

Knowledge Conflicts for LLMs: A Survey

Anonymous ARR submission

Abstract

This survey provides an in-depth analysis of knowledge conflicts for large language models (LLMs), highlighting the complex challenges they encounter when blending contextual and parametric knowledge. Our focus is on three categories of knowledge conflicts: context-memory, inter-context, and intra-memory conflict. These conflicts can significantly impact the trustworthiness and performance of LLMs, especially in real-world applications where noise and misinformation are common. By categorizing these conflicts, exploring the causes, examining the behaviors of LLMs under such conflicts, and reviewing available solutions, this survey aims to shed light on strategies for improving the robustness of LLMs, thereby serving as a valuable resource for advancing research in this evolving area.

1 Introduction

Large language models (LLMs; Brown et al. 2020; Touvron et al. 2023; OpenAI 2024) are renowned for encapsulating a vast repository of world knowledge (Petroni et al., 2019; Roberts et al., 2020), referred to as *parametric knowledge*. These models excel in various knowledge-intensive tasks. Meanwhile, LLMs continue to engage with external *contextual knowledge* after deployed (Pan et al., 2022), including user prompts (Liu et al., 2023a), documents from the Web (Shi et al., 2023c), or tools (Schick et al., 2023; Zhuang et al., 2023).

Integrating contextual knowledge into LLMs enables them to keep abreast of current events (Kasai et al., 2022) and generate more accurate responses (Shuster et al., 2021), yet it risks conflicting due to the rich knowledge sources. The discrepancies among the contexts and the model’s parametric knowledge are referred to as *knowledge conflicts* (Chen et al., 2022; Xie et al., 2023). In this paper, we categorize **three** distinct types of knowledge conflicts, as shown in Figure 1. Contextual

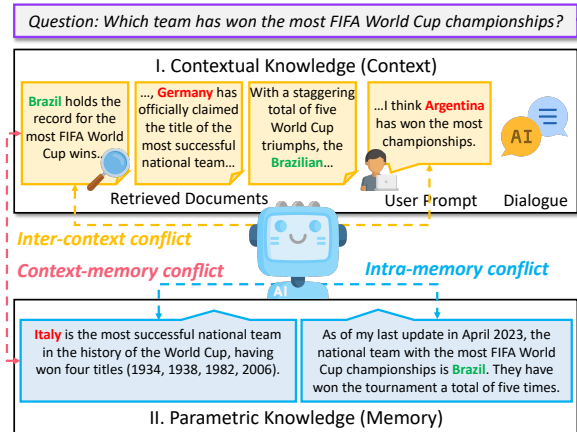


Figure 1: An LLM may encounter three types of knowledge conflicts, stemming from knowledge sources—either contextual (in yellow) or inherent to the LLM’s parameters (in blue). When confronted with a user’s question (in purple) entailing knowledge of complex conflicts, the LLM is required to resolve these discrepancies to deliver accurate responses.

knowledge (*context*, including user prompts, dialogue history, and retrieved documents) can conflict with the parametric knowledge (*memory*), where we term it as **context-memory conflict**. In the meantime, the context might be fraught with noise (Zhang and Choi, 2021) or even deliberately crafted misinformation (Du et al., 2022b). The conflict among contextual knowledge is dubbed as **inter-context conflict**. To reduce uncertainties in responses, the user may pose the question in various forms, resulting in the LLM’s parametric knowledge in divergent responses. This variance may stem from the inconsistencies present in the pre-training data (Huang et al., 2023), which gives rise to what we call **intra-memory conflict**.

Knowledge conflicts attract attention with the advent of LLMs. 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 the context is factually wrong (Pan et al., 2023b). Given the im-

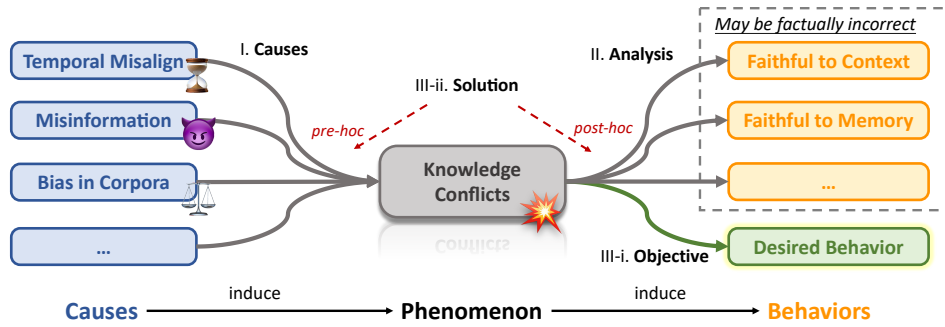


Figure 2: We view knowledge conflict not only as a standalone **phenomenon** but also as a nexus that connects various causal triggers (**causes**) with the **behaviors** of LLMs.

062 applications for the trustworthiness (Du et al., 2022b),
 063 real-time accuracy (Kasai et al., 2022), and robust-
 064 ness (Ying et al., 2023) of LLMs, it is imperative
 065 to delve deeper into understanding such conflicts
 066 (Xie et al., 2023; Wang et al., 2023e). Existing
 067 reviews (Zhang et al., 2023d; Wang et al., 2023a;
 068 Feng et al., 2023) either touch upon knowledge
 069 conflicts as a subtopic within a broader context and
 070 primarily focus on specific scenarios (Feng et al.,
 071 2023). To fill the gap, we aim to provide a compre-
 072 hensive survey encompassing the categorization,
 073 cause and behavior analysis, and solutions for ad-
 074 dressing various knowledge conflicts.

075 We conceptualize the *lifecycle of knowledge con-*
 076 *licts* as both a *cause* leading to various behaviors,
 077 and an *effect* emerges from the intricate nature of
 078 knowledge as in Figure 2. Our research under-
 079 scores the significance of understanding the orig-
 080 ins of these conflicts. Although existing analyses
 081 (Chen et al., 2022; Xie et al., 2023; Wang et al.,
 082 2023e) tend to construct such conflicts artificially,
 083 we posit that these analyses do not sufficiently ad-
 084 dress the interconnectedness of the issue. Going
 085 beyond, we provide a systematic review of mitiga-
 086 tion strategies, which are employed to minimize the
 087 undesirable consequences of knowledge conflicts.
 088 Based on the timing relative to potential conflicts,
 089 such strategies are divided into *pre-hoc* and *post-*
 090 *hoc* strategies. The key distinction between them
 091 lies in whether adjustments are made *before* or *after*
 092 potential conflicts arise. We discuss three kinds of
 093 knowledge conflicts, detailing the causes, analysis
 094 of model behaviors, and available solutions accord-
 095 ing to their respective objectives. The taxonomy of
 096 knowledge conflicts is outlined in Figure 3.

097 2 Context-Memory Conflict

098 LLMs are characterized by fixed parametric knowl-
 099 edge, a result of the substantial pertaining pro-
 100 cess (Sharir et al., 2020; Hoffmann et al., 2022;

101 Smith, 2023). This static parametric knowledge
 102 stands in stark contrast to the dynamic nature of
 103 external information, which evolves at a rapid
 104 pace (De Cao et al., 2021; Kasai et al., 2022).

105 2.1 Causes

106 **Temporal Misalignment.** It *naturally* arises in
 107 models trained on data collected in the past, as
 108 they may not accurately reflect contemporary reali-
 109 ties (Luu et al., 2021; Lazaridou et al., 2021; Liska
 110 et al., 2022). Such misalignment can degrade the
 111 model’s performance on various NLP tasks and
 112 relevancy over time (Luu et al., 2021; Zhang and
 113 Choi, 2021; Dhingra et al., 2022; Kasai et al., 2022;
 114 Cheang et al., 2023), as it may fail to capture new
 115 trends or shifts in language use. Furthermore, the
 116 issue of temporal misalignment is expected to in-
 117 tensify due to the pre-training paradigm and the
 118 escalating costs associated with scaling up mod-
 119 els (Chowdhery et al., 2023; OpenAI, 2024).

120 Prior works tackle temporal misalignment by
 121 focusing on three lines of strategies: *Knowledge*
 122 *editing (KE)* aims to directly update the parametric
 123 knowledge (Sinitsin et al., 2020; Mitchell et al.,
 124 2021; Onoe et al., 2023). *Retrieval-augmented gen-*
 125 *eration (RAG)* fetches relevant documents from
 126 external sources to supplement the model’s knowl-
 127 edge without altering its parameters (Karpukhin
 128 et al., 2020; Guu et al., 2020; Lewis et al., 2020;
 129 Lazaridou et al., 2022; Vu et al., 2023). *Contin-*
 130 *ual learning (CL)* updates the internal knowledge
 131 through continual training on updated data (Lazari-
 132 dou et al., 2021; Jang et al., 2021, 2022). However,
 133 KE can bring in side effects such as knowledge in-
 134 consistency and may enhance the hallucination of
 135 LLMs (Li et al., 2023f; Pinter and Elhadad, 2023).
 136 RAG is inevitable to encounter conflicts since
 137 model parameters are not updated (Chen et al.,
 138 2021; Zhang and Choi, 2021). CL suffers from the
 139 issue of catastrophic forgetting and demands sig-

nificant computational resources (De Lange et al., 2021; He et al., 2021; Wang et al., 2023d).

Misinformation Pollution. Adversaries can exploit this vulnerability by introducing misleading information into retrieved documents (Pan et al., 2023a,b; Weller et al., 2022) and user conversations (Xu et al., 2023). *Prompt injection* attack (Liu et al., 2023b; Greshake et al., 2023; Yi et al., 2023; Xu et al., 2024) is one such technique, where models may inadvertently spread misinformation if they use deceptive inputs (Pan et al., 2023b; Xu et al., 2023). Misinformation undermines the accuracy of automated fact-checking (Du et al., 2022b) and question-answering systems (Pan et al., 2023a,b). Recent studies 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). Recently, there has been growing apprehension regarding the potential generation of misinformation by LLMs (Ayoobi et al., 2023; Kidd and Birhane, 2023; Carlini et al., 2023; Zhou et al., 2023c; 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), which underscores the urgency of addressing the nuanced challenges LLMs pose within contextual misinformation.

Remarks. Temporal misalignment and misinformation pollution are two separate scenarios that give rise to context-memory conflicts. For the former, the up-to-date contextual information is considered accurate. *Conversely*, for the latter, the contextual information contains misinformation and is therefore considered incorrect.

2.2 Analysis of Model Behaviors

We summarize studies on how LLMs behave under context-memory conflicts within open-domain question answering (ODQA) and general setups.

ODQA. Early effort (Longpre et al., 2021) explores how QA models act when the provided contextual information contradicts the memory. An automated framework first 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. Longpre et al. (2021) reveal a tendency of models to over-rely on parametric knowledge. Chen et al. (2022) report differing observations, they note that models pre-

dominantly rely on contextual knowledge in their best-performing settings. This divergence can be attributed to two factors. Firstly, the entity substitution approach (Longpre et al., 2021) potentially reduces the semantic coherence of the perturbed context. Secondly, Chen et al. (2022) utilize multiple evidence rather than one (Longpre et al., 2021). Recently, Tan et al. (2024) examine how large LMs integrate context with generated memory. They observe that LLMs tend to prioritize parametric knowledge thanks to the greater similarity between generated contents and input, as well as the often incomplete nature of retrieved information.

General. LLMs exhibit a complex relationship with conflicting information. While highly receptive to convincing external evidence (Xie et al., 2023), they also demonstrate a strong confirmation bias (Nickerson, 1998), favoring information consistent with their memory. This leads to challenges in resolving such conflicts, as LLMs struggle to pinpoint conflicting segments and provide disentangled responses (Wang et al., 2023e). Research exploring LLMs’ robustness under conflicts reveals a susceptibility to misleading prompts, particularly in commonsense knowledge (Ying et al., 2023). Furthermore, LLMs often deviate from their parametric knowledge when presented with direct conflicts or contextual changes (Qian et al., 2023). Studies investigating LLMs in interactive sessions highlight a tendency to favor logically structured knowledge, even when it is factual wrong (Xu et al., 2023). These findings underscore the need for further research into the interaction between parametric and contextual knowledge for LLMs.

Remarks. Researchers analyze LLMs’ behavior under conflicting knowledge by creating artificial conflicts, initially through entity-level substitutions and later by using LLMs to generate semantically coherent conflicts. While no definitive rule exists for prioritizing contextual or parametric knowledge, LLMs tend to favor information that is semantically coherent over generic conflicting information.

2.3 Solutions

Solutions are organized according to their **objectives**, *i.e.*, the desired behaviors we expect from an LLM when it encounters conflicts. Existing strategies can be categorized into 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 and provide disentangled answers. *Improving factuality* strategies aim for an integrated response leveraging both context and parametric knowledge towards a more truthful solution.

Faithful to Context. Several approaches have been proposed to achieve this goal. Fine-tuning approaches like Knowledge Aware (Li et al., 2022a) incorporate counterfactual and irrelevant contexts into training data to enhance controllability and robustness. Similarly, TrueTeacher (Gekhman et al., 2023) focus on improving factual consistency in summarization by annotating model-generated summaries with LLMs. Prompting strategies (Zhou et al., 2023d) utilize opinion-based prompts and counterfactual demonstrations to enhance LLMs’ adherence to context without additional training. Decoding techniques like Context-aware Decoding (Shi et al., 2023a) amplify the difference in output probabilities with and without context, prioritizing relevant context over prior knowledge. Knowledge plug-in approaches, such as Continuously-updated QA (Lee et al., 2022a), use plug-and-play modules to store updated knowledge, solving knowledge conflicts without affecting the original model. Pre-training methods (Shi et al., 2023b) extend LLMs’ ability to handle long and varied contexts across multiple documents, potentially resolving knowledge conflicts by synthesizing information from broader contexts. Finally, fact validity prediction approaches (Zhang and Choi, 2023) identify and discard outdated facts in LLMs, improving performance on tasks like ODQA by ensuring adherence to up-to-date contextual information.

Discriminating Misinformation. To combat misinformation, various defense strategies have been proposed. Pan et al. (2023b) advocates for misinformation detection and vigilant prompting, aiming to improve the model’s faithfulness to factual information. Xu et al. (2023) employ a system prompt to encourage LLMs to be cautious about misinformation and verify their memorized knowledge before responding, further enhancing faithfulness. Weller et al. (2022) leverage the redundancy of information in large corpora to mitigate knowledge conflicts. Their approach involves query augmentation to retrieve diverse, less likely poisoned passages, then compares the consistency of predicted answers across retrieved contexts. This strategy en-

ures faithfulness by cross-verifying answers from multiple sources. Hong et al. (2023) fine-tune a smaller LM as a discriminator and integrate prompting techniques to enable the model to distinguish between reliable and unreliable information.

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. (2023e) introduce a method to improve 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.

Improving Factuality. Zhang et al. (2023e) 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. Jin et al. (2024a) introduces a contrastive-decoding-based algorithm to maximize the difference between various logits under knowledge conflicts and calibrates the model’s confidence in the truthful answer.

Remarks. Current mitigation approaches for knowledge conflicts are ineffective because they fail to differentiate between the two underlying causes. Blindly prioritizing either faithfulness to context or knowledge is undesirable. Researchers advocate for LLMs that empower users to make informed decisions by providing distinct answers based on both parametric and contextual information (Wang et al., 2023e; Floridi, 2023).

3 Inter-Context Conflict

Inter-context conflicts manifest in LLMs when incorporating conflicting segments among external information sources, a challenge accentuated by the advent of RAG techniques.

3.1 Causes

Misinformation. Similar to context-memory conflict, this type of conflict can also be affected by misinformation and will not be discussed repeatedly.

Outdated Information. It is also important to recognize that facts can evolve. Retrieved documents may contain updated and outdated informa-

tion from the network simultaneously, leading to conflicts between these documents (Chen et al., 2021; Liska et al., 2022; Kasai et al., 2022).

3.2 Analysis of Model Behaviors

Performance Impact. Previous research has shown that LMs can be significantly influenced by misinformation or outdated information within a specific context (Zhang and Choi, 2021; Du et al., 2022b). Pan et al. (2023a) demonstrated that LLMs are susceptible to misinformation attacks, even when the fake articles are generated by models. Chen et al. (2022) investigated how LLMs handle contradictory contexts and found that inconsistencies across knowledge sources have a minimal effect on their confidence levels. These models tend to favor context directly related to the query and context that aligns with their parametric knowledge. Xie et al. (2023) confirmed these findings, showing that LLMs exhibit a bias towards evidence that aligns with their parametric memory and a predisposition towards emphasizing information related to popular entities and answers corroborated by a larger volume of documents. Furthermore, they found that LLMs are sensitive to the order in which data is introduced. Jin et al. (2024a) discovered that LLMs struggle with reasoning as the number of conflicting hops increases.

Detection Ability. Several studies highlight the challenges faced by LMs in identifying contradictions. Zheng et al. (2022) demonstrate that LMs struggle to detect contradictory statements within Chinese conversations. Li et al. (2023a) analyze the performance of LLMs in identifying contradictory documents across various sources, including news (Hermann et al., 2015), stories (Kočíský et al., 2018), and Wikipedia (Merity et al., 2017), finding that the average detection accuracy is low. They also observe that LLMs perform poorly when dealing with contradictions involving subjective emotions or perspectives. Wan et al. (2024) investigate the text features influencing LLMs’ assessment of document credibility in the presence of conflicting information, discovering that models prioritize relevance over stylistic features. Jin et al. (2024a) further highlight the difficulty LLMs encounter in distinguishing truthful information from misinformation, showing a tendency to favor evidence that appears most frequently within the context.

Remarks. Exploring responses to contextual nuances is essential, as variations in training data lead

to differences in behavior. Despite some similarities, LLMs’ methods of identifying misinformation differ significantly from those of humans.

3.3 Solutions

Eliminating Conflict. Several approaches have been proposed to address the challenge of eliminating conflict in text. Specialized models, such as the Pairwise Contradiction Neural Network (Hsu et al., 2021), utilize fine-tuned Sentence-BERT embeddings to determine contradiction probabilities. Pielka et al. (2022) emphasize the importance of integrating linguistic knowledge into the learning process to improve contradiction detection, as models like XLM-RoBERTa struggle with syntactic and semantic features. Wu et al. (2022) propose incorporating topological text representations into language models to enhance contradiction detection, evaluating their approach on the MultiNLI dataset (Williams et al., 2018). General models, such as Chern et al. (2023)’s fact-checking framework, integrate LLMs with various tools to detect factual errors. Leite et al. (2023) leverage LLMs to generate weak labels associated with credibility signals for input text, aggregating these labels through weak supervision techniques to predict veracity.

Improving Robustness. To enhance robustness, Hong et al. (2023) propose a fine-tuning method that trains a discriminator and decoder simultaneously using a shared encoder, alongside strategies involving prompting GPT-3 to identify perturbed documents and integrating the discriminator’s output into prompts. Weller et al. (2022) explore query augmentation by prompting GPT-3 to generate new questions based on the original query, evaluating answer confidence through passage retrieval, and deciding whether to rely on the original prediction or aggregate predictions from high-confidence augmented questions. While both approaches aim for robustness, Hong et al. (2023)’s fine-tuning method demonstrates the most promising results.

Remarks. Strategies for addressing inter-context conflicts primarily rely on model knowledge or leverage external knowledge such as retrieved documents. Moreover, augmenting LLM capabilities with external tools has emerged as a novel paradigm. Exploring the use of external tools to support LLMs in resolving inter-context conflicts is a promising approach. In addition, devising a unified and efficient approach to handle various conflict types remains a formidable challenge.

4 Intra-Memory Conflict

Consistent LLM outputs for identical inputs are essential. However, intra-memory conflicts, where LLMs generate differing responses to similar inputs, undermine their reliability and utility by introducing undesirable uncertainty.

4.1 Causes

The following three factors respectively pertain to training, inference, and knowledge refinement.

Bias in Training Corpora. While LLMs primarily acquire knowledge during pre-training (Zhou et al., 2023a; Kaddour et al., 2023; Naveed et al., 2023; Akyürek et al., 2022; Singhal et al., 2022), the vast and often unreliable nature of internet-sourced training data (Bender et al., 2021; Weidinger et al., 2021) can lead to the memorization and amplification of inaccuracies (Lin et al., 2022; Elazar et al., 2022; Lam et al., 2022; Grosse et al., 2023). This results in LLMs potentially harboring conflicting knowledge within their parameters. Furthermore, LLMs tend to encode superficial associations rather than true comprehension of training data (Li et al., 2022b; Kang and Choi, 2023; Zhao et al., 2023a; Kandpal et al., 2023), leading to predetermined responses based on spurious correlations and potentially divergent answers for semantically equivalent but syntactically distinct prompts.

Decoding Strategy. LLMs generate text by sampling from a probability distribution over potential next tokens. Stochastic sampling methods like top-k and top-p sampling are commonly used for decoding, introducing randomness in the generated content (Jawahar et al., 2020; Massarelli et al., 2020; Fan et al., 2018; Holtzman et al., 2020). However, this randomness can cause intra-memory conflicts, where the model produces different outputs for the same input due to the left-to-right generation pattern and the influence of sampled tokens on subsequent generations (Lee et al., 2022b; Huang et al., 2023; Dziri et al., 2021).

Knowledge Editing. With the exponential increase of model parameters, fine-tuning LLMs become increasingly resource-intensive. In response to this, researchers explore knowledge editing techniques to efficiently modify a small scope of the knowledge in LLMs (Meng et al., 2022; Zhong et al., 2023). Ensuring the consistency of such modification poses a significant challenge. Due to the potential limitations inherent in the editing method, the modified knowledge cannot be generalized ef-

fectively. This can result in LLMs producing inconsistent responses when dealing with the same piece of knowledge in varying situations (Li et al., 2023f; Yao et al., 2023).

Remarks. Intra-memory conflicts in LLMs arise from three main causes at different stages. Training corpus bias is the primary catalyst, causing inconsistencies in the model’s knowledge. The randomness of the decoding process during inference exacerbates these inconsistencies. Additionally, knowledge editing can inadvertently introduce conflicting information.

4.2 Analysis of Model Behaviors

Self-Inconsistency. LLMs exhibit significant self-inconsistency, as evidenced by multiple studies. Elazar et al. (2021) found that BERT, RoBERTa, and ALBERT struggle with knowledge consistency, achieving accuracy rates barely exceeding 50-60%. Hase et al. (2023), using a more diverse dataset, confirmed these findings, highlighting the inconsistency of RoBERTa-base and BART-base in paragraph contexts. Zhao et al. (2023b) revealed that even GPT-4 displays a 13% inconsistency rate in Commonsense Question-Answering tasks, particularly when dealing with uncommon knowledge. Dong et al. (2023) further demonstrated that various open-source LLMs exhibit strong inconsistencies. Li et al. (2023d) identified another aspect of inconsistency, where LLMs may initially answer a question but subsequently deny the answer when asked for confirmation. Li et al. (2022b) attributed this inconsistency in encoder-based models to their reliance on positionally close and highly co-occurring words, leading to the generation of misinformation. Kang and Choi (2023) further explained this phenomenon as a co-occurrence bias, where LLMs prioritize frequently co-occurring words over correct answers, particularly when recalling facts with rarely co-occurring subject-object pairs in the pre-training dataset, even after fine-tuning.

Latent Representation of Knowledge. Contemporary LLMs, built on multi-layer transformer architectures, exhibit a complex inter-memory conflict with distinct knowledge representations scattered across layers. Research suggests that LLMs store low-level information at shallower layers and semantic information at deeper layers (Tenney et al., 2019; Rogers et al., 2020; Wang et al., 2019; Jawahar et al., 2019; Cui et al., 2020). Chuang et al. (2023) demonstrate that factual knowledge is con-

540 concentrated within specific transformer layers, lead- 590
541 ing to inconsistent knowledge across layers. Fur- 591
542 thermore, [Li et al. \(2023c\)](#) highlight a discrepancy 592
543 between knowledge storage and generation accu- 593
544 racy. Their experiments reveal a 40% gap between 594
545 the accuracy of a knowledge probe and the gener- 595
546 ation accuracy, suggesting that while the correct 596
547 knowledge is present within the parameters, it may 597
548 not be effectively expressed during generation. 598

549 **Cross-lingual Inconsistency.** While true knowl- 599
550 edge should be universally accessible regardless 600
551 of language variation ([Ohmer et al., 2023](#)), LLMs 601
552 exhibit cross-lingual inconsistencies ([Ji et al., 2023](#); 602
553 [Xue et al., 2024](#)). This inconsistency arises from 603
554 LLMs storing knowledge related to different lan- 604
555 guages separately within their parameters ([Wang 605](#)
556 [et al., 2023c](#)). [Qi et al. \(2023\)](#) propose RankC, a 606
557 metric for evaluating cross-lingual consistency of 607
558 factual knowledge, and reveals a strong language 608
559 dependence in LLMs, with no improvement in con- 609
560 sistency observed even with larger models. 610

561 **Remarks.** The phenomenon of inter-memory con- 611
562 flict in LLMs predominantly manifests through 612
563 inconsistent responses to semantically identical 613
564 queries. This inconsistency is primarily attributed 614
565 to the suboptimal quality of datasets utilized during 615
566 the pre-training phase. Addressing this challenge 616
567 necessitates the development of efficient and cost- 617
568 effective solutions, which remains a significant hur- 618
569 dle. Additionally, LLMs are characterized by the 619
570 presence of multiple knowledge circuits, which sig- 620
571 nificantly influence their response mechanisms to 621
572 specific inquiries. The exploration and detailed ex- 622
573 amination of these knowledge circuits within LLMs 623
574 represent a promising avenue for future research. 624

575 4.3 Solutions 625

576 **Improving Consistency.** Several approaches have 626
577 been proposed to address the inconsistency issue 627
578 in language models. Fine-tuning methods, such 628
579 as those explored by [Elazar et al. \(2021\)](#) and [Li 629](#)
580 [et al. \(2023d\)](#), aim to improve consistency by in- 630
581 troducing loss functions that penalize inconsis- 631
582 tent outputs or by selectively retaining only consis- 632
583 tent response pairs for training. [Jang and Lukasiewicz 633](#)
584 (2023) propose a plug-in method that leverages in- 634
585 termediate training with word-definition pairs to 635
586 enhance the model’s understanding of symbolic 636
587 meanings, thereby mitigating inconsistency. Out- 637
588 put ensemble approaches, such as those presented 638
589 by [Mitchell et al. \(2022\)](#) and [Zhao et al. \(2023b\)](#),

utilize multiple models to evaluate the consistency 590
of generated outputs. [Mitchell et al. \(2022\)](#) employ 591
a base model for generating potential answers and 592
a relation model for assessing their logical coher- 593
ence, while [Zhao et al. \(2023b\)](#) leverage LLMs to 594
rephrase questions and analyze the divergence of 595
corresponding answers to detect potential incon- 596
sistency. These diverse approaches highlight the 597
ongoing efforts to enhance the consistency and re- 598
liability of language models. 599

600 **Improving Factuality.** [Chuang et al. \(2023\)](#) and 600
601 [Li et al. \(2023c\)](#) propose methods that leverage the 601
602 inconsistency of knowledge across different lay- 602
603 ers. DoLa ([Chuang et al., 2023](#)) utilizes a dynamic 603
604 layer selection strategy, contrasting premature and 604
605 mature layers to determine the next word’s prob- 605
606 ability. ITI ([Li et al., 2023c](#)), on the other hand, 606
607 identifies truth-correlated attention heads based on 607
608 TruthfulQA ([Lin et al., 2022](#)) and shifts activations 608
609 along this direction during inference, repeating this 609
610 process autoregressively for each token. Both ap- 610
611 proaches aim to mitigate factual errors by effec- 611
612 tively utilizing the diverse knowledge representa- 612
613 tions within the model’s layers. 613

614 **Remarks.** The resolution of inter-memory conflict 614
615 in LLMs typically entails three phases: training, 615
616 generation, and post-hoc processing. The train- 616
617 ing phase method mainly focuses on mitigating 617
618 internal inconsistencies among model parameters. 618
619 Conversely, the generation and post-hoc phases 619
620 primarily involve algorithmic interventions aimed 620
621 at alleviating occurrences of inconsistent model 621
622 behavior. Nevertheless, the challenge persists in 622
623 addressing the inconsistency of parameter knowl- 623
624 edge without detrimentally impacting the overall 624
625 performance of LLMs. 625

626 5 Challenges and Future Directions 626

627 **Knowledge Conflicts in the Wild.** While current 627
628 research on knowledge conflicts primarily focuses 628
629 on artificially generated misinformation, real-world 629
630 conflicts often arise in retrieval-augmented LLMs 630
631 due to conflicting information retrieved from the 631
632 web. Existing analyses lack the realism of such 632
633 scenarios, potentially limiting the applicability of 633
634 their findings ([Xie et al., 2023](#); [Wang et al., 2023e](#)). 634
635 Recent work has begun to address this gap by curat- 635
636 ing conflicting documents based on actual Google 636
637 search results for open-ended questions ([Wan et al., 637](#)
638 [2024](#)). Future research should prioritize evaluat- 638
639 ing LLMs in these real-world scenarios to better 639

640 understand their capabilities and limitations.

641 **Solution at a Finer Resolution.** Resolving knowl- 691
642 edge conflicts presents a complex challenge, lack- 692
643 ing a universal solution. Conflicting information 693
644 can stem from misinformation, outdated facts, or 694
645 partially correct data (Uscinski and Butler, 2013; 695
646 Guo et al., 2022). Existing approaches often rely 696
647 on simple prior assumptions (Shi et al., 2023b). A 697
648 more nuanced approach is desired, considering the 698
649 query’s nature, the type of conflict, and user expect- 699
650 ations (Floridi, 2023), *e.g.*, subjective or debatable 700
651 questions inherently lead to conflicts due to multi- 701
652 ple valid answers (Bjerva et al., 2020; Wan et al., 702
653 2024). Future solutions should acknowledge the 703
654 diverse causes, manifestations, and potential user 704
655 expectations, requiring collaboration between NLP 705
656 and social science researchers for comprehensive 706
657 investigation and effective solutions. 707

658 **Evaluation on Downstream Tasks.** While re- 708
659 search on knowledge conflicts primarily focuses 709
660 on evaluating their performance on QA datasets, 710
661 the broader implications of these conflicts remain 711
662 underexplored. Their impact on downstream tasks, 712
663 particularly those demanding high accuracy and 713
664 consistency, such as legal document analysis (Shui 714
665 et al., 2023; Martin et al., 2024), medical diag- 715
666 nosis (Zhou et al., 2023b; Thirunavukarasu et al., 716
667 2023), financial analysis (Zhang et al., 2023a; Li 717
668 et al., 2023e), and educational tools (Caines et al., 718
669 2023; Milano et al., 2023), is crucial. Unresolved 719
670 knowledge conflicts could severely hinder the util- 720
671 ity of these models in such applications. 721

672 **Interplay among the Conflicts.** Current research 722
673 primarily focuses on individual conflict types or 723
674 a combined study of inter-context and context- 724
675 memory conflicts. However, the interplay between 725
676 intra-memory conflict and other types of conflicts 726
677 remains unexplored. Notably, several studies have 727
678 proposed the existence of knowledge circuits in 728
679 LLMs (Chughtai et al., 2024; Huang et al., 2023), 729
680 which are closely related to intra-memory con- 730
681 flict. Understanding this interaction is crucial for 731
682 comprehending the relationship between internal 732
683 knowledge inconsistency and model behavior in 733
684 response to context. Furthermore, exploring the 734
685 synergistic effects of various conflict types could 735
686 reveal underlying mechanisms of knowledge repre- 736
687 sentation and processing in LLMs is vital. 737

688 **Explainability.** While research has focused on ana- 738
689 lyzing LLMs’ outputs when faced with knowledge 739
690 conflicts, the internal mechanisms driving these de-

691 cisions remain underexplored. Studies examining 692
693 model confidence through logits (Xu et al., 2023; 693
694 Jin et al., 2024a; Wang et al., 2024) offer some in- 694
695 sights, but a deeper understanding of how specific 695
696 attention heads or neuron activations contribute to 696
697 conflict resolution is needed. Jin et al. (2024b) 697
698 made progress by investigating the interpretability 698
699 of LLMs through information flow analysis, identi- 699
700 fying memory and context heads with opposing ef- 700
701 fects in later layers. However, further microscopic 701
702 examinations are required to fully comprehend how 702
703 LLMs navigate conflicting information. 703

704 **Multilinguality.** Current research has primarily 704
705 focused on English. Future research should expand 705
706 to address conflicts in non-English texts, lever- 706
707 aging multilingual LLMs like GPT-4 (OpenAI, 707
708 2024) and GLM (Zeng et al., 2022) to account 708
709 for language-specific characteristics. Additionally, 709
710 inter-context conflict, involving documents in dif- 710
711 ferent languages, requires solutions like translation 711
712 systems (Dementieva and Panchenko, 2021), lever- 712
713 aging high-resource language evidence for low- 713
714 resource languages (Xue et al., 2024), or employing 714
715 knowledge distillation techniques. 715

716 **Multimodality.** While current research mainly fo- 716
717 cuses on text modality, potential conflicts arises as 717
718 LLMs evolve to process information across vari- 718
719 ous formats, including text, images (Alayrac et al., 719
720 2022; Li et al., 2023b), video (Ju et al., 2022; Zhang 720
721 et al., 2023b), and audio (Borsos et al., 2023; Wu 721
722 et al., 2023). For example, an audio clip might 722
723 contradict an accompanying document. Future re- 723
724 search should focus on enhancing models’ ability 724
725 to navigate these complex multimodal dynamics, 725
726 developing targeted datasets for training and evalu- 726
727 ation, and exploring user perception of multimodal 727
728 conflicts to improve LLMs. 728

6 Conclusion 728

729 This paper delves into the multifaceted issue of 729
730 knowledge conflicts, analyzing the categorization, 730
731 causes, behavior, and mitigation. We demonstrate 731
732 that the type of conflict significantly influences a 732
733 model’s behavior and that these conflicts exhibit 733
734 complex interplays. Existing solutions, often fo- 734
735 cused on artificial scenarios and relying on priors, 735
736 lack the granularity and breadth needed to address 736
737 the increasing complexity of knowledge conflicts 737
738 in real-world applications. As retrieval-augmented 738
739 LLMs become more prevalent, comprehensive re- 739
740 search on knowledge conflicts is crucial. 740

741 Limitations

742 Considering the rapid expansion of research in the
743 field of knowledge conflict and the abundance of
744 scholarly literature, it is possible that we might
745 have missed some of the most recent or less rele-
746 vant findings. Nevertheless, we have ensured the
747 inclusion of all essential materials in our survey.

748 Ethics Statement

749 We mainly searched for papers published after 2021
750 using key terms including “knowledge conflict”,
751 “knowledge inconsistency”, “knowledge gap”, *inter*
752 *alia*, on Google Scholar and the ACL Anthology.
753 After initially identifying these papers, the authors
754 classified them through reading and continued to
755 track related but overlooked papers using their ci-
756 tations. We also used Google Scholar to follow up
757 on the latest papers citing these to avoid omissions.

758 For the quantitative analysis and comparison sec-
759 tion (§ F), we did not conduct computational exper-
760 iments but simply organized the result reported in
761 other literature as is.

762 References

763 Shourya Aggarwal, Divyanshu Mandowara, Vishwa-
764 jeet Agrawal, Dinesh Khandelwal, Parag Singla, and
765 Dinesh Garg. 2021. [Explanations for Common-](#)
766 [senseQA: New Dataset and Models](#). In *Proceedings*
767 *of the 59th Annual Meeting of the Association for*
768 *Computational Linguistics and the 11th International*
769 *Joint Conference on Natural Language Processing*
770 *(Volume 1: Long Papers)*, pages 3050–3065, Online.
771 Association for Computational Linguistics.

772 Ayush Agrawal, Lester Mackey, and Adam Tauman
773 Kalai. 2023. [Do language models know when](#)
774 [they’re hallucinating references?](#) *ArXiv preprint*,
775 [abs/2305.18248](#).

776 Ekin Akyürek, Tolga Bolukbasi, Frederick Liu, Bin-
777 bin Xiong, Ian Tenney, Jacob Andreas, and Kelvin
778 Guu. 2022. Towards tracing knowledge in language
779 models back to the training data. In *Findings of the*
780 *Association for Computational Linguistics: EMNLP*
781 *2022*, pages 2429–2446.

782 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,
783 Antoine Miech, Iain Barr, Yana Hasson, Karel
784 Lenc, Arthur Mensch, Katherine Millican, Malcolm
785 Reynolds, et al. 2022. Flamingo: a visual language
786 model for few-shot learning. *Advances in neural*
787 *information processing systems*, 35:23716–23736.

788 Navid Ayoobi, Sadat Shahriar, and Arjun Mukherjee.
789 2023. The looming threat of fake and llm-generated
790 linkedin profiles: Challenges and opportunities for

detection and prevention. In *Proceedings of the 34th*
ACM Conference on Hypertext and Social Media,
pages 1–10. 791
792
793

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 confer-*
ence on fairness, accountability, and transparency,
pages 610–623. 794
795
796
797
798
799

Johannes Bjerva, Nikita Bhutani, Behzad Golshan,
Wang-Chiew Tan, and Isabelle Augenstein. 2020.
[SubjQA: A Dataset for Subjectivity and Review Com-](#)
[prehension](#). In *Proceedings of the 2020 Conference*
on Empirical Methods in Natural Language Process-
ing (EMNLP), pages 5480–5494, Online. Association
for Computational Linguistics. 800
801
802
803
804
805
806

Zalán Borsos, Raphaël Marinier, Damien Vincent,
Eugene Kharitonov, Olivier Pietquin, Matt Shar-
ifi, Dominik Roblek, Olivier Teboul, David Grang-
ier, Marco Tagliasacchi, et al. 2023. Audiollm: a
language modeling approach to audio generation.
IEEE/ACM Transactions on Audio, Speech, and Lan-
guage Processing. 807
808
809
810
811
812
813

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
Neelakantan, Pranav Shyam, Girish Sastry, Amanda
Askell, Sandhini Agarwal, Ariel Herbert-Voss,
Gretchen Krueger, Tom Henighan, Rewon Child,
Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
Clemens Winter, Christopher Hesse, Mark Chen, Eric
Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
Jack Clark, Christopher Berner, Sam McCandlish,
Alec Radford, Ilya Sutskever, and Dario Amodei.
2020. [Language models are few-shot learners](#). In *Ad-*
vances in Neural Information Processing Systems 33:
Annual Conference on Neural Information Process-
ing Systems 2020, NeurIPS 2020, December 6-12,
2020, virtual. 814
815
816
817
818
819
820
821
822
823
824
825
826
827
828

Andrew Caines, Luca Benedetto, Shiva Taslimipoor,
Christopher Davis, Yuan Gao, Oeistein Andersen,
Zheng Yuan, Mark Elliott, Russell Moore, Christo-
pher Bryant, et al. 2023. [On the application of large](#)
[language models for language teaching and assess-](#)
[ment technology](#). *ArXiv preprint*, [abs/2307.08393](#). 829
830
831
832
833
834

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*, [abs/2302.10149](#). 835
836
837
838
839

Chi Cheang, Hou Chan, Derek Wong, Xuebo Liu, Zhao-
cong Li, Yanming Sun, Shudong Liu, and Lidia Chao.
2023. Can lms generalize to future data? an empiri-
cal analysis on text summarization. In *Proceedings*
of the 2023 Conference on Empirical Methods in
Natural Language Processing, pages 16205–16217. 840
841
842
843
844
845

Canyu Chen and Kai Shu. 2023a. Can llm-generated
misinformation be detected? In *NeurIPS 2023 Work-*
846
847

848			
849		<i>shop on Instruction Tuning and Instruction Following.</i>	
850	Canyu Chen and Kai Shu. 2023b.	Combating misinformation in the age of llms: Opportunities and challenges. <i>ArXiv preprint</i> , abs/2311.05656.	
851			
852			
853	Hung-Ting Chen, Michael JQ Zhang, and Eunsol Choi.		
854		2022. Rich knowledge sources bring complex knowledge conflicts: Recalibrating models to reflect conflicting evidence. <i>ArXiv preprint</i> , abs/2210.13701.	
855			
856			
857	Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun.		
858		2023. Benchmarking large language models in retrieval-augmented generation. <i>ArXiv preprint</i> , abs/2309.01431.	
859			
860			
861	Wenhu Chen, Xinyi Wang, and William Yang Wang.		
862		2021. A dataset for answering time-sensitive questions. <i>ArXiv preprint</i> , abs/2108.06314.	
863			
864	I Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua		
865	Feng, Chunting Zhou, Junxian He, Graham Neubig,		
866	Pengfei Liu, et al. 2023. Factool: Factuality detection in generative ai—a tool augmented framework for multi-task and multi-domain scenarios. <i>ArXiv preprint</i> , abs/2307.13528.		
867			
868			
869			
870	Tsun-Hin Cheung and Kin-Man Lam. 2023. Factllama:		
871		Optimizing instruction-following language models with external knowledge for automated fact-checking.	
872		In <i>2023 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)</i> , pages 846–853. IEEE.	
873			
874			
875			
876	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,		
877	Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul		
878	Barham, Hyung Won Chung, Charles Sutton, Sebastian		
879	Gehrmann, et al. 2023. Palm: Scaling language		
880	modeling with pathways. <i>Journal of Machine Learning Research</i> , 24(240):1–113.		
881			
882	Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon		
883	Kim, James Glass, and Pengcheng He. 2023. Dola: Decoding by contrasting layers improves factuality in large language models. <i>ArXiv preprint</i> , abs/2309.03883.		
884			
885			
886			
887	Bilal Chughtai, Alan Cooney, and Neel Nanda. 2024.		
888		Summing up the facts: Additive mechanisms behind factual recall in llms. <i>ArXiv preprint</i> , abs/2402.07321.	
889			
890			
891	Christopher Clark, Kenton Lee, Ming-Wei Chang,		
892	Tom Kwiatkowski, Michael Collins, and Kristina		
893	Toutanova. 2019. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 2924–2936, Minneapolis, Minnesota. Association for Computational Linguistics.		
894			
895			
896			
897			
898			
899			
900	Leyang Cui, Sijie Cheng, Yu Wu, and Yue Zhang. 2020.		
901		Does bert solve commonsense task via commonsense knowledge. <i>ArXiv preprint</i> , abs/2008.03945.	
902			
	Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing		
	factual knowledge in language models. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 6491–6506.		903 904 905 906 907
	Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah		
	Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. 2021. A continual learning survey: Defying forgetting in classification tasks. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 44(7):3366–3385.		908 909 910 911 912 913
	Daryna Dementieva and Alexander Panchenko. 2021. Cross-lingual evidence improves monolingual fake news detection. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: Student Research Workshop</i> , pages 310–320, Online. Association for Computational Linguistics.		914 915 916 917 918 919 920 921
	Bhuwan Dhingra, Jeremy R Cole, Julian Martin		
	Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W Cohen. 2022. Time-aware language models as temporal knowledge bases. <i>Transactions of the Association for Computational Linguistics</i> , 10:257–273.		922 923 924 925 926 927
	Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu,		
	Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2023. Chain-of-verification reduces hallucination in large language models. <i>ArXiv preprint</i> , abs/2309.11495.		928 929 930 931 932
	Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Zhifang		
	Sui, and Lei Li. 2023. Statistical knowledge assessment for large language models. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .		933 934 935 936 937
	Li Du, Xiao Ding, Kai Xiong, Ting Liu, and Bing Qin.		
	2022a. e-care: a new dataset for exploring explainable causal reasoning. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 432–446.		938 939 940 941 942
	Yibing Du, Antoine Bosselut, and Christopher D Manning. 2022b. Synthetic disinformation attacks on automated fact verification systems. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pages 10581–10589.		943 944 945 946 947
	Nouha Dziri, Andrea Madotto, Osmar Zaiane, and		
	Avishek Joey Bose. 2021. Neural path hunter: Reducing hallucination in dialogue systems via path grounding. <i>ArXiv preprint</i> , abs/2104.08455.		948 949 950 951
	Yanai Elazar, Nora Kassner, Shauli Ravfogel, Amir		
	Feder, Abhilasha Ravichander, Marius Mosbach, Yonatan Belinkov, Hinrich Schütze, and Yoav Goldberg. 2022. Measuring causal effects of data statistics on language model’s factual predictions. <i>ArXiv preprint</i> , abs/2207.14251.		952 953 954 955 956 957

958	Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Schütze, and Yoav Goldberg. 2021. Measuring and improving consistency in pretrained language models. <i>Transactions of the Association for Computational Linguistics</i> , 9:1012–1031.	1013
959		1014
960		1015
961		1016
962		1017
963		1018
964		1019
965	Hady Elsahar, Pavlos Vougiouklis, Arslan Remaci, Christophe Gravier, Jonathon Hare, Frederique Laforest, and Elena Simperl. 2018. T-REx: A large scale alignment of natural language with knowledge base triples . In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association (ELRA).	1020
966		1021
967		1022
968		1023
969		1024
970		1025
971		1026
972	Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 889–898, Melbourne, Australia. Association for Computational Linguistics.	1027
973		1028
974		1029
975		1030
976		1031
977		1032
978	Zhangyin Feng, Weitao Ma, Weijiang Yu, Lei Huang, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. Trends in integration of knowledge and large language models: A survey and taxonomy of methods, benchmarks, and applications . <i>ArXiv preprint</i> , abs/2311.05876.	1033
979		1034
980		1035
981		1036
982		1037
983		1038
984	Luciano Floridi. 2023. Ai as agency without intelligence: on chatgpt, large language models, and other generative models. <i>Philosophy & Technology</i> , 36(1):15.	1039
985		1040
986		1041
987		1042
988	Zorik Gekhman, Jonathan Herzig, Roei Aharoni, Chen Elkind, and Idan Szpektor. 2023. Trueteacher: Learning factual consistency evaluation with large language models . <i>ArXiv preprint</i> , abs/2305.11171.	1043
989		1044
990		1045
991		1046
992		1047
993		1048
994	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. <i>Transactions of the Association for Computational Linguistics</i> , 9:346–361.	1049
995		1050
996		1051
997		1052
998		1053
999	Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. 2023. More than you’ve asked for: A comprehensive analysis of novel prompt injection threats to application-integrated large language models. <i>arXiv e-prints</i> , pages arXiv–2302.	1054
1000		1055
1001		1056
1002		1057
1003		1058
1004	Roger Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini, Benoit Steiner, Dustin Li, Esin Durmus, Ethan Perez, et al. 2023. Studying large language model generalization with influence functions . <i>ArXiv preprint</i> , abs/2308.03296.	1059
1005		1060
1006		1061
1007		1062
1008		1063
1009	Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking. <i>Transactions of the Association for Computational Linguistics</i> , 10:178–206.	1064
1010		1065
1011		1066
1012		1067
		1068
		1069
		1070
		1071
		1072
		1073
		1074
		1075
		1076
		1077
		1078
		1079
		1080
		1081
		1082
		1083
		1084
		1085
		1086
		1087
		1088
		1089
		1090
		1091
		1092
		1093
		1094
		1095
		1096
		1097
		1098
		1099
		1100
		1101
		1102
		1103
		1104
		1105
		1106
		1107
		1108
		1109
		1110
		1111
		1112
		1113
		1114
		1115
		1116
		1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
		1153
		1154
		1155
		1156
		1157
		1158
		1159
		1160
		1161
		1162
		1163
		1164
		1165
		1166
		1167
		1168
		1169
		1170
		1171
		1172
		1173
		1174
		1175
		1176
		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209
		1210
		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222
		1223
		1224
		1225
		1226
		1227
		1228
		1229
		1230
		1231
		1232
		1233
		1234
		1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
		1259
		1260
		1261
		1262
		1263
		1264
		1265
		1266
		1267
		1268
		1269
		1270
		1271
		1272
		1273
		1274
		1275
		1276
		1277
		1278
		1279
		1280
		1281
		1282
		1283
		1284
		1285
		1286
		1287
		1288
		1289
		1290
		1291
		1292
		1293
		1294
		1295
		1296
		1297
		1298
		1299
		1300
		1301
		1302
		1303
		1304
		1305
		1306
		1307
		1308
		1309
		1310
		1311
		1312
		1313
		1314
		1315
		1316
		1317
		1318
		1319
		1320
		1321
		1322
		1323
		1324
		1325
		1326
		1327
		1328
		1329
		1330
		1331
		1332
		1333
		1334
		1335
		1336
		1337
		1338
		1339
		1340
		1341
		1342
		1343
		1344
		1345
		1346
		1347
		1348
		1349
		1350
		1351
		1352
		1353
		1354
		1355
		1356
		1357
		1358
		1359
		1360
		1361
		1362
		1363
		1364
		1365
		1366
		1367
		1368
		1369
		1370
		1371
		1372
		1373
		1374
		1375
		1376
		1377
		1378
		1379
		1380
		1381
		1382
		1383
		1384
		1385
		1386
		1387
		1388
		1389
		1390
		1391
		1392
		1393
		1394
		1395
		1396
		1397
		1398
		1399
		1400
		1401
		1402
		1403
		1404
		1405
		1406
		1407
		1408
		1409
		1410
		1411
		1412
		1413
		1414
		1415
		1416
		1417
		1418
		1419
		1420
		1421
		1422
		1423
		1424
		1425
		1426
		1427
		1428
		1429
		1430
		1431

1070	Joel Jang, Seonghyeon Ye, Changho Lee, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, and Minjoon Seo. 2022. Temporalwiki: A lifelong benchmark for training and evaluating ever-evolving language models. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 6237–6250.	1126
1071		1127
1072		1128
1073		1129
1074		1130
1075		
1076		
1077	Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, KIM Gyeonghun, Stanley Jungkyu Choi, and Minjoon Seo. 2021. Towards continual knowledge learning of language models. In <i>International Conference on Learning Representations</i> .	
1078		
1079		
1080		
1081		
1082	Myeongjun Erik Jang and Thomas Lukasiewicz. 2023. Improving language models meaning understanding and consistency by learning conceptual roles from dictionary. <i>ArXiv preprint</i> , abs/2310.15541.	
1083		
1084		
1085		
1086	Ganesh Jawahar, Muhammad Abdul-Mageed, and Laks Lakshmanan, V.S. 2020. Automatic detection of machine generated text: A critical survey. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 2296–2309, Barcelona, Spain (Online). International Committee on Computational Linguistics.	
1087		
1088		
1089		
1090		
1091		
1092		
1093	Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. What does BERT learn about the structure of language? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3651–3657, Florence, Italy. Association for Computational Linguistics.	
1094		
1095		
1096		
1097		
1098		
1099	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM Computing Surveys</i> , 55(12):1–38.	
1100		
1101		
1102		
1103		
1104	Bohan Jiang, Zhen Tan, Ayushi Nirmal, and Huan Liu. 2023. Disinformation detection: An evolving challenge in the age of llms. <i>ArXiv preprint</i> , abs/2309.15847.	
1105		
1106		
1107		
1108	Zhuoran Jin, Pengfei Cao, Yubo Chen, Kang Liu, Xiaojian Jiang, Jiexin Xu, Qiuxia Li, and Jun Zhao. 2024a. Tug-of-war between knowledge: Exploring and resolving knowledge conflicts in retrieval-augmented language models. <i>ArXiv preprint</i> , abs/2402.14409.	
1109		
1110		
1111		
1112		
1113	Zhuoran Jin, Pengfei Cao, Hongbang Yuan, Yubo Chen, Jiexin Xu, Huaijun Li, Xiaojian Jiang, Kang Liu, and Jun Zhao. 2024b. Cutting off the head ends the conflict: A mechanism for interpreting and mitigating knowledge conflicts in language models. <i>ArXiv preprint</i> , abs/2402.18154.	
1114		
1115		
1116		
1117		
1118		
1119	Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.	
1120		
1121		
1122		
1123		
1124		
1125		
	Chen Ju, Tengda Han, Kunhao Zheng, Ya Zhang, and Weidi Xie. 2022. Prompting visual-language models for efficient video understanding. In <i>European Conference on Computer Vision</i> , pages 105–124. Springer.	1131
		1132
		1133
		1134
	Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert McHardy. 2023. Challenges and applications of large language models. <i>ArXiv preprint</i> , abs/2307.10169.	1135
		1136
		1137
		1138
		1139
	Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In <i>International Conference on Machine Learning</i> , pages 15696–15707. PMLR.	1140
		1141
		1142
	Cheongwoong Kang and Jaesik Choi. 2023. Impact of co-occurrence on factual knowledge of large language models. <i>ArXiv preprint</i> , abs/2310.08256.	1143
		1144
		1145
		1146
		1147
		1148
		1149
	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6769–6781, Online. Association for Computational Linguistics.	1150
		1151
		1152
		1153
		1154
	Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir Radev, Noah A Smith, Yejin Choi, and Kentaro Inui. 2022. Realtime qa: What’s the answer right now? <i>ArXiv preprint</i> , abs/2207.13332.	1155
		1156
	Celeste Kidd and Abeba Birhane. 2023. How ai can distort human beliefs. <i>Science</i> , 380(6651):1222–1223.	1157
		1158
		1159
		1160
	Miyoung Ko, Ingyu Seong, Hwaran Lee, Joonsuk Park, Minsuk Chang, and Minjoon Seo. 2022. Claimdiff: Comparing and contrasting claims on contentious issues. <i>ArXiv preprint</i> , abs/2205.12221.	1161
		1162
		1163
		1164
		1165
	Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. 2018. The NarrativeQA reading comprehension challenge. <i>Transactions of the Association for Computational Linguistics</i> , 6:317–328.	1166
		1167
		1168
		1169
		1170
		1171
		1172
		1173
		1174
	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research. <i>Transactions of the Association for Computational Linguistics</i> , 7:452–466.	1175
		1176
		1177
		1178
	Tsz Kin Lam, Eva Hasler, and Felix Hieber. 2022. Analyzing the use of influence functions for instance-specific data filtering in neural machine translation. <i>ArXiv preprint</i> , abs/2210.13281.	

1290	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das,	Humza Naveed, Asad Ullah Khan, Shi Qiu, Muham-	1344
1291	Daniel Khashabi, and Hannaneh Hajishirzi. 2023.	mad Saqib, Saeed Anwar, Muhammad Usman, Nick	1345
1292	When not to trust language models: Investigating	Barnes, and Ajmal Mian. 2023. A comprehensive	1346
1293	effectiveness of parametric and non-parametric mem-	overview of large language models . <i>ArXiv preprint</i> ,	1347
1294	ories. In <i>Proceedings of the 61st Annual Meeting of</i>	abs/2307.06435 .	1348
1295	<i>the Association for Computational Linguistics (Vol-</i>		
1296	<i>ume 1: Long Papers)</i> , pages 9802–9822.		
1297	Potsawee Manakul, Adian Liusie, and Mark JF Gales.	Ella Neeman, Roei Aharoni, Or Honovich, Leshem	1349
1298	2023. Selfcheckgpt: Zero-resource black-box hal-	Choshen, Idan Szpektor, and Omri Abend. 2022.	1350
1299	lucination detection for generative large language	Disentqa: Disentangling parametric and contextual	1351
1300	models . <i>ArXiv preprint</i> , abs/2303.08896 .	knowledge with counterfactual question answering .	1352
		<i>ArXiv preprint</i> , abs/2211.05655 .	1353
1301	Katerina Margatina, Shuai Wang, Yogarshi Vyas,	Raymond S Nickerson. 1998. Confirmation bias: A	1354
1302	Neha Anna John, Yassine Benajiba, and Miguel	ubiquitous phenomenon in many guises. <i>Review of</i>	1355
1303	Ballesteros. 2023. Dynamic benchmarking of	<i>general psychology</i> , 2(2):175–220.	1356
1304	masked language models on temporal concept drift		
1305	with multiple views . <i>ArXiv preprint</i> , abs/2302.12297 .	Xenia Ohmer, Elia Bruni, and Dieuwke Hupkes. 2023.	1357
		Separating form and meaning: Using self-consistency	1358
1306	Lauren Martin, Nick Whitehouse, Stephanie Yiu, Lizzie	to quantify task understanding across multiple senses.	1359
1307	Catterson, and Rivindu Perera. 2024. Better call gpt,	<i>CoRR</i> .	1360
1308	comparing large language models against lawyers .		
1309	<i>ArXiv preprint</i> , abs/2401.16212 .	Yasumasa Onoe, Michael JQ Zhang, Shankar Padman-	1361
		abhan, Greg Durrett, and Eunsol Choi. 2023. Can	1362
1310	Luca Massarelli, Fabio Petroni, Aleksandra Piktus,	lms learn new entities from descriptions? challenges	1363
1311	Myle Ott, Tim Rocktäschel, Vassilis Plachouras, Fab-	in propagating injected knowledge . <i>ArXiv preprint</i> ,	1364
1312	rizio Silvestri, and Sebastian Riedel. 2020. How de-	abs/2305.01651 .	1365
1313	coding strategies affect the verifiability of generated		
1314	text . In <i>Findings of the Association for Computa-</i>	OpenAI. 2024. Gpt-4 technical report .	1366
1315	<i>tional Linguistics: EMNLP 2020</i> , pages 223–235,		
1316	Online. Association for Computational Linguistics.	Liangming Pan, Wenhui Chen, Min-Yen Kan, and	1367
		William Yang Wang. 2023a. Attacking open-domain	1368
1317	Kevin Meng, David Bau, Alex Andonian, and Yonatan	question answering by injecting misinformation .	1369
1318	Belinkov. 2022. Locating and editing factual associ-	<i>IJCNLP-AAACL ACL</i> .	1370
1319	ations in gpt . <i>Advances in Neural Information Pro-</i>		
1320	<i>cessing Systems</i> , 35:17359–17372.	Xiaoman Pan, Wenlin Yao, Hongming Zhang, Dian Yu,	1371
		Dong Yu, and Jianshu Chen. 2022. Knowledge-in-	1372
1321	Stephen Merity, Caiming Xiong, James Bradbury, and	context: Towards knowledgeable semi-parametric	1373
1322	Richard Socher. 2017. Pointer sentinel mixture mod-	language models . In <i>The Eleventh International Con-</i>	1374
1323	els . In <i>5th International Conference on Learning</i>	<i>ference on Learning Representations</i> .	1375
1324	<i>Representations, ICLR 2017, Toulon, France, April</i>		
1325	<i>24-26, 2017, Conference Track Proceedings</i> . Open-	Yikang Pan, Liangming Pan, Wenhui Chen, Preslav	1376
1326	Review.net.	Nakov, Min-Yen Kan, and William Yang Wang.	1377
		2023b. On the risk of misinformation pollu-	1378
1327	Silvia Milano, Joshua A McGrane, and Sabina Leonelli.	tion with large language models . <i>ArXiv preprint</i> ,	1379
1328	2023. Large language models challenge the future	abs/2305.13661 .	1380
1329	of higher education . <i>Nature Machine Intelligence</i> ,	Ethan Perez, Sam Ringer, Kamilè Lukošiušė, Karina	1381
1330	5(4):333–334.	Nguyen, Edwin Chen, Scott Heiner, Craig Pettit,	1382
		Catherine Olsson, Sandipan Kundu, Saurav Kada-	1383
1331	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea	vath, et al. 2022. Discovering language model behav-	1384
1332	Finn, and Christopher D Manning. 2021. Fast model	iors with model-written evaluations . <i>ArXiv preprint</i> ,	1385
1333	editing at scale . <i>ArXiv preprint</i> , abs/2110.11309 .	abs/2212.09251 .	1386
		Fabio Petroni, Tim Rocktäschel, Sebastian Riedel,	1387
1334	Eric Mitchell, Joseph J Noh, Siyan Li, William S Arm-	Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and	1388
1335	strong, Ananth Agarwal, Patrick Liu, Chelsea Finn,	Alexander Miller. 2019. Language models as knowl-	1389
1336	and Christopher D Manning. 2022. Enhancing self-	edge bases? In <i>Proceedings of the 2019 Confer-</i>	1390
1337	consistency and performance of pre-trained language	<i>ence on Empirical Methods in Natural Language Pro-</i>	1391
1338	models through natural language inference . <i>ArXiv</i>	<i>cessing and the 9th International Joint Conference</i>	1392
1339	<i>preprint</i> , abs/2211.11875 .	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	1393
		pages 2463–2473, Hong Kong, China. Association	1394
1340	Niels Mündler, Jingxuan He, Slobodan Jenko, and Mar-	for Computational Linguistics.	1395
1341	tin Vechev. 2023. Self-contradictory hallucinations		
1342	of large language models: Evaluation, detection and	Maren Pielka, Felix Rode, Lisa Pucknat, Tobias Deußer,	1396
1343	mitigation . <i>ArXiv preprint</i> , abs/2305.15852 .	and Rafet Sifa. 2022. A linguistic investigation of	1397

1398	machine learning based contradiction detection models: an empirical analysis and future perspectives. In <i>2022 21st IEEE International Conference on Machine Learning and Applications (ICMLA)</i> , pages 1649–1653. IEEE.	Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Scott Wen-tau Yih. 2023a. Trusting your evidence: Hallucinate less with context-aware decoding . <i>ArXiv preprint</i> , abs/2305.14739.	1452
1399			1453
1400			1454
1401			1455
1402			1456
1403	Yuval Pinter and Michael Elhadad. 2023. Emptying the ocean with a spoon: Should we edit models? In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 15164–15172.	Weijia Shi, Sewon Min, Maria Lomeli, Chunting Zhou, Margaret Li, Victoria Lin, Noah A Smith, Luke Zettlemoyer, Scott Yih, and Mike Lewis. 2023b. In-context pretraining: Language modeling beyond document boundaries . <i>ArXiv preprint</i> , abs/2310.10638.	1457
1404			1458
1405			1459
1406			1460
1407	Jirui Qi, Raquel Fernández, and Arianna Bisazza. 2023. Cross-lingual consistency of factual knowledge in multilingual language models . <i>ArXiv preprint</i> , abs/2310.10378.	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2023c. Replug: Retrieval-augmented black-box language models . <i>ArXiv preprint</i> , abs/2301.12652.	1462
1408			1463
1409			1464
1410			1465
1411	Cheng Qian, Xinran Zhao, and Sherry Tongshuang Wu. 2023. "merge conflicts!" exploring the impacts of external distractors to parametric knowledge graphs . <i>ArXiv preprint</i> , abs/2309.08594.		1466
1412			1467
1413			1468
1414			1469
1415	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for SQuAD . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 784–789, Melbourne, Australia. Association for Computational Linguistics.	Ruihao Shui, Yixin Cao, Xiang Wang, and Tat-Seng Chua. 2023. A comprehensive evaluation of large language models on legal judgment prediction . <i>ArXiv preprint</i> , abs/2310.11761.	1470
1416			1471
1417			1472
1418			1473
1419			1474
1420			1475
1421			1476
1422	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. Retrieval augmentation reduces hallucination in conversation . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 3784–3803.	1477
1423			1478
1424			1479
1425			1480
1426			1481
1427			1482
1428	Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How much knowledge can you pack into the parameters of a language model? In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5418–5426, Online. Association for Computational Linguistics.	Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2022. Large language models encode clinical knowledge . <i>ArXiv preprint</i> , abs/2212.13138.	1483
1429			1484
1430			1485
1431			1486
1432			1487
1433			1488
1434	Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in BERTology: What we know about how BERT works . <i>Transactions of the Association for Computational Linguistics</i> , 8:842–866.	Anton Sinitsin, Vsevolod Plokhotnyuk, Dmitriy Pyrkin, Sergei Popov, and Artem Babenko. 2020. Editable neural networks . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	1489
1435			1490
1436			1491
1437			1492
1438	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools . <i>ArXiv preprint</i> , abs/2302.04761.	Craig S. Smith. 2023. What large models cost you – there is no free ai lunch .	1493
1439			1494
1440			1495
1441			1496
1442			1497
1443	Or Sharir, Barak Peleg, and Yoav Shoham. 2020. The cost of training nlp models: A concise overview . <i>ArXiv preprint</i> , abs/2004.08900.	Giovanni Spitale, Nikola Biller-Andorno, and Federico Germani. 2023. Ai model gpt-3 (dis) informs us better than humans . <i>ArXiv preprint</i> , abs/2301.11924.	1498
1444			1499
1445			1500
1446	Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R Bowman, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R Johnston, et al. 2023. Towards understanding sycophancy in language models . <i>ArXiv preprint</i> , abs/2310.13548.	Hexiang Tan, Fei Sun, Wanli Yang, Yuanzhuo Wang, Qi Cao, and Xueqi Cheng. 2024. Blinded by generated contexts: How language models merge generated and retrieved contexts for open-domain qa? <i>ArXiv preprint</i> , abs/2401.11911.	1501
1447			1502
1448			1503
1449			1504
1450			1505
1451			1506

1505	Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. <i>Nature medicine</i> , 29(8):1930–1940.	1559
1506		1560
1507		1561
1508		1562
1509		1563
1510	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>ArXiv preprint</i> , abs/2307.09288.	1564
1511		1565
1512		1566
1513		1567
1514		
1515		
1516	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Musique: Multi-hop questions via single-hop question composition. <i>Transactions of the Association for Computational Linguistics</i> , 10:539–554.	1568
1517		1569
1518		1570
1519		1571
1520		1572
1521	Miles Turpin, Julian Michael, Ethan Perez, and Samuel R Bowman. 2023. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. <i>ArXiv preprint</i> , abs/2305.04388.	1573
1522		1574
1523		1575
1524		1576
1525		1577
1526	Joseph E Uscinski and Ryden W Butler. 2013. The epistemology of fact checking. <i>Critical Review</i> , 25(2):162–180.	1578
1527		1579
1528		1580
1529	Tyler Vergho, Jean-Francois Godbout, Reihaneh Rab-bany, and Kellin Pelrine. 2024. Comparing gpt-4 and open-source language models in misinformation mitigation . <i>ArXiv preprint</i> , abs/2401.06920.	1581
1530		1582
1531		1583
1532		1584
1533	Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, et al. 2023. Freshllms: Refreshing large language models with search engine augmentation . <i>ArXiv preprint</i> , abs/2310.03214.	1585
1534		1586
1535		
1536		
1537		
1538	Alexander Wan, Eric Wallace, and Dan Klein. 2024. What evidence do language models find convincing? <i>ArXiv preprint</i> , abs/2402.11782.	1587
1539		1588
1540		1589
1541	Cunxiang Wang, Shuailong Liang, Yue Zhang, Xiaonan Li, and Tian Gao. 2019. Does it make sense? and why? a pilot study for sense making and explanation . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4020–4026, Florence, Italy. Association for Computational Linguistics.	1590
1542		1591
1543		1592
1544		1593
1545		1594
1546		1595
1547		1596
1548	Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao, Wenyang Gao, Xuming Hu, Zehan Qi, Yidong Wang, Linyi Yang, Jindong Wang, Xing Xie, Zheng Zhang, and Yue Zhang. 2023a. Survey on factuality in large language models: Knowledge, retrieval and domain-specificity .	1597
1549		1598
1550		1599
1551		
1552		
1553		
1554		
1555	Fei Wang, Wenjie Mo, Yiwei Wang, Wenxuan Zhou, and Muhao Chen. 2023b. A causal view of entity bias in (large) language models . <i>ArXiv preprint</i> , abs/2305.14695.	1600
1556		1601
1557		1602
1558		1603
		1604
		1605
		1606
		1607
		1608
		1609
		1610
		1611
		1612
		1613
		1614

1615	models in knowledge conflicts.	<i>ArXiv preprint</i> ,	<i>the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> ,	1670
1616		abs/2305.13300.	pages 543–553.	1671
1617	Nan Xu, Fei Wang, Bangzheng Li, Mingtao Dong, and			1672
1618	Muhao Chen. 2022. Does your model classify entities reasonably? diagnosing and mitigating spurious correlations in entity typing.	<i>ArXiv preprint</i> ,	Jiaxin Zhang, Zhuohang Li, Kamalika Das, Bradley A	1673
1619		abs/2205.12640.	Malin, and Sricharan Kumar. 2023c. Sac³: Reliable hallucination detection in black-box language models via semantic-aware cross-check consistency.	1674
1620			<i>ArXiv preprint</i> ,	1675
1621			abs/2311.01740.	1676
1622	Rongwu Xu, Brian S Lin, Shujian Yang, Tianqi Zhang,			1677
1623	Weiyan Shi, Tianwei Zhang, Zhixuan Fang, Wei		Michael JQ Zhang and Eunsol Choi. 2021. Situatedaq: Incorporating extra-linguistic contexts into qa.	1678
1624	Xu, and Han Qiu. 2023. The earth is flat because...: Investigating llms' belief towards misinformation via persuasive conversation.	<i>arXiv preprint</i>	<i>ArXiv preprint</i> ,	1679
1625		arXiv:2312.09085.	abs/2109.06157.	1680
1626				
1627				
1628	Rongwu Xu, Zehan Qi, and Wei Xu. 2024. Preemptive answer" attacks" on chain-of-thought reasoning.	<i>arXiv preprint</i>	Michael JQ Zhang and Eunsol Choi. 2023. Mitigating temporal misalignment by discarding outdated facts.	1681
1629		arXiv:2405.20902.	<i>ArXiv preprint</i> ,	1682
1630			abs/2305.14824.	1683
1631	Boyang Xue, Hongru Wang, Weichao Wang, Rui Wang,			
1632	Sheng Wang, Zeming Liu, and Kam-Fai Wong. 2024.		Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu,	1684
1633	A comprehensive study of multilingual confidence estimation on large language models.		Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang,	1685
1634			Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei	1686
1635	Boyang Xue, Weichao Wang, Hongru Wang, Fei Mi,		Bi, Freda Shi, and Shuming Shi. 2023d. Siren's song in the ai ocean: A survey on hallucination in large language models.	1687
1636	Rui Wang, Yasheng Wang, Lifeng Shang, Xin Jiang,			1688
1637	Qun Liu, and Kam-Fai Wong. 2023. Improving factual consistency for knowledge-grounded dialogue systems via knowledge enhancement and alignment.		Yunxiang Zhang, Muhammad Khalifa, Lajanugen Logeswaran, Moontae Lee, Honglak Lee, and Lu Wang.	1690
1638			2023e. Merging generated and retrieved knowledge for open-domain qa.	1691
1639			<i>ArXiv preprint</i> ,	1692
1640			abs/2310.14393.	1693
1641				
1642	Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng,		Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu,	1694
1643	Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu		Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei	1695
1644	Zhang. 2023. Editing large language models: Problems, methods, and opportunities.	<i>ArXiv preprint</i> ,	Yin, and Mengnan Du. 2023a. Explainability for large language models: A survey.	1696
1645		abs/2305.13172.	<i>ACM Transactions on Intelligent Systems and Technology.</i>	1697
1646				1698
1647	Jingwei Yi, Yueqi Xie, Bin Zhu, Keegan Hines, Emre		Yukun Zhao, Lingyong Yan, Weiwei Sun, Guoliang	1699
1648	Kiciman, Guangzhong Sun, Xing Xie, and Fangzhao		Xing, Chong Meng, Shuaiqiang Wang, Zhicong	1700
1649	Wu. 2023. Benchmarking and defending against indirect prompt injection attacks on large language models.	<i>ArXiv preprint</i> ,	Cheng, Zhaochun Ren, and Dawei Yin. 2023b.	1701
1650		abs/2312.14197.	Knowing what llms do not know: A simple yet effective self-detection method.	1702
1651			<i>ArXiv preprint</i> ,	1703
1652	Jiahao Ying, Yixin Cao, Kai Xiong, Yidong He, Long		abs/2310.17918.	1704
1653	Cui, and Yongbin Liu. 2023. Intuitive or dependent? investigating llms' robustness to conflicting prompts.	<i>ArXiv preprint</i> ,	Chujie Zheng, Jinfeng Zhou, Yinhe Zheng, Libiao Peng,	1705
1654		abs/2309.17415.	Zhen Guo, Wenquan Wu, Zhengyu Niu, Hua Wu,	1706
1655			and Minlie Huang. 2022. Cdconv: A benchmark for contradiction detection in chinese conversations.	1707
1656	Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang,		<i>ArXiv preprint</i> ,	1708
1657	Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu,		abs/2210.08511.	1709
1658	Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model.			
1659			Zexuan Zhong, Zhengxuan Wu, Christopher D Man-	1710
1660			ning, Christopher Potts, and Danqi Chen. 2023.	1711
1661	Boyuan Zhang, Hongyang Yang, Tianyu Zhou, Muham-		Mquake: Assessing knowledge editing in language models via multi-hop questions.	1712
1662	ad Ali Babar, and Xiao-Yang Liu. 2023a. Enhancing financial sentiment analysis via retrieval augmented large language models.		<i>ArXiv preprint</i> ,	1713
1663			abs/2305.14795.	1714
1664				
1665			Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao	1715
1666			Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu,	1716
1667			Lili Yu, et al. 2023a. Lima: Less is more for alignment.	1717
1668			<i>ArXiv preprint</i> ,	1718
1669			abs/2305.11206.	1719
1667	Hang Zhang, Xin Li, and Lidong Bing. 2023b. Video-llama: An instruction-tuned audio-visual language model for video understanding.		Hongjian Zhou, Boyang Gu, Xinyu Zou, Yiru Li,	1720
1668			Sam S Chen, Peilin Zhou, Junling Liu, Yining Hua,	1721
1669			Chengfeng Mao, Xian Wu, et al. 2023b. A survey of large language models in medicine: Progress, application, and challenge.	1722
			<i>ArXiv preprint</i> ,	1723
			abs/2311.05112.	

1724 Jiawei Zhou, Yixuan Zhang, Qianni Luo, Andrea G
1725 Parker, and Munmun De Choudhury. 2023c. Syn-
1726 thetic lies: Understanding ai-generated misinforma-
1727 tion and evaluating algorithmic and human solutions.
1728 In *Proceedings of the 2023 CHI Conference on Hu-
1729 man Factors in Computing Systems*, pages 1–20.

1730 Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and
1731 Muhao Chen. 2023d. [Context-faithful prompt-](#)
1732 [ing for large language models](#). *ArXiv preprint*,
1733 abs/2303.11315.

1734 Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and
1735 Chao Zhang. 2023. [Toolqa: A dataset for llm ques-](#)
1736 [tion answering with external tools](#). *ArXiv preprint*,
1737 abs/2306.13304.

A Taxonomy of Knowledge Conflicts

Figure 3 outlines the taxonomy we used in organize this survey. To start with, we classify knowledge conflicts into three categories based on the sources: context-memory conflict (§ 2), inter-context conflict (§ 3), and intra-memory conflict (§ 4). Within each type of conflict, we sequentially present its causes, analysis of LLMs’ behaviors, and possible mitigation solutions. Each specific issue is further categorized according to its internal characteristics (e.g., solutions are categorized based on the characteristics of the strategies engaged).

B Datasets of Knowledge Conflicts

We list notable datasets employed in investigating the three types of knowledge conflict in Table 1. It is worth noting that for all context-memory datasets, extra attention should be paid to their applicability. This is because these datasets always need to be based on model-specific memories as a baseline when constructing conflicting knowledge. Obviously, this parameterized knowledge varies from model to model, greatly reducing the reusability of these datasets. Furthermore, the value of these datasets is further diminished by the existence of model variants from different *knowledge cutoff date* (e.g., OpenAI’s GPT-4 family of models). The parameterized knowledge varies from variant to variant due to different cutoff date.

C Detailed Solutions for Context-Memory Conflict

C.1 Faithful to Context

Fine-tuning. Li et al. (2022a) argue 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 FineTuning (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

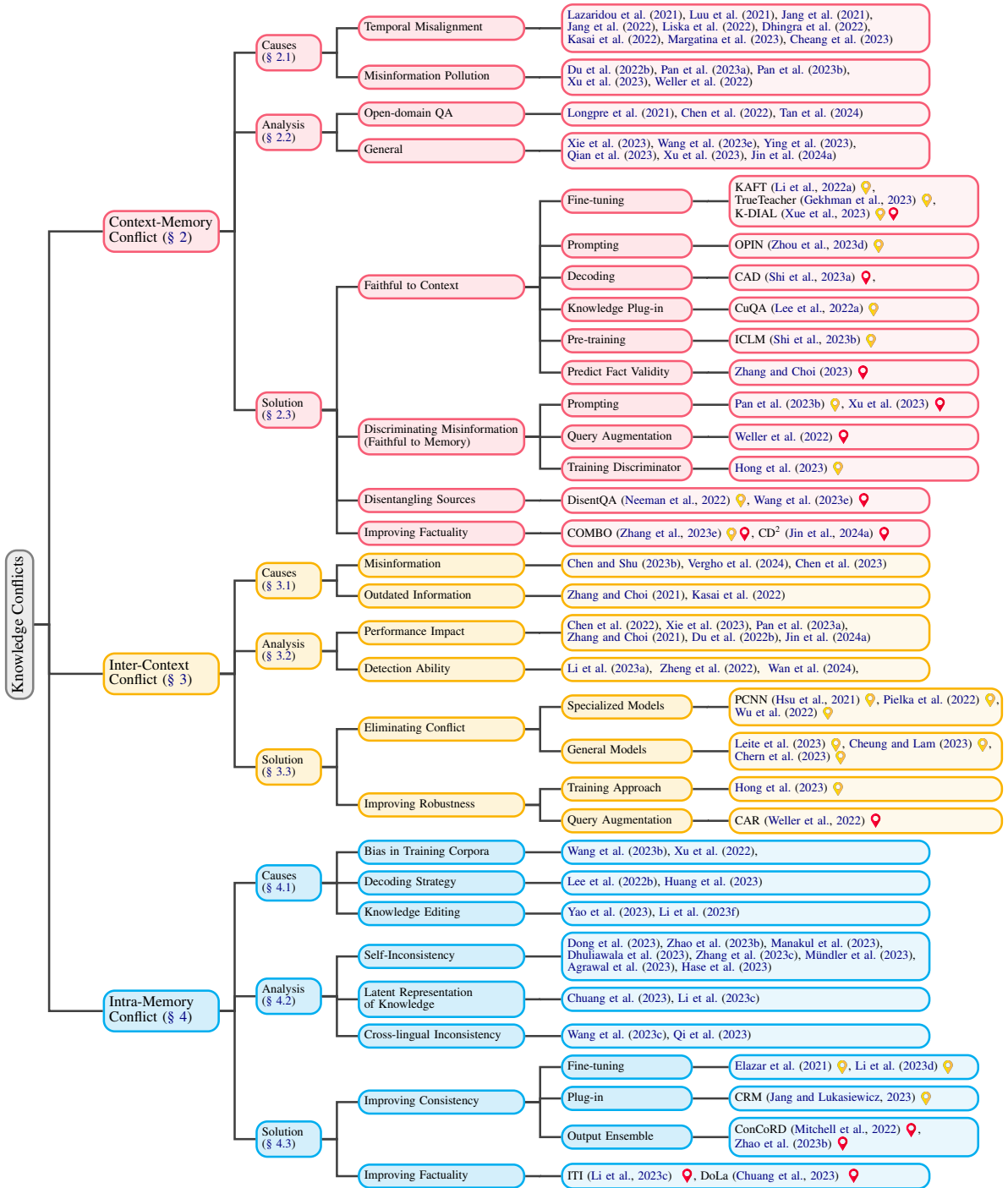


Figure 3: Taxonomy of knowledge conflicts. We mainly list works in the era of large language models. 🟡 denotes pre-hoc solution and 🔴 denotes post-hoc solution.

dialogue systems via direct knowledge enhancement and reinforcement learning for factual consistency (RLFC) for aligning responses accurately with provided factual knowledge.

Prompting. Zhou et al. (2023d) explores enhancing LLMs’ adherence to context through specialized prompting strategies, specifically opinion-based prompts and counterfactual demonstrations. These techniques are shown to significantly im-

prove LLMs’ performance in context-sensitive tasks by ensuring they remain faithful to relevant context, without additional training.

Decoding. Shi et al. (2023a) introduce Context-aware Decoding (CAD) to reduce hallucinations by amplifying the difference in output probabilities with and without context, which is similar to the concept of contrastive decoding (Li et al., 2022c). CAD enhances faithfulness in LLMs by effectively

1795
1796
1797
1798
1799
1800
1801
1802
1803

Dataset	Approach ¹	Base ²	Size	Conflict
Xie et al. (2023)	Gen	PopQA (2023), STRATEGYQA ((Geva et al., 2021))	20,091	CM ³
KC (2023e)	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
WikiContradiction (2021)	Hum	Wikipedia	2,210	IC
ClaimDiff (2022)	Hum	N/A	2,941	IC
Pan et al. (2023a)	Gen,Sub	SQuAD v1.1 (2016)	52,189	IC
CONTRADOC (2023a)	Gen	CNN-DailyMail (2015), NarrativeQA (2018), WikiText (2017)	449	IC
CONFLICTINGQA (2024)	Gen	N/A	238	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), and 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, one should select a subset of the dataset that aligns with the tested model’s knowledge when using **CM** datasets.

Table 1: Datasets on evaluating a large language model’s behavior when encountering knowledge conflicts. **CM**: context-memory conflict, **IC**: inter-context conflict, **IM**: intra-memory conflict.

prioritizing relevant context over the model’s prior knowledge, especially in tasks with conflicting information.

Knowledge Plug-in. Lee et al. (2022a) propose 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.

C.2 Discriminating Misinformation (Faithful to Memory)

Prompting. To address misinformation pollution, Pan et al. (2023b) propose 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) utilize 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) leverage 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 the cross-verification of answers in multiple sources.

Training Discriminator. Hong et al. (2023) fine-tune a smaller LM as a 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.

C.3 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. (2023e) introduce a method to improve 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.

1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877

1878
1879

1880

1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905

1906
1907
1908
1909
1910
1911
1912
1913

C.4 Improving Factuality

Zhang et al. (2023e) 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. **Jin et al. (2024a)** introduce a contrastive-decoding-based algorithm, namely CD², which maximizes the difference between various logits under knowledge conflicts and calibrates the model’s confidence in the truthful answer.

D Detailed Solutions for Inter-Context Conflict

D.1 Eliminating Conflict

Specialized Models. **Hsu et al. (2021)** develop a model named Pairwise Contradiction Neural Network (PCNN), leveraging fine-tuned Sentence-BERT embeddings to calculate contradiction probabilities of articles. **Pielka et al. (2022)** suggest incorporating linguistic knowledge into the learning process based on the discovery that XLM-RoBERTa struggles to effectively grasp the syntactic and semantic features that are vital for accurate contradiction detection. **Wu et al. (2022)** propose an innovative approach that integrates topological representations of text into language models to enhance the contradiction detection ability and evaluated their methods on the MultiNLI dataset (**Williams et al., 2018**).

General Models. **Chern et al. (2023)** propose a fact-checking framework that integrates LLMs with various tools, including Google Search, Google Scholar, code interpreters, and Python, for detecting factual errors in texts. **Leite et al. (2023)** employ LLMs to generate weak labels associated with predefined credibility signals for the input text and aggregate these labels through weak supervision techniques to make predictions regarding the veracity of the input.

D.2 Improving Robustness

Training Approach. **Hong et al. (2023)** present a novel fine-tuning method that involves training a discriminator and a decoder simultaneously using a shared encoder. Additionally, the authors introduce two other strategies 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. Their experimental results indicate that the fine-tuning method yields the most promising results.

Query Augmentation. **Weller et al. (2022)** explore a query augmentation technique that prompts GPT-3 to formulate new questions derived from the original inquiry. They then assess the confidence for each question’s answer by referencing the corresponding passages retrieved. Based on the confidence, they decide whether to rely on the original question’s prediction or aggregate predictions from the augmented questions with high confidence scores.

E Detailed Solutions for Intra-Memory Conflict

E.1 Improving Consistency

Fine-tuning. **Elazar et al. (2021)** propose a consistency loss function and train the language model with the combination of the consistency loss and standard MLM loss. **Li et al. (2023d)** utilize one language model in dual capacities: as a generator to produce responses and as a validator to evaluate the accuracy of these responses. The process involves querying the generator for a response, which is subsequently assessed by the validator for accuracy. Only those pairs of responses deemed consistent are retained. This subset of consistent pairs is then used to fine-tune the model, aiming to increase the generation likelihood of consistent response pairs.

Plug-in. **Jang and Lukasiewicz (2023)** leverage the technique of intermediate training, utilizing word-definition pairs from dictionaries to retrain language models and improve their comprehension of symbolic meanings. Subsequently, they propose an efficient parameter integration approach, which amalgamates these enhanced parameters with those of existing language models. This method aims to rectify the models’ inconsistent behavior by bolstering their capacity to understand meanings.

Output Ensemble. **Mitchell et al. (2022)** propose a method to mitigate the inconsistency of language models by leveraging a two-model architecture, involving the utilization of a base model responsible for generating a set of potential answers, followed by a relation model that evaluates the logical coherence among these answers. The final answer is selected by considering both the base model’s and the relation model’s beliefs. **Zhao et al. (2023b)** introduce a method to detect whether a question may

1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927

1928
1929

1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963

1964 cause inconsistency for LLMs. Specifically, they
1965 first use LLMs to rephrase the original question and
1966 obtain corresponding answers. They then cluster
1967 these answers and examine the divergence. The
1968 detection is determined based on the divergence
1969 level.

1970 E.2 Improving Factuality

1971 [Chuang et al. \(2023\)](#) propose a novel contrastive
1972 decoding approach named DoLa. Specifically, the
1973 authors develop a dynamic layer selection strategy,
1974 choosing the appropriate premature layers and ma-
1975 ture layers. The next word’s output probability is
1976 then determined by computing the difference in log
1977 probabilities of the premature layers and the mature
1978 layers. [Li et al. \(2023c\)](#) devise a similar method
1979 named ITI. They first identify a sparse set of atten-
1980 tion heads that exhibit high linear probing accuracy
1981 for truthfulness, as measured by TruthfulQA ([Lin
1982 et al., 2022](#)). During the inference phase, ITI
1983 shifts activations along the truth-correlated direc-
1984 tion, which is obtained through knowledge probing.
1985 This intervention is repeated autoregressively for
1986 every token during completion. Both DoLa and ITI
1987 address the inconsistency of knowledge across the
1988 model’s different layers to reduce factual errors.

1989 F Quantitative Analysis and Comparison

1990 In the context of a survey paper, while it is benefi-
1991 cial to include quantitative results and analyses con-
1992 cerning the impact of knowledge conflicts across
1993 various types of conflicts and the performance com-
1994 parison of different mitigation strategies, it is not a
1995 strict requirement. We acknowledge the *complex-
1996 ity and impracticality* involved in conducting such
1997 quantitative experiments, particularly due to the use
1998 of disparate datasets in behavioral analyses, as well
1999 as the variance in the inherent knowledge of LLMs
2000 across different knowledge cut-off snapshots, as
2001 detailed in § B.

2002 Moreover, establishing a “fair” comparison
2003 within the mitigation strategies segment poses its
2004 own set of challenges, given the diversity in objec-
2005 tives influenced by various assumed priors, such
2006 as the perceived accuracy of context or inherent
2007 knowledge, as discussed in the main text. De-
2008 spite these intricacies, we opt to present quantita-
2009 tive results by compiling existing evaluations from
2010 a range of papers. *It is imperative, however, to
2011 approach this analysis with caution, recognizing
2012 that original authors may have employed different*

*datasets, LLMs variants, or even pursued contrast-
ing objectives.*

2013 F.1 Quantitative Results on the Impact of 2016 Knowledge Conflicts

2017 The comparison of quantitative results on the im-
2018 pact of the three types of knowledge conflicts is
2019 shown in [Table 2](#). We pick the results of represen-
2020 tative behavior analysis literature for comparison.

2021 F.2 Quantitative Results on the Effectiveness 2022 of Mitigation Strategies

2023 The effectiveness of various mitigation strategies
2024 is quantitatively compared in [Table 3](#). It is impor-
2025 tant to note that our analysis is limited to works
2026 addressing *three predominant types of mitigating
2027 objectives* within the context of memory conflicts.
2028 This selection is deliberate, as other types of miti-
2029 gating objectives in different conflict categories do
2030 not yet have a substantial body of work that would
2031 allow for a meaningful cross-method comparison.

Reference	Model	Dataset	Quantitative Results
<i>Context-memory conflict</i>			
Pan et al. (2023b)	ChatGPT	NQ-1500 and CovidNews	Misinformation in the context can lead to a significant degradation (up to 87%) in the performance.
Xie et al. (2023)	ChatGPT, GPT-4, PaLM2, Qwen, Llama2, and Vicuna	POPQA and STRATEGYQA	For entity substitution-based counter-memory, only ChatGPT, GPT-4, and PaLM2 over 60% probability of choosing parametric memory. For generation-based counter-memory, all models have more than 80% probability of choosing context knowledge.
Xu et al. (2023)	ChatGPT, GPT-4, Llama2, and Vicuna	Farm, BoolQ, TruthfulQA and NQ	In multiple rounds of dialogue, as the number of counter-memory context increases, the cumulative proportion of belief alteration of LLMs spans from 20.7% to 78.2%
<i>Inter-context conflict</i>			
Jin et al. (2024a)	ChatGPT, Llama2, Baichuan2, FLAN-UL2 and FLAN-T5	NQ, TriviaQA, PopQA, and MuSiQue	When faced with conflicting evidence, ChatGPT’s recall declined the least, but more than 10%.
Chen et al. (2023)	ChatGPT, ChatGLM, Vicuna, Qwen, and BELLE	RGB	As the noise in evidence increases, the performance of models will gradually decrease. When the noise rate exceeds 0.8, the performance of all models decreases by more than 20%.
Li et al. (2023a)	GPT-4, ChatGPT, PaLM2, and Llama2	CONTRADOC	Faced with self-contradictory documents, gpt4 has a more than 70% probability of determining the occurrence of a contradiction, while other models are less than 50%.
<i>Intra-memory conflict</i>			
Mündler et al. (2023)	GPT-4, ChatGPT, Llama2, and Vicuna	MainTestSet	LLMs create contradictory content, with a probability of between 15.7% and 22.9%. More powerful models create fewer contradictory results.
Zhao et al. (2023b)	ChatGPT, GPT-4, Vicuna, and Llama2	FaVIQ, ComQA, GSM-8K, SVAMP, ARCChallenge, and CommonsenseQA	The findings of their research reveal that even GPT-4 can exhibit an inconsistency rate of 32% in FaVIQ.

Table 2: Comparison of quantitative results on the impact of various types of knowledge conflicts.

Reference	Model	Dataset	Quantitative Results
<i>Faithful to context</i>			
Shi et al. (2023a)	Llama, OPT, GPT-Neo, and FLAN	NQ-SWAP, MemoTrap, and NQ	Their method improves GPT-Neo 20B by 54.4% on Memotrap and by 128% on NQ-SWAP where LLMs need to adhere to the given context.
Zhou et al. (2023d)	ChatGPT and Llama2	MRC and Re-TACRED	Compared to the zero-shot base prompts, their prompting method leads to a reduction of 32.2% for maintaining parametric knowledge for MRC and a 10.9% reduction for Re-TACRED on GPT-3.5. Similarly, on Llama2, there is a 39.4% reduction for MRC and a 57.3% reduction for Re-TACRED.
<i>Discriminating misinformation</i>			
Hong et al. (2023)	ChatGPT and FiD	NQ and TQA	The authors train a discriminator with about 80% F1 score and use it to improve models performance above 5%.
Pan et al. (2023b)	ChatGPT	NQ-1500 and CovidNews	The author’s mitigation method improves the accuracy by more than 10%.
<i>Disentangling sources</i>			
Wang et al. (2023e)	ChatGPT	KNOWLEDGE CONFLICT	The authors’ method achieved over 80% F1 score on contextual knowledge conflict detection.

Table 3: Comparison of quantitative results on the effectiveness of various mitigation strategies *w.r.t.* their objectives.