impact on LLM behavior, along with the desirable behaviors and related mitigation strategies.

Wang et al., 2023g) tend to construct such conflicts artificially, we posit that these analyses do not sufficiently address the interconnectedness of the issue.

Going beyond reviewing and analyzing causes and effects, we delve deeper to provide a systematic review of mitigation strategies, which are employed to minimize the undesirable consequences of knowledge conflicts, *i.e.*, to encourage the model to exhibit *desired behaviors that conform to specific objectives* (it should be noted that these *objectives may differ* based on the particular scenario). Based on the timing relative to potential conflicts, strategies are divided into two primary categories: *pre-hoc* and *post-hoc* strategies. The key distinction between pre-hoc and post-hoc approaches lies in whether adjustments are made *before* or *after* potential conflicts arise<sup>2</sup>. The taxonomy of knowledge conflicts is outlined in Figure 3. We sequentially discuss the three kinds of knowledge conflicts, detailing for each the causes, analysis, and available mitigation strategies, which are organized according to their respective objectives.

### 3 Context-Memory Conflict

Context-memory conflict is the most extensively studied one among the three types of conflict. LLMs often have fixed parametric knowledge due to the prohibitive costs of training (Sharir et al., 2020; Hoffmann et al., 2022; Smith, 2023), while external information continues to evolve rapidly (De Cao et al., 2021; Kasai et al., 2022).

<sup>2</sup>Another interpretation is that pre-hoc strategy is proactive while post-hoc is reactive.

#### 3.1 Causes

The ultimate reason for context-memory conflict is the knowledge disparity between the context and parametric knowledge. We consider two main causes: temporal misalignment (Lazaridou et al., 2021; Luu et al., 2021; Dhingra et al., 2022) and misinformation pollution (Du et al., 2022b; Pan et al., 2023a).

**Temporal Misalignment.** Temporal misalignment is *natural* since a model trained on data collected in the past may not accurately reflect current or future states (*i.e.*, the contextual knowledge after the deployment) (Luu et al., 2021; Lazaridou et al., 2021; Liska et al., 2022). Such misalignment can lead to decreased performance and relevancy of the model’s outputs over time, as it may not account for new trends, changes in language use, cultural shifts, or updates in knowledge. Researchers have observed that temporal misalignment relegates the model’s performance on various NLP tasks (Luu et al., 2021; Zhang and Choi, 2021; Dhingra et al., 2022; Kasai et al., 2022; Cheang et al., 2023). Temporal Misalignment only seems to get worse in the future due to the paradigm of pre-training and the increased training costs that accompany model enlargement (Chowdhery et al., 2023; OpenAI, 2023b).

Prior work tries to address temporal misalignment by focusing on three lines: *Knowledge editing (KE)* focuses on updating the parametric knowledge of an existing pre-trained model directly (Sinitsin et al., 2019; De Cao et al., 2021; Mitchell et al., 2021; Onoe et al., 2023). *Retrieval-augmented generation (RAG)* leverages a retrieval model to fetch relevant documents from external sources

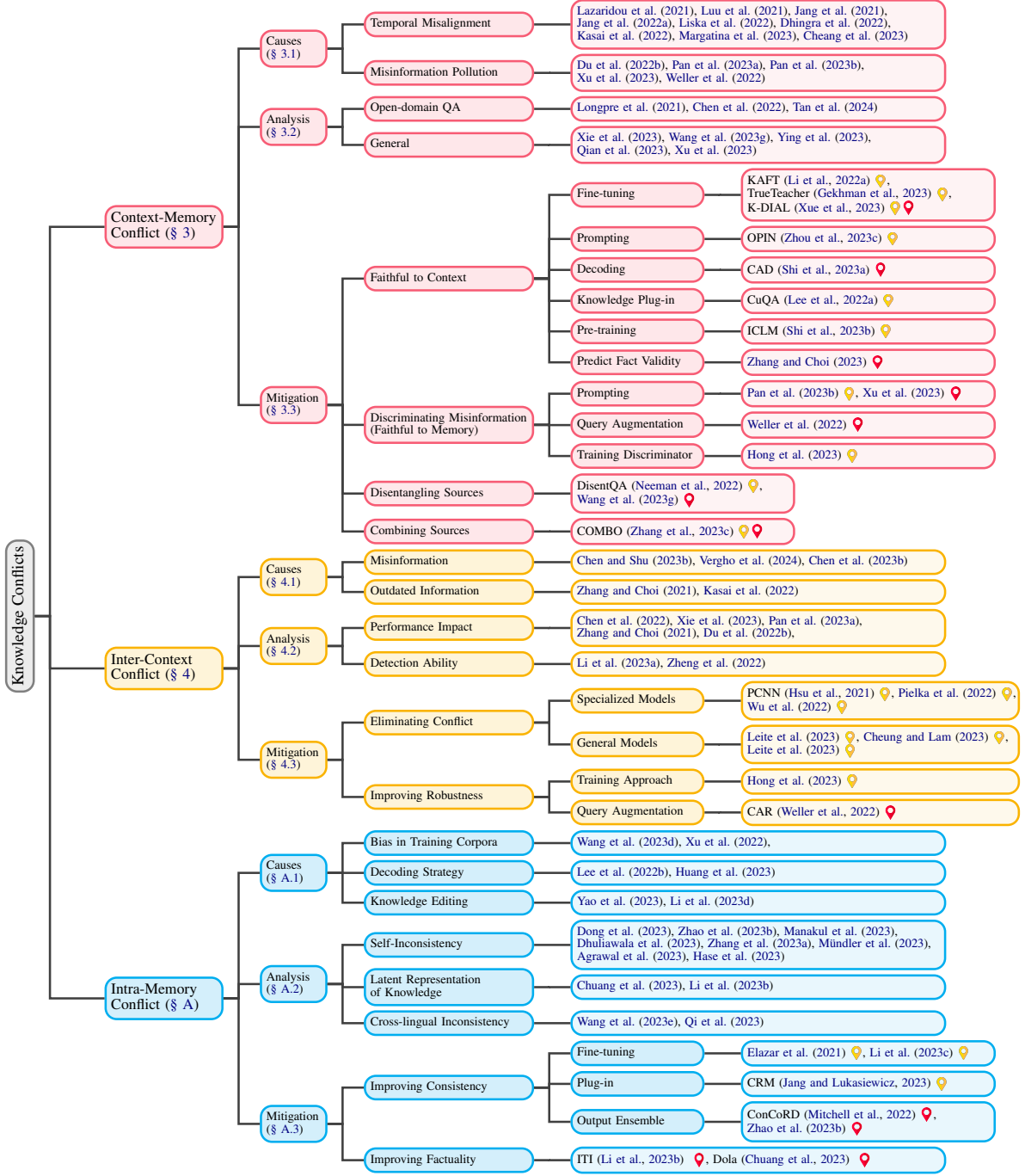


Figure 3: Taxonomy of knowledge conflicts. We mainly list works in the era of large language models. ♀ denotes pre-hoc mitigation strategy and ♀ denotes post-hoc mitigation strategy.

(e.g., database, Internet) to aid the model and maintains its parameters unchanged (Karpukhin et al., 2020; Guu et al., 2020; Lewis et al., 2020; Lazaridou et al., 2022; Borgeaud et al., 2022; Peng et al., 2023; Vu et al., 2023). *Continue learning* (CL) aims to update the internal knowledge through continual pre-training on new and updated data (Lazaridou et al., 2021; Jang et al., 2021, 2022a).

However, these methods on mitigating temporal misalignment are not magic bullets. KE can

bring in side effects of knowledge conflict, leading to knowledge inconsistency (i.e., a sort of intra-memory conflict) and may even enhance the hallucination of LLMs (Li et al., 2023d; Pinter and Elhadad, 2023). CL suffers from catastrophic forgetting issues and also is computationally-intensive (De Lange et al., 2021; He et al., 2021; Wang et al., 2023f). For RAG, it is inevitable to encounter knowledge conflicts since the model’s parameters are not updated (Chen et al., 2021; Zhang



and Choi, 2021).

**Misinformation Pollution.** Misinformation pollution emerges as another contributor to context-memory conflict, particularly for time-invariant knowledge (Jang et al., 2021) that the model has learned accurately. Adversaries exploit this vulnerability by injecting false or misleading information into both the Web corpus of retrieved documents (Pan et al., 2023a,b; Weller et al., 2022) and user conversations (Xu et al., 2023). The latter poses a practical threat, as adversaries can leverage techniques such as *prompt injection* attacks (Liu et al., 2023b; Greshake et al., 2023; Yi et al., 2023). In these adversarial contexts, where the information is factually incorrect, the model may face severe consequences if it unquestioningly accepts opinions present in the context (Xie et al., 2023; Pan et al., 2023b; Xu et al., 2023).

Researchers observe that fabricated, malicious misinformation can markedly decrease the accuracy of automated fact checkers (Du et al., 2022b) and ODQA models (Pan et al., 2023a,b). Recent studies also highlight the model’s tendency to align with user opinions, *a.k.a.*, *sycophancy*, further exacerbating the issue (Perez et al., 2022; Turpin et al., 2023; Wei et al., 2023; Sharma et al., 2023).

Furthermore, in the current landscape of LLMs, there is growing apprehension in the NLP community regarding the potential generation of misinformation by LLMs (Ayoobi et al., 2023; Kidd and Birhane, 2023; Carlini et al., 2023; Zhou et al., 2023b; Spitale et al., 2023; Chen and Shu, 2023b). Researchers acknowledge the challenges associated with detecting misinformation generated by LLMs (Tang et al., 2023; Chen and Shu, 2023a; Jiang et al., 2023). This underscores the urgency of addressing the nuanced challenges posed by LLMs in the context of contextual misinformation.

## 3.2 Analysis

*How do LLMs navigate context-memory conflicts?*

This section will detail the relevant research, although they present quite different answers. Depending on the scenario, we first introduce the Open-domain question answering (ODQA) setup and then focus on general setups.

**ODQA.** In earlier ODQA literature, Longpre et al. (2021) explore how QA models act when the provided contextual information contradicts the learned information. The authors create an automated framework that identifies QA instances with

named entity answers, then substitutes mentions of the entity in the gold document with an alternate entity, thus creating the conflict context. This study reveals a tendency of these models to over-rely on parametric knowledge. Chen et al. (2022) revisits this setup while reporting differing observations, they note that models predominantly rely on contextual knowledge in their best-performing settings. They attribute this divergence in findings to two factors. Firstly, the entity substitution approach used by Longpre et al. (2021) potentially reduces the semantic coherence of the perturbed passages. Secondly, Longpre et al. (2021) based their research on single evidence passages, as opposed to Chen et al. (2022), who utilize multiple ones. Recently, with the emergence of really “large” language models such as ChatGPT (Ouyang et al., 2022; OpenAI, 2023a) and Llama (Touvron et al., 2023), *inter alia*, researchers re-examined this issue. Tan et al. (2024) examine how LLMs blend retrieved context with generated knowledge in the ODQA setup, and discover models tend to favor the parametric knowledge, influenced by the greater resemblance of these generated contexts to the input questions and the often incomplete nature of the retrieved information, especially within the scope of conflicting sources.

**General.** Xie et al. (2023) leverage LLMs to generate conflicting context alongside the memorized knowledge. They find that LLMs are highly receptive to external evidence, even when it conflicts with their parametric, provided that the external knowledge is coherent and convincing. Meanwhile, they also identify a strong confirmation bias (Nickerson, 1998) in LLMs, *i.e.*, the models tend to favor information consistent with their internal memory, even when confronted with conflicting external evidence. Wang et al. (2023g) posit that the desired behaviors when an LLM encounters conflicts should be to pinpoint the conflicts and provide distinct answers. They find while LLMs perform well in identifying the existence of knowledge conflicts, they struggle to determine the specific conflicting segments and produce a response with distinct answers amidst conflicting information. Ying et al. (2023) analyze the robustness of LLMs under conflicts with a focus on two perspectives: factual robustness (the ability to identify correct facts from prompts or memory) and decision style (categorizing LLMs’ behavior as intuitive, dependent, or rational-based on cognitive theory). The study finds that LLMs

are highly susceptible to misleading prompts, especially in the context of commonsense knowledge. Qian et al. (2023) evaluate the potential interaction between parametric and external knowledge more systematically, cooperating knowledge graph (KG). They reveal that LLMs often deviate from their parametric knowledge when presented with direct conflicts or detailed contextual changes. Xu et al. (2023) study how large language models (LLMs) respond to knowledge conflicts during interactive sessions. Their findings suggest LLMs tend to favor logically structured knowledge, even when it contradicts factual accuracy.

### 3.3 Mitigation

Mitigation strategies are organized according to their **objectives**, *i.e.*, the desired behaviors we expect from an LLM when it encounters conflicts. Strategies are categorized to the following objectives: *Faithful to context* strategies aim to align with contextual knowledge, focusing on context prioritization. *Discriminating misinformation* strategies encourage skepticism towards dubious context in favor of parametric knowledge. *Disentangling sources* strategies treat context and knowledge separately. *Combining sources* strategies aim for an integrated response leveraging both context and parametric knowledge.

**Faithful to Context.** *Fine-tuning.* Li et al. (2022a) argues that an LLM should prioritize context for task-relevant information and rely on internal knowledge when the context is unrelated. They name the two properties controllability and robustness. They introduce Knowledge Aware Fine-Tuning (KAFT) to strengthen the two properties by incorporating counterfactual and irrelevant contexts to standard training datasets. Gekhman et al. (2023) introduce TrueTeacher, which focuses on improving factual consistency in summarization by annotating model-generated summaries with LLMs. This approach helps in maintaining faithfulness to the context of the original documents, ensuring that generated summaries remain accurate without being misled by irrelevant or incorrect details. DIAL (Xue et al., 2023) focuses on improving factual consistency in dialogue systems via direct knowledge enhancement and RLFC for aligning responses accurately with provided factual knowledge.

*Prompting.* Zhou et al. (2023c) explores enhancing LLMs’ adherence to context through specialized

prompting strategies, specifically opinion-based prompts and counterfactual demonstrations. These techniques are shown to significantly improve LLMs’ performance in context-sensitive tasks by ensuring they remain faithful to relevant context, without additional training.

*Decoding.* Shi et al. (2023a) introduces Context-aware Decoding (CAD) to reduce hallucinations by amplifying the difference in output probabilities with and without context. CAD enhances faithfulness in LLMs by effectively prioritizing relevant context over the model’s prior knowledge, especially in tasks with conflicting information.

*Knowledge Plug-in.* Lee et al. (2022a) proposes Continuously-updated QA (CuQA) for improving LMs’ ability to integrate new knowledge. Their approach uses plug-and-play modules to store updated knowledge, ensuring the original model remains unaffected. Unlike traditional continue pre-training or fine-tuning approaches, CuQA can solve knowledge conflicts.

*Pre-training.* ICLM (Shi et al., 2023b) is a new pre-training method that extends LLMs’ ability to handle long and varied contexts across multiple documents. This approach could potentially aid in resolving knowledge conflicts by enabling models to synthesize information from broader contexts, thus improving their understanding and application of relevant knowledge.

*Predict Fact Validity.* Zhang and Choi (2023) addresses knowledge conflict by introducing fact duration prediction to identify and discard outdated facts in LLMs. This approach improves model performance on tasks like ODQA by ensuring adherence to up-to-date contextual information.

**Discriminating Misinformation (Faithful to Memory).** *Prompting.* To address misinformation pollution, Pan et al. (2023b) proposes defense strategies such as misinformation detection and vigilant prompting, aiming to enhance the model’s ability to remain faithful to factual, parametric information amidst potential misinformation. Similarly, Xu et al. (2023) utilizes a system prompt to remind the LLM to be cautious about potential misinformation and to verify its memorized knowledge before responding. This approach aims to enhance the LLM’s ability to maintain faithfulness.

*Query Augmentation.* Weller et al. (2022) leverages the redundancy of information in large corpora to defend misinformation pollution. Their method involves query augmentation to find a diverse set of

less likely poisoned passages, coupled with a confidence method named Confidence from Answer Redundancy (CAR), which compares the predicted answer’s consistency across retrieved contexts. This strategy mitigates knowledge conflicts by ensuring the model’s faithfulness through cross-verification of answers in multiple sources.

**Training Discriminator.** Hong et al. (2023) fine-tune a smaller LM as discriminator and combine prompting techniques to develop the model’s ability to discriminate between reliable and unreliable information, helping the model remain faithful when confronted with misleading context.

**Disentangling Sources.** DisentQA (Neeman et al., 2022) trains a model that predicts two types of answers for a given question: one based on contextual knowledge and one on parametric knowledge. Wang et al. (2023g) introduce a method to improve Large Language Models’ (LLMs) handling of knowledge conflicts. Their approach is a three-step process designed to help LLMs detect conflicts, accurately identify the conflicting segments, and generate distinct, informed responses based on the conflicting data, aiming for more precise and nuanced model outputs.

**Combining Sources.** Zhang et al. (2023c) propose COMBO, a framework that pairs compatible generated and retrieved passages to resolve discrepancies. It uses discriminators trained on silver labels to assess passage compatibility, improving ODQA performance by leveraging both LLM-generated (parametric) and external retrieved knowledge.

## 4 Inter-Context Conflict

### 4.1 Causes

**Misinformation.** Misinformation has long been a significant concern in the modern digital age (Shu et al., 2017; Zubiaga et al., 2018; Kumar and Shah, 2018; Meel and Vishwakarma, 2020; Fung et al., 2022; Wang et al., 2023b). Currently, the use of retrieve augment generate (RAG) LLMs has emerged as a novel paradigm in the field. However, the incorporation of retrieved documents as external knowledge introduces a noteworthy concern – the potential for misinformation like fake news within these retrieved documents (Chen et al., 2023b). In the past, there have been instances of AI being utilized for the generation of misinformation (Zhou et al., 2023b; Verghe et al., 2024; Weidinger et al., 2021). With the formidable generative capabilities of LLMs, this issue has been further exacerbated,

contributing to a significant surge in misinformation generated by LLMs (Chen and Shu, 2023b; Menczer et al., 2023; Barrett et al., 2023; Bengio et al., 2023; Wang et al., 2023c; Solaiman et al., 2023; Weidinger et al., 2023; Ferrara, 2023; Goldstein et al., 2023).

**Outdated Information.** In addition to the challenge of misinformation, it is important to recognize that facts can evolve over time. The retrieved documents may contain updated and outdated information from the network simultaneously (Chen et al., 2021; Liska et al., 2022; Zhang and Choi, 2021; Kasai et al., 2022).

### 4.2 Analysis

**Performance Impact.** Previous research has empirically demonstrated that the performance of a pre-trained language model can be significantly influenced by the presence of misinformation (Zhang and Choi, 2021) or outdated information (Du et al., 2022b) within a specific context. In a recent study, Pan et al. (2023a) introduced a misinformation attack strategy involving the creation of fabricated versions of Wikipedia articles, which are subsequently inserted into the authentic Wikipedia corpus. Their research findings revealed that existing language models are susceptible to misinformation attacks, irrespective of whether the fake articles are manually crafted or generated by models. Notably, when the retrieval dataset includes 4% of misinformation, the model’s Exact Match (EM) performance drops by approximately 10%. To gain a deeper understanding of how LLMs behave when encountering contradictory contexts, Chen et al. (2022) primarily conducted experiments using Fusion-in-Decoder on the NQ-Open (Kwiatkowski et al., 2019) and TriviaQA (Joshi et al., 2017). They found that contradictions within knowledge sources have a marginal impact on models’ confidence and models tend to prioritize context that is more relevant to the query and context that contains answers consistent with the model’s parametric knowledge. Xie et al. (2023) conducted experiments on both closed-source LLMs and open-source LLMs in POPQA (Mallen et al., 2022) and StrategyQA (Geva et al., 2021). The results obtained were in line with those of Chen et al. (2022), indicating that LLMs exhibit a significant bias to evidence that aligns with the model’s parametric memory. They also found that LLMs tend to place stronger emphasis on facts related to more popular



entities and answers supported by more documents in the context, and are highly sensitive to the order in which information is presented.

**Detection Ability.** In addition to assessing the performance of LLMs when confronted with contradictory contexts, several studies also investigate their capacity to identify such contradictions. [Zheng et al. \(2022\)](#) examines the performance of various models including BERT, RoBERTa, and ERNIE in detecting the contradiction within Chinese conversations. Their experimental findings reveal that identifying contradictory statements within a conversation is a significant challenge for these models. [Li et al. \(2023a\)](#) analyse the performance of GPT-4, ChatGPT, PaLM-2, and LLaMAv2 in identifying contradictory documents within news articles ([Hermann et al., 2015](#)), stories ([Kočíský et al., 2018](#)), and wikipedia ([Merity et al., 2016](#)). The authors found that, even for GPT-4, the average detection accuracy remains at around 70%. The study also revealed that LLMs encounter particular challenges when dealing with specific types of contradiction, notably those related to subjective emotions or perspectives and the length of the documents and the range of self-contradictions have a slight impact on the detection performance.

### 4.3 Mitigation

**Eliminating Conflict. Specialized Models.** [Hsu et al. \(2021\)](#) develop a model named Pairwise Contradiction Neural Network (PCNN), which generates contradiction probabilities based on sentence representations obtained through fine-tuned Sentence-BERT. They conducted on Wikipedia various proportions of misinformation to demonstrate the effectiveness of their approach. [Pielka et al. \(2022\)](#) discovered that XLM-RoBERTa struggles to effectively grasp the syntactic and semantic features that contribute to incorrect contradiction detection. They suggested incorporating linguistic knowledge into the learning process. [Wu et al. \(2022\)](#) developed a novel approach that integrates topological representations of text into deep learning models. They performed experiments using three widely-used models including BERT, ESIM, and CBOW on MultiNLI dataset ([Williams et al., 2017](#)). The results provide compelling evidence that their approach is effective.

**General Models.** [Chern et al. \(2023\)](#) proposed a fact-checking framework that integrates LLMs including GPT4, GPT3.5 and FLAN-T5-XXL with

various tools, such as Google Search, Google Scholar, code interpreters, and Python, for detecting factual errors in texts. The authors conducted experiments using RoSE ([Liu et al., 2022](#)) and Fact-Prompts. [Cheung and Lam \(2023\)](#) combined the search engine with Llama to predict the veracity of claims. They conducted experiments using two datasets, RAWFC and LIAR. [Leite et al. \(2023\)](#) employ LLMs to produce weak labels associated with 18 credibility signals for the input text and aggregate these labels through weak supervision techniques to make predictions regarding the veracity of the input. Effectiveness of their method on the FA-KES and EUvsDisinfo datasets.

**Improving Robustness. Training Approach.** [Hong et al. \(2023\)](#) presents a novel fine-tuning method that involves training a discriminator alongside the decoder using the same encoder of FiD. Additionally, the author introduces two other methods to improve the robustness of the model including prompting GPT-3 to identify perturbed documents before generating responses and integrating the discriminator’s output into the prompt for GPT-3. They conduct experiments on NQ-Open with entity replacement operations ([Longpre et al., 2021](#)) and their experimental results indicate that the fine-tuning method yields the most promising results.

**Query Augmentation.** [Weller et al. \(2022\)](#) first prompt GPT-3 to generate new questions based on the original question and then measure the confidence for each query and its corresponding retrieved passages. This confidence is used to determine whether to use the original question’s prediction or to opt for a majority vote based on the predictions obtained from the augmented questions that exhibit a high degree of confidence. They verify the effectiveness of their method on Natural Questions and TriviaQA.

## 5 Conclusion

Through this survey, we have highlighted the importance of investigating knowledge conflicts, shedding light on their categorization, causes, behavioral patterns, and potential mitigation approaches. We have observed that existing studies have started to recognize the significance of knowledge conflicts but lack a systematic exploration dedicated solely to this topic. Our survey aims to fill this gap by providing a comprehensive review, synthesizing relevant literature, and offering insights into potential avenues for further research and resolution.



## Limitations

Given the rapid growth of research in the field of knowledge conflict and the abundance of research literature, it is inevitable that we may overlook some most recent or less relative findings. However, we have included all the most essential materials in our survey.

## Ethics Statement

No Consideration.

## References

- Abdalghani Abujabal, Rishiraj Saha Roy, Mohamed Yahya, and Gerhard Weikum. 2018. Comqa: A community-sourced dataset for complex factoid question answering with paraphrase clusters. *arXiv preprint arXiv:1809.09528*.
- Shourya Aggarwal, Divyanshu Mandowara, Vishwajeet Agrawal, Dinesh Khandelwal, Parag Singla, and Dinesh Garg. 2021. Explanations for commonsenseqa: New dataset and models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3050–3065.
- Ayush Agrawal, Lester Mackey, and Adam Tauman Kalai. 2023. Do language models know when they’re hallucinating references? *arXiv preprint arXiv:2305.18248*.
- Ekin Akyürek, Tolga Bolukbasi, Frederick Liu, Binbin Xiong, Ian Tenney, Jacob Andreas, and Kelvin Guu. 2022. Towards tracing knowledge in language models back to the training data. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2429–2446.
- Navid Ayoobi, Sadat Shahriar, and Arjun Mukherjee. 2023. The looming threat of fake and llm-generated linkedin profiles: Challenges and opportunities for detection and prevention. In *Proceedings of the 34th ACM Conference on Hypertext and Social Media*, pages 1–10.
- Clark Barrett, Brad Boyd, Elie Bursztein, Nicholas Carlini, Brad Chen, Jihye Choi, Amrita Roy Chowdhury, Mihai Christodorescu, Anupam Datta, Soheil Feizi, et al. 2023. Identifying and mitigating the security risks of generative ai. *Foundations and Trends® in Privacy and Security*, 6(1):1–52.
- Henning Bartsch, Ole Jorgensen, Domenic Rosati, Jason Hoelscher-Obermaier, and Jacob Pfau. 2023. Self-consistency of large language models under ambiguity. *arXiv preprint arXiv:2310.13439*.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the

- dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages 610–623.
- Yoshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Yuval Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, Gillian Hadfield, et al. 2023. Managing ai risks in an era of rapid progress. *arXiv preprint arXiv:2310.17688*.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millan, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Nicholas Carlini, Matthew Jagielski, Christopher A Choquette-Choo, Daniel Paleka, Will Pearce, Hyrum Anderson, Andreas Terzis, Kurt Thomas, and Florian Tramèr. 2023. Poisoning web-scale training datasets is practical. *arXiv preprint arXiv:2302.10149*.
- Tyler A Chang and Benjamin K Bergen. 2023. Language model behavior: A comprehensive survey. *arXiv preprint arXiv:2303.11504*.
- Chi Cheang, Hou Chan, Derek Wong, Xuebo Liu, Zhaocong Li, Yanming Sun, Shudong Liu, and Lidia Chao. 2023. Can lms generalize to future data? an empirical analysis on text summarization. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16205–16217.
- Canyu Chen and Kai Shu. 2023a. Can llm-generated misinformation be detected? In *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*.
- Canyu Chen and Kai Shu. 2023b. Combating misinformation in the age of llms: Opportunities and challenges. *arXiv preprint arXiv:2311.05656*.
- 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*.
- Jiangjie Chen, Wei Shi, Ziquan Fu, Sijie Cheng, Lei Li, and Yanghua Xiao. 2023a. Say what you mean! large language models speak too positively about negative commonsense knowledge. *arXiv preprint arXiv:2305.05976*.
- Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2023b. Benchmarking large language models in retrieval-augmented generation. *arXiv preprint arXiv:2309.01431*.

747	Liang Chen, Yang Deng, Yatao Bian, Zeyu Qin, Bingzhe	Bhuwan Dhingra, Jeremy R Cole, Julian Martin	802
748	Wu, Tat-Seng Chua, and Kam-Fai Wong. 2023c. Be-	Eisenschlos, Daniel Gillick, Jacob Eisenstein, and	803
749	beyond factuality: A comprehensive evaluation of large	William W Cohen. 2022. Time-aware language mod-	804
750	language models as knowledge generators. In <i>Pro-</i>	els as temporal knowledge bases. <i>Transactions of the</i>	805
751	<i>ceedings of the 2023 Conference on Empirical Meth-</i>	<i>Association for Computational Linguistics</i> , 10:257–	806
752	<i>ods in Natural Language Processing</i> , pages 6325–	273.	807
753	6341.		
754	Wenhu Chen, Xinyi Wang, and William Yang Wang.	Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu,	808
755	2021. A dataset for answering time-sensitive ques-	Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Ja-	809
756	tions. <i>arXiv preprint arXiv:2108.06314</i> .	son Weston. 2023. Chain-of-verification reduces hal-	810
757		lucination in large language models. <i>arXiv preprint</i>	811
758	I Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua	<i>arXiv:2309.11495</i> .	812
759	Feng, Chunting Zhou, Junxian He, Graham Neubig,		
760	Pengfei Liu, et al. 2023. Factool: Factuality detec-	Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Zhifang	813
761	tion in generative ai—a tool augmented framework	Sui, and Lei Li. 2023. Statistical knowledge assess-	814
762	for multi-task and multi-domain scenarios. <i>arXiv</i>	ment for large language models. In <i>Thirty-seventh</i>	815
	<i>preprint arXiv:2307.13528</i> .	<i>Conference on Neural Information Processing Sys-</i>	816
		<i>tems</i> .	817
763	Tsun-Hin Cheung and Kin-Man Lam. 2023. Factllama:	Li Du, Xiao Ding, Kai Xiong, Ting Liu, and Bing Qin.	818
764	Optimizing instruction-following language models	2022a. e-care: a new dataset for exploring explain-	819
765	with external knowledge for automated fact-checking.	able causal reasoning. In <i>Proceedings of the 60th</i>	820
766	In <i>2023 Asia Pacific Signal and Information Pro-</i>	<i>Annual Meeting of the Association for Computational</i>	821
767	<i>cessing Association Annual Summit and Conference</i>	<i>Linguistics (Volume 1: Long Papers)</i> , pages 432–446.	822
768	( <i>APSIPA ASC</i> ), pages 846–853. IEEE.		
769	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,	Yibing Du, Antoine Bosselut, and Christopher D Man-	823
770	Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul	ning. 2022b. Synthetic disinformation attacks on	824
771	Barham, Hyung Won Chung, Charles Sutton, Sebas-	automated fact verification systems. In <i>Proceedings</i>	825
772	tian Gehrmann, et al. 2023. Palm: Scaling language	<i>of the AAAI Conference on Artificial Intelligence</i> ,	826
773	modeling with pathways. <i>Journal of Machine Learn-</i>	volume 36, pages 10581–10589.	827
774	<i>ing Research</i> , 24(240):1–113.		
775	Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon	Nouha Dziri, Andrea Madotto, Osmar Zaiane, and	828
776	Kim, James Glass, and Pengcheng He. 2023. Dola:	Avishek Joey Bose. 2021. Neural path hunter: Re-	829
777	Decoding by contrasting layers improves factu-	ducing hallucination in dialogue systems via path	830
778	ality in large language models. <i>arXiv preprint</i>	grounding. <i>arXiv preprint arXiv:2104.08455</i> .	831
779	<i>arXiv:2309.03883</i> .		
780	Christopher Clark, Kenton Lee, Ming-Wei Chang,	Yanai Elazar, Nora Kassner, Shauli Ravfogel, Amir	832
781	Tom Kwiatkowski, Michael Collins, and Kristina	Feder, Abhilasha Ravichander, Marius Mosbach,	833
782	Toutanova. 2019. Boolq: Exploring the surprising	Yonatan Belinkov, Hinrich Schütze, and Yoav Gold-	834
783	difficulty of natural yes/no questions. In <i>Proceedings</i>	berg. 2022. Measuring causal effects of data statis-	835
784	<i>of the 2019 Conference of the North American Chap-</i>	tics on language model’s factual’predictions. <i>arXiv</i>	836
785	<i>ter of the Association for Computational Linguistics:</i>	<i>preprint arXiv:2207.14251</i> .	837
786	<i>Human Language Technologies, Volume 1 (Long and</i>		
787	<i>Short Papers)</i> , pages 2924–2936.	Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhi-	838
788	Leyang Cui, Sijie Cheng, Yu Wu, and Yue Zhang. 2020.	lasha Ravichander, Eduard Hovy, Hinrich Schütze,	839
789	Does bert solve commonsense task via commonsense	and Yoav Goldberg. 2021. Measuring and improving	840
790	knowledge. <i>arXiv preprint arXiv:2008.03945</i> , 4.	consistency in pretrained language models. <i>Transac-</i>	841
791	Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Edit-	<i>tions of the Association for Computational Linguis-</i>	842
792	ing factual knowledge in language models. In <i>Pro-</i>	<i>tics</i> , 9:1012–1031.	843
793	<i>ceedings of the 2021 Conference on Empirical Meth-</i>		
794	<i>ods in Natural Language Processing</i> , pages 6491–	Hady Elsahar, Pavlos Vougiouklis, Arslan Remaci,	844
795	6506.	Christophe Gravier, Jonathon Hare, Frederique Lafor-	845
796	Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah	est, and Elena Simperl. 2018. T-rex: A large scale	846
797	Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh,	alignment of natural language with knowledge base	847
798	and Tinne Tuytelaars. 2021. A continual learning sur-	triples. In <i>Proceedings of the Eleventh International</i>	848
799	vey: Defying forgetting in classification tasks. <i>IEEE</i>	<i>Conference on Language Resources and Evaluation</i>	849
800	<i>transactions on pattern analysis and machine intelli-</i>	( <i>LREC 2018</i> ).	850
801	<i>gence</i> , 44(7):3366–3385.	Angela Fan, Mike Lewis, and Yann Dauphin. 2018.	851
		Hierarchical neural story generation. <i>arXiv preprint</i>	852
		<i>arXiv:1805.04833</i> .	853
		Zhangyin Feng, Weitao Ma, Weijiang Yu, Lei Huang,	854
		Haotian Wang, Qianglong Chen, Weihua Peng, Xi-	855
		aocheng Feng, Bing Qin, et al. 2023. Trends in inte-	856
		gration of knowledge and large language models: A	857

858	survey and taxonomy of methods, benchmarks, and	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasu-	914
859	applications. <i>arXiv preprint arXiv:2311.05876</i> .	pat, and Mingwei Chang. 2020. Retrieval augmented	915
860	Emilio Ferrara. 2023. Genai against humanity: Ne-	language model pre-training. In <i>International confer-</i>	916
861	furious applications of generative artificial intelli-	<i>ence on machine learning</i> , pages 3929–3938. PMLR.	917
862	gence and large language models. <i>arXiv preprint</i>		
863	<i>arXiv:2310.00737</i> .	Peter Hase, Mona Diab, Asli Celikyilmaz, Xian Li, Zor-	918
864	Yi R Fung, Kung-Hsiang Huang, Preslav Nakov, and	nitsa Kozareva, Veselin Stoyanov, Mohit Bansal, and	919
865	Heng Ji. 2022. The battlefield of combating misinfor-	Srinivasan Iyer. 2023. Methods for measuring, up-	920
866	mation and coping with media bias. In <i>Proceedings</i>	dating, and visualizing factual beliefs in language	921
867	<i>of the 28th ACM SIGKDD Conference on Knowledge</i>	models. In <i>Proceedings of the 17th Conference of</i>	922
868	<i>Discovery and Data Mining</i> , pages 4790–4791.	<i>the European Chapter of the Association for Compu-</i>	923
869	Wee Chung Gan and Hwee Tou Ng. 2019. Improv-	<i>tational Linguistics</i> , pages 2706–2723.	924
870	ing the robustness of question answering systems to		
871	question paraphrasing. In <i>Proceedings of the 57th</i>	Tianxing He, Jun Liu, Kyunghyun Cho, Myle Ott, Bing	925
872	<i>annual meeting of the association for computational</i>	Liu, James Glass, and Fuchun Peng. 2021. Analyzing	926
873	<i>linguistics</i> , pages 6065–6075.	the forgetting problem in pretrain-finetuning of open-	927
874	Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony	domain dialogue response models. In <i>Proceedings</i>	928
875	Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vin-	<i>of the 16th Conference of the European Chapter of</i>	929
876	cent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan,	<i>the Association for Computational Linguistics: Main</i>	930
877	et al. 2023a. Rarr: Researching and revising what	<i>Volume</i> , pages 1121–1133.	931
878	language models say, using language models. In		
879	<i>Proceedings of the 61st Annual Meeting of the As-</i>	Karl Moritz Hermann, Tomas Kocisky, Edward Grefen-	932
880	<i>sociation for Computational Linguistics (Volume 1:</i>	stette, Lasse Espeholt, Will Kay, Mustafa Suleyman,	933
881	<i>Long Papers)</i> , pages 16477–16508.	and Phil Blunsom. 2015. Teaching machines to read	934
882	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,	and comprehend. <i>Advances in neural information</i>	935
883	Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen	<i>processing systems</i> , 28.	936
884	Wang. 2023b. Retrieval-augmented generation for		
885	large language models: A survey. <i>arXiv preprint</i>	Jordan Hoffmann, Sebastian Borgeaud, Arthur Men-	937
886	<i>arXiv:2312.10997</i> .	sch, Elena Buchatskaya, Trevor Cai, Eliza Ruther-	938
887	Zorik Gekhman, Jonathan Herzig, Roei Aharoni, Chen	ford, Diego de Las Casas, Lisa Anne Hendricks,	939
888	Elkind, and Idan Szpektor. 2023. Trueteacher: Learn-	Johannes Welbl, Aidan Clark, et al. 2022. Train-	940
889	ing factual consistency evaluation with large lan-	ing compute-optimal large language models. <i>arXiv</i>	941
890	guage models. <i>arXiv preprint arXiv:2305.11171</i> .	<i>preprint arXiv:2203.15556</i> .	942
891	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot,	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and	943
892	Dan Roth, and Jonathan Berant. 2021. Did aristotle	Yejin Choi. 2019. The curious case of neural text	944
893	use a laptop? a question answering benchmark with	degeneration. <i>arXiv preprint arXiv:1904.09751</i> .	945
894	implicit reasoning strategies. <i>Transactions of the</i>		
895	<i>Association for Computational Linguistics</i> , 9:346–	Giwon Hong, Jeonghwan Kim, Junmo Kang, Sung-	946
896	361.	Hyon Myaeng, and Joyce Jiyoung Whang. 2023. Dis-	947
897	Josh A Goldstein, Girish Sastry, Micah Musser, Ree-	cern and answer: Mitigating the impact of misinfor-	948
898	nee DiResta, Matthew Gentzel, and Katerina Sedova.	mation in retrieval-augmented models with discrimi-	949
899	2023. Generative language models and automated	nators. <i>arXiv preprint arXiv:2305.01579</i> .	950
900	influence operations: Emerging threats and potential		
901	mitigations. <i>arXiv preprint arXiv:2301.04246</i> .	Cheng Hsu, Cheng-Te Li, Diego Saez-Trumper, and Yi-	951
902	Kai Greshake, Sahar Abdelnabi, Shailesh Mishra,	Zhan Hsu. 2021. Wikicontradiction: Detecting self-	952
903	Christoph Endres, Thorsten Holz, and Mario Fritz.	contradiction articles on wikipedia. In <i>2021 IEEE</i>	953
904	2023. More than you’ve asked for: A comprehen-	<i>International Conference on Big Data (Big Data)</i> ,	954
905	sive analysis of novel prompt injection threats to	pages 427–436. IEEE.	955
906	application-integrated large language models. <i>arXiv</i>		
907	<i>e-prints</i> , pages arXiv–2302.	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong,	956
908	Roger Grosse, Juhan Bae, Cem Anil, Nelson El-	Zhangyin Feng, Haotian Wang, Qianglong Chen,	957
909	hage, Alex Tamkin, Amirhossein Tajdini, Benoit	Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023.	958
910	Steiner, Dustin Li, Esin Durmus, Ethan Perez, et al.	A survey on hallucination in large language models:	959
911	2023. Studying large language model general-	Principles, taxonomy, challenges, and open questions.	960
912	ization with influence functions. <i>arXiv preprint</i>	<i>arXiv preprint arXiv:2311.05232</i> .	961
913	<i>arXiv:2308.03296</i> .	Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Worts-	962
		man, Suchin Gururangan, Ludwig Schmidt, Han-	963
		naneh Hajishirzi, and Ali Farhadi. 2022. Edit-	964
		ing models with task arithmetic. <i>arXiv preprint</i>	965
		<i>arXiv:2212.04089</i> .	966
		Joel Jang, Seonghyeon Ye, Changho Lee, Sohee Yang,	967
		Joongbo Shin, Janghoon Han, Gyeonghun Kim, and	968
		Minjoon Seo. 2022a. Temporalwiki: A lifelong	969

970	benchmark for training and evaluating ever-evolving	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick	1023
971	language models. In <i>Proceedings of the 2022 Con-</i>	Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and	1024
972	<i>ference on Empirical Methods in Natural Language</i>	Wen-tau Yih. 2020. Dense passage retrieval for open-	1025
973	<i>Processing</i> , pages 6237–6250.	domain question answering. In <i>Proceedings of the</i>	1026
974	Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin,	<i>2020 Conference on Empirical Methods in Natural</i>	1027
975	Janghoon Han, KIM Gyeonghun, Stanley Jungkyu	<i>Language Processing (EMNLP)</i> , pages 6769–6781.	1028
976	Choi, and Minjoon Seo. 2021. Towards continual	Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi,	1029
977	knowledge learning of language models. In <i>Internat-</i>	Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir	1030
978	<i>ional Conference on Learning Representations</i> .	Radev, Noah A Smith, Yejin Choi, and Kentaro Inui.	1031
979	Myeongjun Jang, Deuk Sin Kwon, and Thomas	2022. Realtime qa: What’s the answer right now?	1032
980	Lukasiewicz. 2022b. Becel: Benchmark for consis-	<i>arXiv preprint arXiv:2207.13332</i> .	1033
981	tency evaluation of language models. In <i>Proceedings</i>	Nora Kassner, Oyvind Tafjord, Hinrich Schütze, and	1034
982	<i>of the 29th International Conference on Computa-</i>	Peter Clark. 2021. Beliefbank: Adding memory to a	1035
983	<i>tional Linguistics</i> , pages 3680–3696.	pre-trained language model for a systematic notion	1036
984	Myeongjun Erik Jang and Thomas Lukasiewicz. 2023.	of belief. <i>arXiv preprint arXiv:2109.14723</i> .	1037
985	Improving language models meaning understanding	Celeste Kidd and Abeba Birhane. 2023. How ai can dis-	1038
986	and consistency by learning conceptual roles from	tort human beliefs. <i>Science</i> , 380(6651):1222–1223.	1039
987	dictionary. <i>arXiv preprint arXiv:2310.15541</i> .	Miyoung Ko, Ingyu Seong, Hwaran Lee, Joonsuk Park,	1040
988	Ganesh Jawahar, Muhammad Abdul-Mageed, and	Minsuk Chang, and Minjoon Seo. 2022. Claimdiff:	1041
989	Laks VS Lakshmanan. 2020. Automatic detection	Comparing and contrasting claims on contentious	1042
990	of machine generated text: A critical survey. <i>arXiv</i>	issues. <i>arXiv preprint arXiv:2205.12221</i> .	1043
991	<i>preprint arXiv:2011.01314</i> .	Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris	1044
992	Ganesh Jawahar, Benoît Sagot, and Djamé Seddah.	Dyer, Karl Moritz Hermann, Gábor Melis, and Ed-	1045
993	2019. What does bert learn about the structure of	ward Grefenstette. 2018. The narrativeqa reading	1046
994	language? In <i>ACL 2019-57th Annual Meeting of the</i>	comprehension challenge. <i>Transactions of the Asso-</i>	1047
995	<i>Association for Computational Linguistics</i> .	<i>ciation for Computational Linguistics</i> , 6:317–328.	1048
996	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan	Srikanth Kumar and Neil Shah. 2018. False information	1049
997	Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea	on web and social media: A survey. <i>arXiv preprint</i>	1050
998	Madotto, and Pascale Fung. 2023. Survey of halluci-	<i>arXiv:1804.08559</i> .	1051
999	nation in natural language generation. <i>ACM Comput-</i>	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	1052
1000	<i>ing Surveys</i> , 55(12):1–38.	field, Michael Collins, Ankur Parikh, Chris Alberti,	1053
1001	Bohan Jiang, Zhen Tan, Ayushi Nirmal, and Huan	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-	1054
1002	Liu. 2023. Disinformation detection: An evolving	ton Lee, et al. 2019. Natural questions: a benchmark	1055
1003	challenge in the age of llms. <i>arXiv preprint</i>	for question answering research. <i>Transactions of the</i>	1056
1004	<i>arXiv:2309.15847</i> .	<i>Association for Computational Linguistics</i> , 7:453–	1057
1005	Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke	466.	1058
1006	Zettlemoyer. 2017. Triviaqa: A large scale distantly	Tsz Kin Lam, Eva Hasler, and Felix Hieber. 2022. An-	1059
1007	supervised challenge dataset for reading comprehen-	alyzing the use of influence functions for instance-	1060
1008	sion. In <i>Proceedings of the 55th Annual Meeting of</i>	specific data filtering in neural machine translation.	1061
1009	<i>the Association for Computational Linguistics (Vol-</i>	<i>arXiv preprint arXiv:2210.13281</i> .	1062
1010	<i>ume 1: Long Papers)</i> , pages 1601–1611.	Angeliki Lazaridou, Elena Gribovskaya, Wojciech	1063
1011	Jean Kaddour, Joshua Harris, Maximilian Mozes, Her-	Stokowiec, and Nikolai Grigorev. 2022. Internet-	1064
1012	bie Bradley, Roberta Raileanu, and Robert McHardy.	augmented language models through few-shot	1065
1013	2023. Challenges and applications of large language	prompting for open-domain question answering.	1066
1014	models. <i>arXiv preprint arXiv:2307.10169</i> .	<i>arXiv preprint arXiv:2203.05115</i> .	1067
1015	Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric	Angeliki Lazaridou, Adhi Kuncoro, Elena Gribovskaya,	1068
1016	Wallace, and Colin Raffel. 2023. Large language	Devang Agrawal, Adam Liska, Tayfun Terzi, Mai	1069
1017	models struggle to learn long-tail knowledge. In <i>Inter-</i>	Gimenez, Cyprien de Masson d’Autume, Tomas Ko-	1070
1018	<i>national Conference on Machine Learning</i> , pages	ciskiy, Sebastian Ruder, et al. 2021. Mind the gap:	1071
1019	15696–15707. PMLR.	Assessing temporal generalization in neural language	1072
1020	Cheongwoong Kang and Jaesik Choi. 2023. Impact	models. <i>Advances in Neural Information Processing</i>	1073
1021	of co-occurrence on factual knowledge of large lan-	<i>Systems</i> , 34:29348–29363.	1074
1022	guage models. <i>arXiv preprint arXiv:2310.08256</i> .	Kyungjae Lee, Wookje Han, Seung-won Hwang,	1075
		Hwaran Lee, Joonsuk Park, and Sang-Woo Lee.	1076



1077	2022a. Plug-and-play adaptation for continuously-	Susannah Young, et al. 2022. Streamingqa: A bench-	1132
1078	updated qa. In <i>Findings of the Association for Com-</i>	mark for adaptation to new knowledge over time in	1133
1079	<i>putational Linguistics: ACL 2022</i> , pages 438–447.	question answering models. In <i>International Con-</i>	1134
1080	Nayeon Lee, Wei Ping, Peng Xu, Mostofa Patwary,	<i>ference on Machine Learning</i> , pages 13604–13622.	1135
1081	Pascale N Fung, Mohammad Shoeybi, and Bryan	PMLR.	1136
1082	Catanzaro. 2022b. Factuality enhanced language	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang,	1137
1083	models for open-ended text generation. <i>Advances in</i>	Hiroaki Hayashi, and Graham Neubig. 2023a. Pre-	1138
1084	<i>Neural Information Processing Systems</i> , 35:34586–	train, prompt, and predict: A systematic survey of	1139
1085	34599.	prompting methods in natural language processing.	1140
1086	João A Leite, Olesya Razuvaevskaya, Kalina	<i>ACM Computing Surveys</i> , 55(9):1–35.	1141
1087	Bontcheva, and Carolina Scarton. 2023. Detect-	Yi Liu, Gelei Deng, Yuekang Li, Kailong Wang, Tian-	1142
1088	ing misinformation with llm-predicted credibility	wei Zhang, Yepang Liu, Haoyu Wang, Yan Zheng,	1143
1089	signals and weak supervision. <i>arXiv preprint</i>	and Yang Liu. 2023b. Prompt injection attack	1144
1090	<i>arXiv:2309.07601</i> .	against llm-integrated applications. <i>arXiv preprint</i>	1145
1091	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	<i>arXiv:2306.05499</i> .	1146
1092	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	Shayne Longpre, Kartik Perisetla, Anthony Chen,	1147
1093	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	Nikhil Ramesh, Chris DuBois, and Sameer Singh.	1148
1094	täschel, et al. 2020. Retrieval-augmented generation	2021. Entity-based knowledge conflicts in question	1149
1095	for knowledge-intensive nlp tasks. <i>Advances in Neu-</i>	answering. In <i>Proceedings of the 2021 Conference</i>	1150
1096	<i>ral Information Processing Systems</i> , 33:9459–9474.	<i>on Empirical Methods in Natural Language Process-</i>	1151
1097	Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin	<i>ing</i> , pages 7052–7063.	1152
1098	Wang, Michal Lukasik, Andreas Veit, Felix Yu,	Kelvin Luu, Daniel Khashabi, Suchin Gururangan, Kar-	1153
1099	and Sanjiv Kumar. 2022a. Large language models	ishma Mandyam, and Noah A Smith. 2021. Time	1154
1100	with controllable working memory. <i>arXiv preprint</i>	waits for no one! analysis and challenges of temporal	1155
1101	<i>arXiv:2211.05110</i> .	misalignment. <i>arXiv preprint arXiv:2111.07408</i> .	1156
1102	Jierui Li, Vipul Raheja, and Dhruv Kumar. 2023a. Con-	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi	1157
1103	tradoc: Understanding self-contradictions in docu-	Das, Hannaneh Hajishirzi, and Daniel Khashabi.	1158
1104	ments with large language models. <i>arXiv preprint</i>	2022. When not to trust language models: Inves-	1159
1105	<i>arXiv:2311.09182</i> .	tigating effectiveness and limitations of paramet-	1160
1106	Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter	ric and non-parametric memories. <i>arXiv preprint</i>	1161
1107	Pfister, and Martin Wattenberg. 2023b. Inference-	<i>arXiv:2212.10511</i> , 7.	1162
1108	time intervention: Eliciting truthful answers from a	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das,	1163
1109	language model. <i>arXiv preprint arXiv:2306.03341</i> .	Daniel Khashabi, and Hannaneh Hajishirzi. 2023.	1164
1110	Shaobo Li, Xiaoguang Li, Lifeng Shang, Zhenhua Dong,	When not to trust language models: Investigating	1165
1111	Chengjie Sun, Bingquan Liu, Zhenzhou Ji, Xin Jiang,	effectiveness of parametric and non-parametric mem-	1166
1112	and Qun Liu. 2022b. How pre-trained language mod-	ories. In <i>Proceedings of the 61st Annual Meeting of</i>	1167
1113	els capture factual knowledge? a causal-inspired anal-	<i>the Association for Computational Linguistics (Vol-</i>	1168
1114	ysis. <i>arXiv preprint arXiv:2203.16747</i> .	<i>ume 1: Long Papers</i> ), pages 9802–9822.	1169
1115	Xiang Lisa Li, Vaishnavi Shrivastava, Siyan Li, Tat-	Potsawee Manakul, Adian Liusie, and Mark JF Gales.	1170
1116	sunori Hashimoto, and Percy Liang. 2023c. Bench-	2023. Selfcheckgpt: Zero-resource black-box hal-	1171
1117	marking and improving generator-validator consis-	lucination detection for generative large language	1172
1118	tency of language models. <i>arXiv preprint</i>	models. <i>arXiv preprint arXiv:2303.08896</i> .	1173
1119	<i>arXiv:2310.01846</i> .	Katerina Margatina, Shuai Wang, Yogarshi Vyas,	1174
1120	Zhoubo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang,	Neha Anna John, Yassine Benajiba, and Miguel	1175
1121	Xi Chen, and Huajun Chen. 2023d. Unveiling the pit-	Ballesteros. 2023. Dynamic benchmarking of	1176
1122	falls of knowledge editing for large language models.	masked language models on temporal concept	1177
1123	<i>arXiv preprint arXiv:2310.02129</i> .	drift with multiple views. <i>arXiv preprint</i>	1178
1124	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022.	<i>arXiv:2302.12297</i> .	1179
1125	Truthfulqa: Measuring how models mimic human	Luca Massarelli, Fabio Petroni, Aleksandra Piktus,	1180
1126	falsehoods. In <i>Proceedings of the 60th Annual Meet-</i>	Myle Ott, Tim Rocktäschel, Vassilis Plachouras, Fab-	1181
1127	<i>ing of the Association for Computational Linguistics</i>	rizio Silvestri, and Sebastian Riedel. 2019. How de-	1182
1128	<i>(Volume 1: Long Papers)</i> , pages 3214–3252.	coding strategies affect the verifiability of generated	1183
1129	Adam Liska, Tomas Kocisky, Elena Gribovskaya, Tay-	text. <i>arXiv preprint arXiv:1911.03587</i> .	1184
1130	fun Terzi, Eren Sezener, Devang Agrawal, D’Autume		
1131	Cyprien De Masson, Tim Scholtes, Manzil Zaheer,		

1185	Priyanka Meel and Dinesh Kumar Vishwakarma. 2020.	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	1238
1186	Fake news, rumor, information pollution in social me-	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	1239
1187	dia and web: A contemporary survey of state-of-the-	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	1240
1188	arts, challenges and opportunities. <i>Expert Systems</i>	2022. Training language models to follow instruc-	1241
1189	<i>with Applications</i> , 153:112986.	tions with human feedback. <i>Advances in Neural</i>	1242
		<i>Information Processing Systems</i> , 35:27730–27744.	1243
1190	Filippo Menczer, David Crandall, Yong-Yeol Ahn, and	Liangming Pan, Wenhui Chen, Min-Yen Kan, and	1244
1191	Apu Kapadia. 2023. Addressing the harms of ai-	William Yang Wang. 2023a. Attacking open-domain	1245
1192	generated inauthentic content. <i>Nature Machine Intel-</i>	question answering by injecting misinformation.	1246
1193	<i>ligence</i> , 5(7):679–680.	<i>IJCNLP-AAACL ACL</i> .	1247
1194	Kevin Meng, David Bau, Alex Andonian, and Yonatan	Xiaoman Pan, Wenlin Yao, Hongming Zhang, Dian Yu,	1248
1195	Belinkov. 2022. Locating and editing factual associ-	Dong Yu, and Jianshu Chen. 2022. Knowledge-in-	1249
1196	ations in gpt. <i>Advances in Neural Information Pro-</i>	context: Towards knowledgeable semi-parametric	1250
1197	<i>cessing Systems</i> , 35:17359–17372.	language models. In <i>The Eleventh International Con-</i>	1251
		<i>ference on Learning Representations</i> .	1252
1198	Stephen Merity, Caiming Xiong, James Bradbury, and	Yikang Pan, Liangming Pan, Wenhui Chen, Preslav	1253
1199	Richard Socher. 2016. Pointer sentinel mixture mod-	Nakov, Min-Yen Kan, and William Yang Wang.	1254
1200	els. <i>arXiv preprint arXiv:1609.07843</i> .	2023b. On the risk of misinformation pollu-	1255
		tion with large language models. <i>arXiv preprint</i>	1256
1201	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea	<i>arXiv:2305.13661</i> .	1257
1202	Finn, and Christopher D Manning. 2021. Fast model	Jungsoo Park, Sewon Min, Jaewoo Kang, Luke Zettle-	1258
1203	editing at scale. <i>arXiv preprint arXiv:2110.11309</i> .	moyer, and Hannaneh Hajishirzi. 2021. Faviq:	1259
		Fact verification from information-seeking questions.	1260
1204	Eric Mitchell, Joseph J Noh, Siyan Li, William S Arm-	<i>arXiv preprint arXiv:2107.02153</i> .	1261
1205	strong, Ananth Agarwal, Patrick Liu, Chelsea Finn,	Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng,	1262
1206	and Christopher D Manning. 2022. Enhancing self-	Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou	1263
1207	consistency and performance of pre-trained language	Yu, Weizhu Chen, et al. 2023. Check your facts and	1264
1208	models through natural language inference. <i>arXiv</i>	try again: Improving large language models with	1265
1209	<i>preprint arXiv:2211.11875</i> .	external knowledge and automated feedback. <i>arXiv</i>	1266
		<i>preprint arXiv:2302.12813</i> .	1267
1210	Niels Mündler, Jingxuan He, Slobodan Jenko, and Mar-	Ethan Perez, Sam Ringer, Kamilė Lukošiuotė, Karina	1268
1211	tin Vechev. 2023. Self-contradictory hallucinations	Nguyen, Edwin Chen, Scott Heiner, Craig Pettit,	1269
1212	of large language models: Evaluation, detection and	Catherine Olsson, Sandipan Kundu, Saurav Kada-	1270
1213	mitigation. <i>arXiv preprint arXiv:2305.15852</i> .	vath, et al. 2022. Discovering language model behav-	1271
1214	Humza Naveed, Asad Ullah Khan, Shi Qiu, Muham-	iors with model-written evaluations. <i>arXiv preprint</i>	1272
1215	mad Saqib, Saeed Anwar, Muhammad Usman, Nick	<i>arXiv:2212.09251</i> .	1273
1216	Barnes, and Ajmal Mian. 2023. A comprehensive	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel,	1274
1217	overview of large language models. <i>arXiv preprint</i>	Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and	1275
1218	<i>arXiv:2307.06435</i> .	Alexander Miller. 2019. Language models as knowl-	1276
		edge bases? In <i>Proceedings of the 2019 Confer-</i>	1277
1219	Ella Neeman, Roei Aharoni, Or Honovich, Leshem	<i>ence on Empirical Methods in Natural Language Pro-</i>	1278
1220	Choshen, Idan Szpektor, and Omri Abend. 2022.	<i>cessing and the 9th International Joint Conference</i>	1279
1221	Disentqa: Disentangling parametric and contextual	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	1280
1222	knowledge with counterfactual question answering.	pages 2463–2473.	1281
1223	<i>arXiv preprint arXiv:2211.05655</i> .	Maren Pielka, Felix Rode, Lisa Pucknat, Tobias Deuß,	1282
		and Rafet Sifa. 2022. A linguistic investigation of	1283
1224	Raymond S Nickerson. 1998. Confirmation bias: A	machine learning based contradiction detection mod-	1284
1225	ubiquitous phenomenon in many guises. <i>Review of</i>	els: an empirical analysis and future perspectives.	1285
1226	<i>general psychology</i> , 2(2):175–220.	In <i>2022 21st IEEE International Conference on Ma-</i>	1286
		<i>chine Learning and Applications (ICMLA)</i> , pages	1287
1227	Xenia Ohmer, Elia Bruni, and Dieuwke Hupkes. 2023.	1649–1653. IEEE.	1288
1228	Separating form and meaning: Using self-consistency	Yuval Pinter and Michael Elhadad. 2023. Emptying	1289
1229	to quantify task understanding across multiple senses.	the ocean with a spoon: Should we edit models?	1290
1230	<i>CoRR</i> .	In <i>Findings of the Association for Computational</i>	1291
		<i>Linguistics: EMNLP 2023</i> , pages 15164–15172.	1292
1231	Yasumasa Onoe, Michael JQ Zhang, Shankar Padman-		
1232	abhan, Greg Durrett, and Eunsol Choi. 2023. Can		
1233	lms learn new entities from descriptions? challenges		
1234	in propagating injected knowledge. <i>arXiv preprint</i>		
1235	<i>arXiv:2305.01651</i> .		
1236	OpenAI. 2023a. <a href="#">Chatgpt</a> .		
1237	OpenAI. 2023b. <a href="#">Gpt-4 technical report</a> .		

1293	Jirui Qi, Raquel Fernández, and Arianna Bisazza.	Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia	1347
1294	2023. Cross-lingual consistency of factual knowl-	Tsvetkov, Luke Zettlemoyer, and Scott Wen-tau	1348
1295	edge in multilingual language models. <i>arXiv preprint</i>	Yih. 2023a. Trusting your evidence: Hallucinate	1349
1296	<i>arXiv:2310.10378</i> .	less with context-aware decoding. <i>arXiv preprint</i>	1350
		<i>arXiv:2305.14739</i> .	1351
1297	Cheng Qian, Xinran Zhao, and Sherry Tongshuang Wu.	Weijia Shi, Sewon Min, Maria Lomeli, Chunting Zhou,	1352
1298	2023. "merge conflicts!" exploring the impacts of	Margaret Li, Victoria Lin, Noah A Smith, Luke	1353
1299	external distractors to parametric knowledge graphs.	Zettlemoyer, Scott Yih, and Mike Lewis. 2023b. In-	1354
1300	<i>arXiv preprint arXiv:2309.08594</i> .	context pretraining: Language modeling beyond doc-	1355
		ument boundaries. <i>arXiv preprint arXiv:2310.10638</i> .	1356
1301	Ella Rabinovich, Samuel Ackerman, Orna Raz, Eitan	Weijia Shi, Sewon Min, Michihiro Yasunaga, Min-	1357
1302	Farchi, and Ateret Anaby-Tavor. 2023. Predicting	joon Seo, Rich James, Mike Lewis, Luke Zettle-	1358
1303	question-answering performance of large language	moyer, and Wen-tau Yih. 2023c. Replug: Retrieval-	1359
1304	models through semantic consistency. <i>arXiv preprint</i>	augmented black-box language models. <i>arXiv</i>	1360
1305	<i>arXiv:2311.01152</i> .	<i>preprint arXiv:2301.12652</i> .	1361
1306	Harsh Raj, Vipul Gupta, Domenic Rosati, and Sub-	Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and	1362
1307	habrata Majumdar. 2023. Semantic consistency for	Huan Liu. 2017. Fake news detection on social me-	1363
1308	assuring reliability of large language models. <i>arXiv</i>	dia: A data mining perspective. <i>ACM SIGKDD ex-</i>	1364
1309	<i>preprint arXiv:2308.09138</i> .	<i>plorations newsletter</i> , 19(1):22–36.	1365
1310	Harsh Raj, Domenic Rosati, and Subhabrata Majum-	Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela,	1366
1311	dar. 2022. Measuring reliability of large language	and Jason Weston. 2021. Retrieval augmentation	1367
1312	models through semantic consistency. <i>arXiv preprint</i>	reduces hallucination in conversation. In <i>Findings</i>	1368
1313	<i>arXiv:2211.05853</i> .	<i>of the Association for Computational Linguistics:</i>	1369
		<i>EMNLP 2021</i> , pages 3784–3803.	1370
1314	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018.	Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mah-	1371
1315	Know what you don’t know: Unanswerable ques-	davi, Jason Wei, Hyung Won Chung, Nathan Scales,	1372
1316	tions for squad. In <i>Proceedings of the 56th Annual</i>	Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl,	1373
1317	<i>Meeting of the Association for Computational Lin-</i>	et al. 2022. Large language models encode clinical	1374
1318	<i>guistics (Volume 2: Short Papers)</i> . Association for	knowledge. <i>arXiv preprint arXiv:2212.13138</i> .	1375
1319	Computational Linguistics.	Anton Sinitin, Vsevolod Plokhhotnyuk, Dmitry Pyrkin,	1376
1320	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and	Sergei Popov, and Artem Babenko. 2019. Editable	1377
1321	Percy Liang. 2016. Squad: 100,000+ questions	neural networks. In <i>International Conference on</i>	1378
1322	for machine comprehension of text. <i>arXiv preprint</i>	<i>Learning Representations</i> .	1379
1323	<i>arXiv:1606.05250</i> .	Craig S. Smith. 2023. <a href="#">What large models cost you –</a>	1380
		<a href="#">there is no free ai lunch</a> .	1381
1324	Adam Roberts, Colin Raffel, and Noam Shazeer. 2020.	Irene Solaiman, Zeerak Talat, William Agnew, Lama	1382
1325	How much knowledge can you pack into the param-	Ahmad, Dylan Baker, Su Lin Blodgett, Hal	1383
1326	eters of a language model? In <i>Proceedings of the</i>	Daumé III, Jesse Dodge, Ellie Evans, Sara Hooker,	1384
1327	<i>2020 Conference on Empirical Methods in Natural</i>	et al. 2023. Evaluating the social impact of genera-	1385
1328	<i>Language Processing (EMNLP)</i> , pages 5418–5426.	tive ai systems in systems and society. <i>arXiv preprint</i>	1386
		<i>arXiv:2306.05949</i> .	1387
1329	Anna Rogers, Olga Kovaleva, and Anna Rumshisky.	Giovanni Spitalé, Nikola Biller-Andorno, and Federico	1388
1330	2021. A primer in bertology: What we know about	Germani. 2023. Ai model gpt-3 (dis) informs us bet-	1389
1331	how bert works. <i>Transactions of the Association for</i>	ter than humans. <i>arXiv preprint arXiv:2301.11924</i> .	1390
1332	<i>Computational Linguistics</i> , 8:842–866.	Hexiang Tan, Fei Sun, Wanli Yang, Yuanzhuo Wang,	1391
1333	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta	Qi Cao, and Xueqi Cheng. 2024. Blinded by gen-	1392
1334	Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola	erated contexts: How language models merge gen-	1393
1335	Cancedda, and Thomas Scialom. 2023. Toolformer:	erated and retrieved contexts for open-domain qa?	1394
1336	Language models can teach themselves to use tools.	<i>arXiv preprint arXiv:2401.11911</i> .	1395
1337	<i>arXiv preprint arXiv:2302.04761</i> .	Ruixiang Tang, Yu-Neng Chuang, and Xia Hu. 2023.	1396
1338	Or Sharir, Barak Peleg, and Yoav Shoham. 2020. The	The science of detecting llm-generated texts. <i>arXiv</i>	1397
1339	cost of training nlp models: A concise overview.	<i>preprint arXiv:2303.07205</i> .	1398
1340	<i>arXiv preprint arXiv:2004.08900</i> .	Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. Bert	1399
1341	Mrinank Sharma, Meg Tong, Tomasz Korbak, David	rediscovered the classical nlp pipeline. <i>arXiv preprint</i>	1400
1342	Duvenaud, Amanda Askell, Samuel R Bowman,	<i>arXiv:1905.05950</i> .	1401
1343	Newton Cheng, Esin Durmus, Zac Hatfield-Dodds,		
1344	Scott R Johnston, et al. 2023. Towards understand-		
1345	ing sycophancy in language models. <i>arXiv preprint</i>		
1346	<i>arXiv:2310.13548</i> .		

1402	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	learning in artificial intelligence. <i>Nature Machine</i>	1458
1403	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	<i>Intelligence</i> , pages 1–13.	1459
1404	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti		
1405	Bhosale, et al. 2023. Llama 2: Open founda-	Yike Wang, Shangbin Feng, Heng Wang, Weijia	1460
1406	tion and fine-tuned chat models. <i>arXiv preprint</i>	Shi, Vidhisha Balachandran, Tianxing He, and Yu-	1461
1407	<i>arXiv:2307.09288</i> .	lia Tsvetkov. 2023g. Resolving knowledge con-	1462
		licts in large language models. <i>arXiv preprint</i>	1463
1408	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot,	<i>arXiv:2310.00935</i> .	1464
1409	and Ashish Sabharwal. 2022. Musique: Multi-		
1410	hop questions via single-hop question composition.	Jerry Wei, Da Huang, Yifeng Lu, Denny Zhou, and	1465
1411	<i>Transactions of the Association for Computational</i>	Quoc V Le. 2023. Simple synthetic data reduces	1466
1412	<i>Linguistics</i> , 10:539–554.	sycophancy in large language models. <i>arXiv preprint</i>	1467
		<i>arXiv:2308.03958</i> .	1468
1413	Miles Turpin, Julian Michael, Ethan Perez, and		
1414	Samuel R Bowman. 2023. Language models don’t	Laura Weidinger, John Mellor, Maribeth Rauh, Conor	1469
1415	always say what they think: Unfaithful explana-	Griffin, Jonathan Uesato, Po-Sen Huang, Myra	1470
1416	tions in chain-of-thought prompting. <i>arXiv preprint</i>	Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh,	1471
1417	<i>arXiv:2305.04388</i> .	et al. 2021. Ethical and social risks of harm from	1472
		language models. <i>arXiv preprint arXiv:2112.04359</i> .	1473
1418	Tyler Vergho, Jean-Francois Godbout, Reihaneh Rab-		
1419	bany, and Kellin Pelrine. 2024. Comparing gpt-4	Laura Weidinger, Maribeth Rauh, Nahema Marchal, Ar-	1474
1420	and open-source language models in misinformation	rianna Manzini, Lisa Anne Hendricks, Juan Mateos-	1475
1421	mitigation. <i>arXiv preprint arXiv:2401.06920</i> .	Garcia, Stevie Bergman, Jackie Kay, Conor Grif-	1476
		fin, Ben Bariach, et al. 2023. Sociotechnical safety	1477
1422	Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry	evaluation of generative ai systems. <i>arXiv preprint</i>	1478
1423	Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny	<i>arXiv:2310.11986</i> .	1479
1424	Zhou, Quoc Le, et al. 2023. Freshllms: Refreshing		
1425	large language models with search engine augmenta-	Orion Weller, Aleem Khan, Nathaniel Weir, Dawn	1480
1426	tion. <i>arXiv preprint arXiv:2310.03214</i> .	Lawrie, and Benjamin Van Durme. 2022. Defending	1481
		against misinformation attacks in open-domain ques-	1482
1427	Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru	tion answering. <i>arXiv preprint arXiv:2212.10002</i> .	1483
1428	Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao,		
1429	Wenyang Gao, Xuming Hu, Zehan Qi, Yidong Wang,	Adina Williams, Nikita Nangia, and Samuel R Bow-	1484
1430	Linyi Yang, Jindong Wang, Xing Xie, Zheng Zhang,	man. 2017. A broad-coverage challenge corpus for	1485
1431	and Yue Zhang. 2023a. <a href="#">Survey on factuality in large</a>	sentence understanding through inference. <i>arXiv</i>	1486
1432	<a href="#">language models: Knowledge, retrieval and domain-</a>	<i>preprint arXiv:1704.05426</i> .	1487
1433	<a href="#">specificity</a> .		
1434	Cunxiang Wang, Zhikun Xu, Qipeng Guo, Xiangkun	Xiangcheng Wu, Xi Niu, and Ruhani Rahman. 2022.	1488
1435	Hu, Xuefeng Bai, Zheng Zhang, and Yue Zhang.	Topological analysis of contradictions in text. In	1489
1436	2023b. <a href="#">Exploiting Abstract Meaning Representation</a>	<i>Proceedings of the 45th International ACM SIGIR</i>	1490
1437	<a href="#">for open-domain question answering</a> . In <i>Findings of the</i>	<i>Conference on Research and Development in Informa-</i>	1491
1438	<i>Association for Computational Linguistics: ACL</i>	<i>tion Retrieval</i> , pages 2478–2483.	1492
1439	2023, pages 2083–2096, Toronto, Canada. Associa-		
1440	tion for Computational Linguistics.	Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and	1493
		Yu Su. 2023. Adaptive chameleon or stubborn	1494
1441	Cunxiang Wang, Haofei Yu, and Yue Zhang. 2023c.	sloth: Unraveling the behavior of large language	1495
1442	<a href="#">RFID: Towards rational fusion-in-decoder for open-</a>	models in knowledge conflicts. <i>arXiv preprint</i>	1496
1443	<a href="#">domain question answering</a> . In <i>Findings of the As-</i>	<i>arXiv:2305.13300</i> .	1497
1444	<i>sociation for Computational Linguistics: ACL 2023</i> ,		
1445	pages 2473–2481, Toronto, Canada. Association for	Nan Xu, Fei Wang, Bangzheng Li, Mingtao Dong, and	1498
1446	Computational Linguistics.	Muhao Chen. 2022. Does your model classify en-	1499
		tities reasonably? diagnosing and mitigating spu-	1500
1447	Fei Wang, Wenjie Mo, Yiwei Wang, Wenxuan Zhou,	rious correlations in entity typing. <i>arXiv preprint</i>	1501
1448	and Muhao Chen. 2023d. A causal view of enti-	<i>arXiv:2205.12640</i> .	1502
1449	ty bias in (large) language models. <i>arXiv preprint</i>		
1450	<i>arXiv:2305.14695</i> .	Rongwu Xu, Brian S Lin, Shujian Yang, Tianqi Zhang,	1503
		Weiyan Shi, Tianwei Zhang, Zhixuan Fang, Wei	1504
1451	Jiaan Wang, Yunlong Liang, Zengkui Sun, Yuxuan Cao,	Xu, and Han Qiu. 2023. The earth is flat be-	1505
1452	and Jiarong Xu. 2023e. Cross-lingual knowledge	cause...: Investigating llms’ belief towards misinfor-	1506
1453	editing in large language models. <i>arXiv preprint</i>	mation via persuasive conversation. <i>arXiv preprint</i>	1507
1454	<i>arXiv:2309.08952</i> .	<i>arXiv:2312.09085</i> .	1508
1455	Liyuan Wang, Xingxing Zhang, Qian Li, Mingtian	Boyang Xue, Weichao Wang, Hongru Wang, Fei Mi,	1509
1456	Zhang, Hang Su, Jun Zhu, and Yi Zhong. 2023f. In-	Rui Wang, Yasheng Wang, Lifeng Shang, Xin Jiang,	1510
1457	corporating neuro-inspired adaptability for continual	Qun Liu, and Kam-Fai Wong. 2023. Improving fac-	1511
		tual consistency for knowledge-grounded dialogue	1512



1513	systems via knowledge enhancement and alignment.	large language models: A survey. <i>ACM Transactions</i>	1568
1514	In <i>Findings of the Association for Computational</i>	<i>on Intelligent Systems and Technology</i> .	1569
1515	<i>Linguistics: EMNLP 2023</i> , pages 7829–7844.		
1516	Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng,	Yukun Zhao, Lingyong Yan, Weiwei Sun, Guoliang	1570
1517	Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu	Xing, Chong Meng, Shuaiqiang Wang, Zhicong	1571
1518	Zhang. 2023. Editing large language models: Prob-	Cheng, Zhaochun Ren, and Dawei Yin. 2023b.	1572
1519	lems, methods, and opportunities. <i>arXiv preprint</i>	Knowing what llms do not know: A simple yet	1573
1520	<i>arXiv:2305.13172</i> .	effective self-detection method. <i>arXiv preprint</i>	1574
		<i>arXiv:2310.17918</i> .	1575
1521	Jingwei Yi, Yueqi Xie, Bin Zhu, Keegan Hines, Emre	Chujie Zheng, Jinfeng Zhou, Yinhe Zheng, Libiao Peng,	1576
1522	Kiciman, Guangzhong Sun, Xing Xie, and Fangzhao	Zhen Guo, Wenquan Wu, Zhengyu Niu, Hua Wu,	1577
1523	Wu. 2023. Benchmarking and defending against indi-	and Minlie Huang. 2022. Cdconv: A benchmark	1578
1524	rect prompt injection attacks on large language mod-	for contradiction detection in chinese conversations.	1579
1525	els. <i>arXiv preprint arXiv:2312.14197</i> .	<i>arXiv preprint arXiv:2210.08511</i> .	1580
1526	Jiahao Ying, Yixin Cao, Kai Xiong, Yidong He, Long	Zexuan Zhong, Zhengxuan Wu, Christopher D Man-	1581
1527	Cui, and Yongbin Liu. 2023. Intuitive or dependent?	ning, Christopher Potts, and Danqi Chen. 2023.	1582
1528	investigating llms’ robustness to conflicting prompts.	Mquake: Assessing knowledge editing in language	1583
1529	<i>arXiv preprint arXiv:2309.17415</i> .	models via multi-hop questions. <i>arXiv preprint</i>	1584
		<i>arXiv:2305.14795</i> .	1585
1530	Wenhao Yu, Dan Iter, Shuohang Wang, Yichong	Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao	1586
1531	Xu, Mingxuan Ju, Soumya Sanyal, Chenguang	Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu,	1587
1532	Zhu, Michael Zeng, and Meng Jiang. 2022. Gen-	Lili Yu, et al. 2023a. Lima: Less is more for align-	1588
1533	erate rather than retrieve: Large language mod-	ment. <i>arXiv preprint arXiv:2305.11206</i> .	1589
1534	els are strong context generators. <i>arXiv preprint</i>		
1535	<i>arXiv:2209.10063</i> .	Jiawei Zhou, Yixuan Zhang, Qianni Luo, Andrea G	1590
1536	Jiaxin Zhang, Zhuohang Li, Kamalika Das, Bradley A	Parker, and Munmun De Choudhury. 2023b. Syn-	1591
1537	Malin, and Sricharan Kumar. 2023a. Sac <sup>3</sup> : Reliable	thetic lies: Understanding ai-generated misinformation	1592
1538	hallucination detection in black-box language models	and evaluating algorithmic and human solutions.	1593
1539	via semantic-aware cross-check consistency. <i>arXiv</i>	In <i>Proceedings of the 2023 CHI Conference on Hu-</i>	1594
1540	<i>preprint arXiv:2311.01740</i> .	<i>man Factors in Computing Systems</i> , pages 1–20.	1595
1541	Michael JQ Zhang and Eunsol Choi. 2021. Situatedqa:	Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and	1596
1542	Incorporating extra-linguistic contexts into qa. <i>arXiv</i>	Muhao Chen. 2023c. Context-faithful prompt-	1597
1543	<i>preprint arXiv:2109.06157</i> .	ing for large language models. <i>arXiv preprint</i>	1598
1544	Michael JQ Zhang and Eunsol Choi. 2023. Mitigating	<i>arXiv:2303.11315</i> .	1599
1545	temporal misalignment by discarding outdated facts.	Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and	1600
1546	<i>arXiv preprint arXiv:2305.14824</i> .	Chao Zhang. 2023. Toolqa: A dataset for llm ques-	1601
1547	Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen,	tion answering with external tools. <i>arXiv preprint</i>	1602
1548	Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing	<i>arXiv:2306.13304</i> .	1603
1549	Liu, and William B Dolan. 2020. Dialogpt: Large-	Arkaitz Zubiaga, Ahmet Aker, Kalina Bontcheva, Maria	1604
1550	scale generative pre-training for conversational re-	Liakata, and Rob Procter. 2018. Detection and res-	1605
1551	sponse generation. In <i>Proceedings of the 58th An-</i>	olution of rumours in social media: A survey. <i>ACM</i>	1606
1552	<i>annual Meeting of the Association for Computational</i>	<i>Computing Surveys (CSUR)</i> , 51(2):1–36.	1607
1553	<i>Linguistics: System Demonstrations</i> , pages 270–278.		
1554	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu,	<b>A Intra-Memory Conflict</b>	1608
1555	Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang,	Recently, LLMs have gained widespread utilization	1609
1556	Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei	across various domains, especially in knowledge-	1610
1557	Bi, Freda Shi, and Shuming Shi. 2023b. <a href="#">Siren’s song</a>	intensive question-answering systems (Gao et al.,	1611
1558	<a href="#">in the ai ocean: A survey on hallucination in large</a>	2023b; Yu et al., 2022; Petroni et al., 2019; Chen	1612
1559	<a href="#">language models</a> .	et al., 2023c). The deployment of LLMs securely	1613
1560	Yunxiang Zhang, Muhammad Khalifa, Lajanugen	and dependably hinges upon ensuring the consis-	1614
1561	Logeswaran, Moontae Lee, Honglak Lee, and	tency of their outputs when presented with expres-	1615
1562	Lu Wang. 2023c. Merging generated and retrieved	sions conveying similar meanings or intents. How-	1616
1563	knowledge for open-domain qa. <i>arXiv preprint</i>	ever, a significant challenge emerges in the form of	1617
1564	<i>arXiv:2310.14393</i> .	intra-memory conflict within LLMs. Intra-memory	1618
1565	Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu,	conflict pertains to the phenomenon wherein the	1619
1566	Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei		
1567	Yin, and Mengnan Du. 2023a. Explainability for		

Datasets	Approach <sup>1</sup>	Base <sup>2</sup>	Size	Conflict
Xie et al. (2023)	Gen	PopQA (2023), STRATEGYQA ((Geva et al., 2021))	20,091	CM <sup>3</sup>
KC (2023g)	Sub	N/A (LLM generated)	9,803	CM
KRE (2023)	Gen	MuSiQue (2022), SQuAD2.0 (2018), ECQA (2021), e-CARE (2022a)	11,684	CM
Farm (2023)	Gen	BoolQ (2019), NQ (2019), TruthfulQA (2022)	1,952	CM
Tan et al. (2024)	Gen	NQ (2019), TriviaQA (2017)	14,923	CM
Pan et al. (2023a)	Gen,Sub	SQuAD 1.1 (2016)	52,189	IC
CONTRADOC (2023a)	Gen	CNN-DailyMail (2015), NarrativeQA (2018) WikiText (2016)	449	IC
ClaimDiff (2022)	Hum	N/A	2,941	IC
WikiContradiction (2021)	Hum	Wikipedia	2,210	IC
PARAREL (2021)	Hum	T-REx (2018)	328	IM

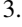
1. Approach refers to how the conflicts are crafted, including entity-level substitution (Sub), generative approaches employing an LLM (Gen), human annotation (Hum).
2. Base refers to the base dataset(s) that serve as the foundation for generating conflicts, if applicable.
3.  For **CM** datasets, conflicts are derived from a *certain* model’s parametric knowledge, which can vary between models. Therefore, application requires selecting a subset of the dataset that aligns with the tested model’s knowledge.

Table 1: Datasets on evaluating a large language model’s behavior when encounter knowledge conflicts.

model’s parameters encompass multiple, potentially conflicting versions of knowledge. As a result, the model may exhibit unpredictable behavior and different outputs for inputs that, while syntactically varied, convey the same semantic information (Chang and Bergen, 2023; Chen et al., 2023a; Raj et al., 2023; Rabinovich et al., 2023; Raj et al., 2022; Bartsch et al., 2023). Intra-memory conflict substantially diminishes the practicality and efficacy of LLMs.

### A.1 Causes

Intra-memory conflicts within LLMs can be attributed to three primary factors: training corpus bias (Wang et al., 2023d; Xu et al., 2022), decoding strategies Lee et al. (2022b); Huang et al. (2023), and knowledge editing (Yao et al., 2023; Li et al., 2023d). These factors respectively pertain to the training phase, the inference phase, and subsequent knowledge refinement.

**Bias in Training Corpora.** Recent research demonstrates that the primary phase for knowledge acquisition in LLMs predominantly occurs in the pre-training stage (Zhou et al., 2023a; Kaddour et al., 2023; Naveed et al., 2023; Akyürek et al., 2022; Singhal et al., 2022). Pre-training data is primarily crawled from the internet, which exhibits a diverse range of data quality, potentially including inaccurate or misleading information (Bender et al., 2021; Weidinger et al., 2021). When LLMs are trained on data containing incorrect knowledge, they may memorize and inadvertently amplify these inaccuracies (Lin et al., 2022; Elazar et al., 2022; Lam et al., 2022; Grosse et al., 2023), leading to a situation where conflicting knowledge

coexists within the parameters of LLMs concurrently.

Moreover, prior works indicate that LLMs may possess a propensity for encoding superficial associations prevalent within their training data, as opposed to genuinely comprehending the underlying knowledge contained therein (Li et al., 2022b; Kang and Choi, 2023; Zhao et al., 2023a; Kandpal et al., 2023). Consequently, when the training dataset exhibits a bias towards spurious correlations, it can result in the LLMs displaying a propensity to generate predetermined responses rooted in those spurious correlations. Due to the dependency of spurious correlations, LLMs may provide divergent answers when presented with prompts exhibiting distinct syntactic structures but conveying equivalent semantic meaning, thereby leading to instances of intra-memory conflicts.

**Decoding Strategy.** The direct output of LLMs is a probability distribution representing the possible next tokens. Sampling is a crucial step in determining the generated content from this distribution. Currently, there are various proposed sampling techniques, including greedy sampling, top-p sampling, top-k sampling, and others (Jawahar et al., 2020; Massarelli et al., 2019). These techniques can be categorized into two main groups: deterministic sampling and stochastic sampling. Stochastic sampling stands as the prevailing decoding strategy employed by LLMs (Fan et al., 2018; Holtzman et al., 2019). However, the stochastic nature of sampling introduces uncertainty into the generated content. Furthermore, due to the intrinsic left-to-right generation pattern of LLMs, the selection of the sampling token can wield a significant influ-

ence over the content of subsequent generations. The use of stochastic sampling may lead LLMs to produce entirely different content, even when provided with the same context, causing intra-memory conflict (Lee et al., 2022b; Huang et al., 2023; Dziri et al., 2021).

**Knowledge Editing.** With the dramatic increase of model parameters, fine-tuning LLMs become increasingly challenging and resource-intensive. In response to this challenge, researchers have turned to knowledge editing techniques as a means to efficiently modify the small scope of knowledge learned of LLMs (Meng et al., 2022; Ilharco et al., 2022; Zhong et al., 2023). Ensuring the consistency of modifications poses a significant challenge. Due to the potential defects inherent in the editing method, the modified knowledge cannot be generalized effectively. Consequently, the model exhibits variations in its responses across different contexts. It may adapt its knowledge to specific situations while maintaining consistency in others (Li et al., 2023d; Yao et al., 2023). Intra-memory conflict is primarily considered a side effect in the context of knowledge editing.

## A.2 Analysis

**Self-Inconsistency.** Elazar et al. (2021) developed a method for assessing the knowledge consistency of a model, focusing specifically on knowledge triples. The authors primarily conducted experiments using BERT, RoBERTa, and ALBERT, and their findings indicate that these models exhibit poor consistency, with guaranteed accuracy ranging from only 50% to 60% on the test data. Hase et al. (2023) employed the same indicators of Elazar et al. (2021), but they utilized a more diverse dataset. Their study also revealed that the consistency of RoBERTa-base and BART-base within the paraphrase context was lacking. Zhao et al. (2023b) first reformulated the questions and then assessed the consistency of the LLM’s responses to these reformulated questions. The findings of their research revealed that even GPT-4 exhibits a notable inconsistency rate of 13% when applied to Commonsense Question-Answering tasks. They further found that LLMs are more likely to produce inconsistencies in the face of uncommon knowledge. Dong et al. (2023) conducted experiments on 20 open-source LLMs and found that all of these models exhibit strong inconsistencies. The authors also found that fine-tuning LLMs on data collected from

a more knowledgeable model could augment its knowledge. Li et al. (2023c) explored an additional aspect of inconsistency that LLMs can give an initial answer to a question, but it may subsequently contradict that answer when asked if it’s correct. The authors conducted experiments focusing on Close-Book Question Answering and revealed that GPT-4 is consistent on 95% of cases, while Alpaca-30B only displays consistency in 50% of cases.

To further analyze the inconsistency exhibited by LLMs, a study conducted by Li et al. (2022b) revealed that encoder-based models tend to generate missing factual words more relying on positionally close and highly co-occurring words, rather than knowledge-dependent words. This phenomenon arises due to these models’ tendency to overlearn inappropriate associations from the training dataset. Kang and Choi (2023) demonstrated that LLMs are prone to co-occurrence bias, where they favor frequently co-occurring words over the correct answer. Furthermore, their research highlighted that LLMs face challenges in recalling facts in cases where the subject and object rarely appear together in the pre-training dataset, even though these facts are encountered during fine-tuning.

**Latent Representation of Knowledge.** Contemporary large language models all employ a multi-layer transformer structure, giving rise to a unique form of inter-memory conflict, *i.e.* the presence of distinct knowledge representations across different layers of a model. In the past, numerous researchers have proposed that a language model would store low-level information at a shallow level, and semantic information at a high level (Tenney et al., 2019; Rogers et al., 2021; Jawahar et al., 2019; Cui et al., 2020). Chuang et al. (2023) explored this aspect within the context of LLMs and discovered that factual knowledge in LLMs is typically concentrated within specific transformer layers and different layers of inconsistent knowledge. Moreover, Li et al. (2023b) discovered that the correct knowledge is indeed stored within the parameters of the large model, but it may not be accurately expressed during the generation process. The authors conducted two experiments on the same LLaMa 7B, one focused on the generation accuracy, and the other utilizing a knowledge probe to examine the knowledge containment. The results of these experiments revealed a substantial 40% disparity between the knowledge probe accuracy and the generation accuracy.



**Cross-lingual Inconsistency.** The meaning of true knowledge is not influenced by surface form (Ohmer et al., 2023). Knowledge held by LLMs should also possess this characteristic. However, unlike humans, LLMs maintain distinct sets of knowledge for various languages (Ji et al., 2023), leading to potential inconsistencies in their knowledge across different languages. Wang et al. (2023e) analyzed LLMs’ knowledge expressed in one language after implementing knowledge editing on this knowledge expressed in another language. Their findings suggest that LLMs face difficulties when attempting to extend the edited knowledge to other languages and exhibit inconsistent behaviors. These findings indicate that knowledge related to different languages is stored separately within the model parameters, introducing a risk of intra-memory conflict of the model. Qi et al. (2023) conducted a more direct study. They propose a metric named RankC for evaluating the cross-lingual consistency of factual knowledge of LLMs. They employed this metric for analyzing multiple models and revealed that the knowledge learned by LLMs is not language-agnostic. Instead, a strong language dependence exists in the knowledge of LLMs. Additionally, they found that increasing the model size does not lead to an improvement in cross-lingual consistency.

### A.3 Mitigation

#### A.3.1 Improving Consistency

*Fine-tuning.* Elazar et al. (2021) introduces a new benchmark called PARALLEL, which consists of 328 paraphrases that describe 38 binary relationships. The author proposed the consistency loss function and leveraged both T-REx (Elsahar et al., 2018) and PARAREL to train BERT with the consistency loss and standard MLM loss. The fine-tuned BERT yielded impressive results when evaluated on the corresponding test dataset. However, it is worth noting that when applying the same fine-tuning approach to the SQuAD dataset (Gan and Ng, 2019), no significant impact was observed. Li et al. (2023c) initially employ the model as both a generator and a validator and queries the generator to acquire a response, following which the validator is consulted to assess the accuracy of the generated response. The paired responses from both the generator and the validator are filtered, retaining only those pairs that exhibit consistency, which is used to fine-tune the same model to enhance

the likelihood of consistent pairs. They use their method to finetune Alpaca-30B with TriviaQA and demonstrate significant consistency enhancement in TriviaQA and natural questions.

*Plug-in.* Jang and Lukasiewicz (2023) first employ the intermediate training technique to retrain PLMs with word definition pairs in a dictionary to enhance the models’ understanding of symbol meanings. They then introduce a training-efficient parameter integration method that merges the acquired parameters with those of other existing PLMs. They conduct experiments using RoBERTa as the backbone model and evaluate the effectiveness of their method in BECEL (Jang et al., 2022b). *Output Ensemble.* Mitchell et al. (2022) propose a method to mitigate the inconsistency of LMs. They use a base model to create possible answers and a relation model to assess the logical relationships between these answers. The final answer is selected by considering both the base model’s and the relation model’s beliefs. The effectiveness of this approach on the language model is primarily demonstrated through experiments conducted on the BeliefBank QA dataset (Kassner et al., 2021). Instead, Zhao et al. (2023b) introduces a method to detect whether a question may cause inconsistency for LLMs. Specifically, they first use LLMs to rephrase the original question and obtain corresponding answers. They then cluster these answers and examine the divergence. The detection is determined based on the divergence level. They conduct comprehensive experiments using open-source and closed-source LLMs to validate the effectiveness of their proposed method on datasets including FaVI (Park et al., 2021) and ComQA (Abujabal et al., 2018).

#### A.3.2 Improving Factuality

Chuang et al. (2023) proposed a novel contrastive decoding approach named Dola. Specifically, the author initially developed a dynamic layer selection strategy, choosing the appropriate premature layers and mature layers. The next word’s output probability is then determined by computing the difference in log probabilities of the premature layers and the mature layers. The author primarily conducted experiments on LLaMa in Truthful QA, Strategy QA, GSM8K, and FACTOR. The experiments provided substantial evidence that Dola consistently enhances the truthfulness of models and mitigates inconsistency issues. Li et al. (2023b) proposed a similar method named ITI. ITI first identifies a



sparse set of attention heads that exhibit high linear probing accuracy for truthfulness, as measured by the TruthfulQA (Lin et al., 2022). During the inference phase, ITI shifts activations along the truth-correlated direction, obtained through knowledge probing. This intervention is repeated autoregressively for every token. The authors also verified their method using LLaMA-7B on TruthfulQA and found that ITI achieves a significant improvement in the factual knowledge accuracy of LLMs.