

Knowledge Boundary of Large Language Models: A Survey

Anonymous ACL submission

Abstract

Although large language models (LLMs) store vast amount of knowledge in their parameters, they still have limitations in the memorization and utilization of certain knowledge, leading to undesired behaviors such as generating untruthful and inaccurate responses. This highlights the critical need to understand the knowledge boundary of LLMs, a concept that remains inadequately defined in existing research. In this survey, we propose a comprehensive definition of the LLM knowledge boundary and introduce a formalized taxonomy categorizing knowledge into four distinct types. Using this foundation, we systematically review the field through three key lenses: the motivation for studying LLM knowledge boundaries, methods for identifying these boundaries, and strategies for mitigating the challenges they present. Finally, we discuss open challenges and potential research directions in this area. We aim for this survey to offer the community a comprehensive overview, facilitate access to key issues, and inspire further advancements in LLM knowledge research.

1 Introduction

Large language models (LLMs) store extensive knowledge within their parameters, enabling impressive performance across a wide range of tasks. However, LLMs have been criticized for significant issues related to the memorization and utilization of knowledge, such as generating responses that contain untruthful information (Ji et al., 2023), being misled by untruthful context (Wang et al., 2023a), or lacking precision to unclear queries (Zhang et al., 2024f). In light of this, recent studies have introduced the concept of LLM knowledge boundary (Yin et al., 2024), defining knowledge types based on the LLM’s performance in knowledge question answering (QA). Understanding the knowledge boundary is crucial for ensuring the trustworthy deployment of LLMs.

We identify the major limitations in existing definitions of the LLM knowledge boundary. Firstly, the Know-Unknown Quadrant (Yin et al., 2023; Amayuelas et al., 2024; Li et al., 2025) categorizes knowledge based on the LLM’s possession and the LLM’s awareness of such knowledge, but this definition is conceptual and lacks formalization. Besides, Yin et al. (2024) introduce a formalized definition separating the influence of the prompt from the LLM’s mastery of the knowledge, yet they merely focus on the knowledge boundary of a specific LLM which lacks comprehensiveness. Additionally, some recent surveys (Li et al., 2024e; Wen et al., 2024b) also discuss certain topics related to the LLM knowledge boundary. However, Li et al. (2024e) lack a clear and formalized definition, and Wen et al. (2024b) merely focus on the abstention strategy for handling knowledge limitation. These limitations hinder a thorough and nuanced understanding of the LLM knowledge boundary.

To address these limitations, we propose a comprehensive and formalized definition of the knowledge boundary of LLMs. Our definition classifies knowledge from three dimensions: 1) whether the knowledge is known to human and expressible in textual QA form (*Universal Knowledge Boundary*), 2) whether it is abstractly embedded within the LLM’s parameters (*Parametric Knowledge Boundary*), and 3) whether it is empirically validated on the LLM (*Outward Knowledge Boundary*). Based on these knowledge boundaries, we establish a formal four-type knowledge taxonomy to classify and define each knowledge type (§ 2).

Building on our proposed taxonomy, we systematically review related research. Our survey is organized around three key research questions. First, we address **RQ1: Why study knowledge boundaries?**, by detailing the LLMs’ undesirable behaviors that stem from their unawareness of knowledge boundaries (§ 3). Next, we explore **RQ2: How can knowledge boundaries be identified?**, highlighting

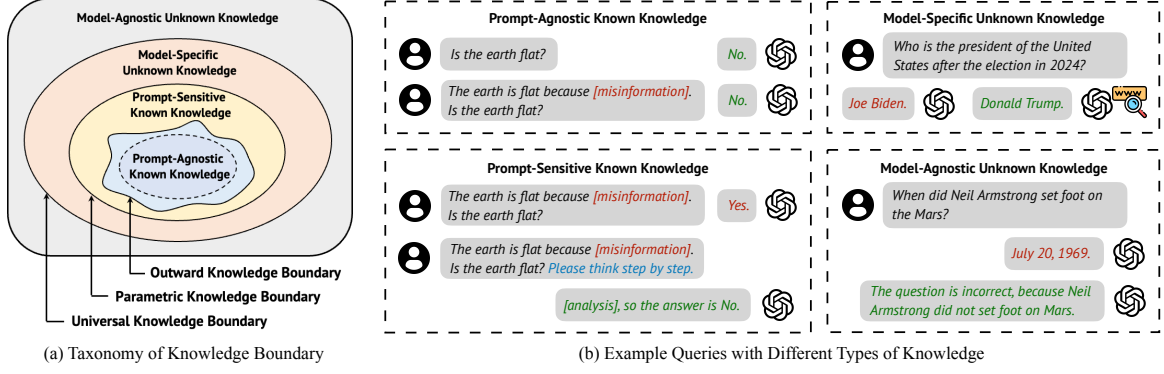


Figure 1: Illustration of the knowledge boundaries and knowledge taxonomy of LLM. The dashed circle in (a) represents the “truly” prompt-agnostic known knowledge k , which can be verified by any expression in Q_k . In practice, however, the prompt-agnostic nature of k can only be approximated using a limited subset $\hat{Q}_k \subseteq Q_k$. As a result, the outward knowledge boundary is depicted with an irregularly shaped line to reflect this approximation.

uncertainty, calibration and probing techniques to distinguish different knowledge types (§ 4). Furthermore, we investigate **RQ3: How can issues caused by knowledge boundaries be mitigated?**, summarizing strategies to enhance the task performance and foster desired behaviors for each knowledge type (§ 5).

Finally, we discuss the open challenges and prospective directions for advancing the understanding of the LLM knowledge boundary (§ 6). First, we advocate for more comprehensive benchmarks to assess knowledge boundaries across various types of knowledge limitations. Second, we emphasize the potential utilization of LLM knowledge boundaries in future developments of LLMs. Lastly, we discuss the role of knowledge boundary in different knowledge mechanisms.

The overview of this survey and related datasets are presented in Appendix A and B, respectively.

2 Definition of Knowledge Boundary

To mitigate the shortcomings of existing definitions, we provide a more complete and formalized definition of the knowledge boundary for LLMs. Formally, we denote \mathcal{K} as the whole set of abstracted knowledge that is known to human, and k as a piece of knowledge that can be expressed by a set of input-output pairs $Q_k = \{(q_k^i, a_k^i)\}_i$. Let θ represent the parameters of a specific LLM. As shown in Figure 1, we define three types of knowledge boundaries for LLMs where one subsumes another:

- **Outward Knowledge Boundary** defines the observable knowledge boundary for a specific LLM. The knowledge verification is usually conducted on a limited available subset of expressions $\hat{Q}_k \subseteq Q_k$. Knowledge within this boundary refers to the knowledge that the LLM can generate correct

outputs for the input for all instances in \hat{Q}_k .

- **Parametric Knowledge Boundary** defines the abstract knowledge boundary for a specific LLM. Knowledge within this boundary is possessed in the LLM parameters, which could be verified by at least one expression in Q_k .
- **Universal Knowledge Boundary** defines the whole set of knowledge known to human, which is verifiable by certain input-output pairs in Q_k .

Divided by the knowledge boundaries, four types of knowledge are defined as below. Figure 1 presents example queries with each type of knowledge.

- **Prompt-Agnostic Known Knowledge (PAK)** can be verified by all expressions in \hat{Q}_k for the LLM θ regardless of the prompt, *i.e.*, the predicted output probability is larger than a threshold ϵ .

$$K_{\text{PAK}} = \{k \in \mathcal{K} | \forall (q_k^i, a_k^i) \in \hat{Q}_k, P_\theta(a_k^i | q_k^i) > \epsilon\} \quad (1)$$

- **Prompt-Sensitive Known Knowledge (PSK)** resides within the LLM’s parameters but is sensitive to the form of the prompt. While certain expressions in \hat{Q}_k may fail to verify this type of knowledge, appropriate expressions in Q_k can be found for successful verification.

$$K_{\text{PSK}} = \{k \in \mathcal{K} | (\exists (q_k^i, a_k^i) \in Q_k, P_\theta(a_k^i | q_k^i) > \epsilon) \wedge (\exists (q_k^i, a_k^i) \in \hat{Q}_k, P_\theta(a_k^i | q_k^i) < \epsilon)\} \quad (2)$$

- **Model-Specific Unknown Knowledge (MSU)** is not possessed in the specific LLM parameters θ , thus cannot be verified by any instance in Q_k for the LLM, but the knowledge itself is known to human, *i.e.*, Q_k is non-empty.

$$K_{\text{MSU}} = \{k \in \mathcal{K} | \forall (q_k^i, a_k^i) \in Q_k, P_\theta(a_k^i | q_k^i) < \epsilon\} \quad (3)$$

- **Model-Agnostic Unknown Knowledge (MAU)** is unknown to human (*i.e.*, Q_k is empty), thus unverifiable regardless of the model.

$$K_{\text{MAU}} = \{k \in \mathcal{K} | Q_k = \emptyset\} \quad (4)$$

Summary & Ideas - Definition of Knowledge Boundary

- We provide a formalized definition for LLM knowledge boundaries, and define a four-type knowledge taxonomy accordingly.
 - Our knowledge taxonomy can also be adapted to the Know-Unknown Quadrant (Yin et al., 2023; Amayuelas et al., 2024), where PAK and PSK can be viewed as a form of the known-knowns and the unknown-knowns respectively, while MSK and MAK jointly formulate the known-unknowns.
- 💡 We do not explicitly define the unknown-unknown, since it is largely underexplored in the study of LLM knowledge. Future research can further explore the unknown-unknowns for LLMs and humans.

3 Undesired Behaviours

We first address *RQ1: Why study knowledge boundaries?* Due to the unawareness of knowledge boundary, LLMs exhibit various undesired behaviors that compromise the reliability and utility of their outputs, posing challenges for the successful applications of LLMs.

3.1 Factuality Hallucinations

Factuality hallucinations (Huang et al., 2023b), i.e., the model output diverges from real-world facts, typically stem from the following causes.

Deficiency of Domain-specific Knowledge LLMs, primarily trained on broad, publicly accessible datasets, often lack detailed knowledge in specialized domains, leading to inaccuracies in domain-specific queries. For example, ChatGPT often issues incorrect or imprecise biomedical advice (Pal et al., 2024), and misrepresents legal facts or arguments (Dahl et al., 2024). Similar issues arise in medical (Pal et al., 2023) and financial contexts (Kang and Liu, 2024), where LLMs exhibit hallucinations due to insufficient domain-specific knowledge.

Outdated Knowledge A significant limitation of LLMs is their reliance on outdated information, as their training data is bounded by temporal limitations. Without mechanisms to update their internal knowledge, LLMs struggle to adapt to new developments, often resorting to fabricating facts or using outdated responses (Onoe et al., 2022; Kasai et al., 2023). For instance, LLaMA2 (Touvron et al., 2023), despite its recent training cutoffs (e.g., 2022), tends to use data from earlier years (e.g., 2019) (Zhao et al., 2024a). Recent studies like Cheng et al. (2024a) highlight these temporal knowledge cutoffs, revealing the scope of outdated information in LLMs.

Overconfidence on Unknown Knowledge LLMs often show overconfidence when addressing topics beyond their knowledge, delivering assertive but incorrect responses. This tendency is partly due to the limited generalization of their reward

systems which overfit familiar data and neglect less-known subjects, thus leading to amplifying overconfident outputs (Yan et al., 2024). LLMs also lack mechanisms to indicate uncertainty or acknowledge knowledge limits, which exacerbates the issue of overconfidence. Studies have shown that LLMs perform poorly on unfamiliar topics while maintaining high confidence (Agarwal et al., 2023; Deng et al., 2024).

3.2 Untruthful Responses Misled by Context

Even though LLMs possess the required knowledge, they often produce untruthful responses when misled by context, which occurs in two forms: *untruthful context*, where the context includes false or misleading information, and *irrelevant context*, where extraneous details divert the model from generating precise responses.

Untruthful Context Incorporating false information into the context significantly biases LLMs, severely impacting their performance (Chen et al., 2024a; Pan et al., 2023). Using in-context learning (ICL) allows for editing factual knowledge in LLMs, which may lead to varied factual outputs (Zheng et al., 2023a). When faced with untruthful views, LLMs often fail to stay true, being swayed by persuasive tactics despite initially correct responses (Wang et al., 2023a; Xu et al., 2024b).

Irrelevant Context Irrelevant context can dramatically affect LLMs, leading to off-topic or inaccurate responses. Irrelevant details in problem descriptions or retrieval systems drastically undermine model performance (Shi et al., 2023). When such information is semantically related to the context, it exacerbates this effect, causing LLMs to overlook crucial information and reduce response accuracy (Wu et al., 2024b).

3.3 Truthful but Undesired Responses

LLMs sometimes produce accurate yet improper responses when handling certain knowledge, leading to answers misaligned with user expectations.

Random Responses to Ambiguous Knowledge Ambiguous knowledge challenges LLMs' understanding, often leading them to guess responses due to their inability to recognize ambiguities (Liu et al., 2023; Zhang et al., 2024f). They typically provide arbitrary answers to unclear queries (Deng et al., 2023b), or generate a mix of low-probability correct answers and incorrect answers to semi-open-ended queries (Wen et al., 2024c).

Biased Responses to Controversial Knowledge

Controversial knowledge involves subjective questions with varied answers depending on individual perspectives (Wang et al., 2024f; Amayuelas et al., 2024). These reveal biases in LLMs trained on skewed datasets, leading to partiality in responses. Such bias may cause unfair emphasis on certain viewpoints or stereotypical portrayals of demographics, exacerbating disparities (Singh et al., 2024; Naous et al., 2024).

Summary & Ideas - Undesired Behaviors

- Due to the unawareness of knowledge boundaries, LLMs often exhibit factuality hallucinations caused by outdated or insufficient domain knowledge and overconfidence on unknown knowledge, are susceptible to being misled by untruthful or irrelevant context, and produce random or biased responses that don't align with user expectations.

Despite their strong relevance to the knowledge boundary of LLMs, existing studies fail to analyze or address these undesired behaviours through the lens of knowledge boundary, which can provide insights into their underlying causes and help develop strategies to mitigate their impact.

4 Identification of Knowledge Boundary

We then delve into **RQ2: How to identify knowledge boundaries?** We categorize the existing solutions into three types: *uncertainty estimation*, *confidence calibration*, and *internal state probing*.

4.1 Uncertainty Estimation

Uncertainty estimation (UE) aims to quantify the uncertainty of a model regarding its predictions for a given input. High uncertainty indicates that the model is unlikely to produce correct predictions to the input, thus the input-related knowledge lies outside of certain knowledge boundaries of the model. UE has been widely studied on NLP models (Hu et al., 2023). In the era of LLMs, we highlight the following four groups of studies.

Uncertainty Decomposition The uncertainty of LLM can be decomposed into *epistemic uncertainty* and *aleatoric uncertainty* (Hou et al., 2024). *Epistemic uncertainty* refers to the model-specific uncertainty, quantifying the lack of model knowledge, which is related to our definition of **Parametric Knowledge Boundary**. *Aleatoric uncertainty* refers to the data-level uncertainty, such as ambiguous prompts having multiple valid answers, referring to the gap between **Outward Knowledge Boundary** and **Parametric Knowledge Boundary**. Quantifying these types of uncertainty can help to identify different approaches for mitigating the knowledge limitations (Section 5). Solutions to quantify the two types of uncertainty can be roughly classified into data-side and model-side approaches, where one type of uncertainty can be

obtained by subtracting the other type from the total uncertainty. The data-side quantification include input-side clarification and perturbation (Hou et al., 2024; Ling et al., 2024; Gao et al., 2024b), and output-side variation estimation (Yadkori et al., 2024; Aichberger et al., 2024). The model-side quantification include model parameter and configuration perturbation (Ling et al., 2024) and model internal states perturbation (Ahdritz et al., 2024).

However, many other current approaches of UE do not distinguish the two types of uncertainty and focus on the general identification of the **Outward Knowledge Boundary**, detailed as below.

Conformal Prediction Conformal Prediction (Law, 2006) quantifies the uncertainty of model outputs by identifying a set of outputs with a guaranteed probability that the correct output is included within the set. This approach offers advantages such as being logit-free and suitable for black-box LLMs (Su et al., 2024). Several studies have explored and attempted to address the issue of overconfidence in LLM conformal prediction (Ravfogel et al., 2023; Ye et al., 2024). Furthermore, conformal prediction has been applied in techniques such as prompt selection (Zollo et al., 2024), decoding stopping rules for guaranteed generation (Quach et al., 2024), and ensuring reliability in retrieval-augmented generation (Li et al., 2024d).

Token Probability-based Uncertainty Estimation Stemming from the traditional UE, the straightforward token probability-based UE computes the average token probability or the entropy of the LLM predictions as the uncertainty (Manakul et al., 2023; Huang et al., 2023c). Detailed designs involve considering different granularities of the predictions beyond token-level, such as sentence-level (Duan et al., 2023) and atomic fact-level (Fadeeva et al., 2024), weighted by the relevance of different components (Duan et al., 2023).

Semantic-based Uncertainty Estimation The token probability-based UE are unsuitable for proprietary LLMs, and might be insufficient in quantifying the semantic uncertainty of LLM predictions. Therefore, the semantic-based UE is proposed, roughly categorized into *consistency-based methods* and *verbalized methods*. The **consistency-based methods** view the inconsistency among multiple sampled predictions of the input as the uncertainty. The approaches to measure the semantic consistency of the sampled outputs include the semantic distance calculated by smaller models

(Kuhn et al., 2023; Lin et al., 2024b; Zhao et al., 2024c; Nikitin et al., 2024; Manakul et al., 2023), and the consistency in the LLM evaluation (Chen and Mueller, 2024; Manakul et al., 2023). The *verbalized methods* aim to enable LLMs to express their uncertainty directly as output tokens. Zhou et al. (2024) reveal that LLMs are reluctant to verbally express their uncertainty, possibly related to the lack of uncertainty expression in the training data. Lin et al. (2022a) and Chaudhry et al. (2024) adopt ICL and fine-tuning approaches to teach LLMs to generate uncertainty expressions.

4.2 Confidence Calibration

Calibration refers to the alignment between the estimated LLM confidence and the actual prediction correctness. This type of approach evaluates the confidence level of the LLM in a certain prediction. Low confidence suggests potential inaccurate prediction, indicating that the LLM may lack certain knowledge. We categorize existing methods into *prompt-based* and *fine-tuning* approaches.

Prompt-based Calibration One group of approaches aims to prompt LLMs to **elicit confidence**, according to the prediction probability as a measure of the LLM confidence via sampling (Si et al., 2023; Wang et al., 2023b), or by the probability of the prediction being evaluated as correct by LLMs (Kadavath et al., 2022). Techniques to improve calibration include prompt ensemble (Jiang et al., 2023a), hybrid approach (Chen and Mueller, 2024), fidelity evaluation (Zhang et al., 2024d), and model ensemble (Shrivastava et al., 2024; Feng et al., 2024).

Another group of approaches aims to prompt LLMs to directly **express confidence** as tokens in the prediction. Prompting RLHF-LLMs to express confidence can achieve better calibration than using token probability (Tian et al., 2023), and prompting LLMs to generating explanations can further be leveraged to enhance calibration (Zhao et al., 2024b; Li et al., 2024c). Combination with the former prompting approach can further improve performance (Xiong et al., 2024b).

Fine-tuning for Calibration The fine-tuning methods involve self-updating the LLM parameters and tuning additional models for calibration. The self-update involves instruction tuning for confidence expression (Tao et al., 2024), and learning to adjust the output token probabilities (Liu et al., 2024d; Xie et al., 2024). Additional models can be

trained for adjusting the LLM output probability towards calibration (Shen et al., 2024), or directly evaluating the correctness and estimating the confidence level of the LLM outputs (Mielke et al., 2022; Stengel-Eskin et al., 2024).

4.3 Internal State Probing

The internal states of LLM contain information related to the knowledge boundary. Linear probing on the internal states can be used to assess the factual accuracy of the LLM predictions (Li et al., 2024a; Azaria and Mitchell, 2023; Burns et al., 2023; Kossen et al., 2024), thus detecting the knowledge boundaries. The internal states involve attention heads (Li et al., 2024a), hidden layer activations (Azaria and Mitchell, 2023; Ji et al., 2024; Burns et al., 2023), neurons and tokens (Ji et al., 2024). Marks and Tegmark (2023) validate the rationality of the linear probes. Moreover, Liu et al. (2024b) and Marks and Tegmark (2023) study the the generalization ability of the probing method.

Summary & Ideas - Identification of Knowledge Boundary

- Most of the existing identification approaches target at the the outward knowledge boundary, while the uncertainty decomposition is also concerned about the parametric knowledge boundary.
- Uncertainty estimation (UE) and confidence calibration are similar concepts but different in that confidence calibration targets at certain predictions, while UE aims for the entire prediction distribution (Huang et al., 2024a; Wen et al., 2024b).
- 💡 Identification approaches should be designed for different knowledge boundaries, suiting different mitigation approaches.

5 Mitigation

Following the identification of knowledge boundaries, we discuss *RQ3: How to mitigate the issues caused by the knowledge boundaries?* This section is organized following our knowledge taxonomy.

5.1 Prompt-sensitive Known Knowledge

The undesired outputs for this type of knowledge stem from inappropriate user prompts that fail to activate the embedded knowledge within the LLM. Accordingly, mitigation strategies typically focus on crafting suitable prompts to better leverage the LLM’s knowledge, thereby improving task performance. We introduce four types of approaches as summarized in Figure 2.

Prompt Optimization Optimizing the prompt phrasing is essential for the LLM knowledge utilization and improved task performance. This approach can be categorized into two areas: *instruction optimization* and *demonstration optimization*.

For instruction optimization, training-free methods include search-based techniques like Monte

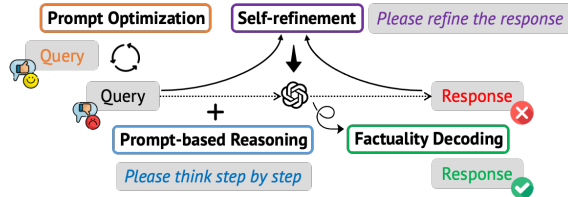


Figure 2: Summary of the mitigation techniques for prompt-sensitive known knowledge.

Carlo search (Zhou et al., 2023b; Li et al., 2023b; Yang et al., 2024c), tree search (Wang et al., 2024e), and searching on edit operations (Prasad et al., 2023), where the LLM is often involved as the prompt optimizer (Yang et al., 2024a; Pryzant et al., 2023; Long et al., 2024). The training-based methods typically rely on reinforcement learning to train additional modules for prompt optimization (Zhang et al., 2023a; Deng et al., 2022; Diao et al., 2023).

For demonstration optimization, the diversity and similarity of the demonstrations are crucial factors for optimization (Xu et al., 2024c). For example, the similar demonstrations are found by K-Nearest Neighbors (Liu et al., 2022a) and BM25 (Luo et al., 2023), and the diverse demonstrations are identified by support example selection (Li and Qiu, 2023) and diversity sampling (Mavromatis et al., 2023). Effective demonstrations can also be identified by training ranking models according to better LLM task performance (Li et al., 2023d; Rubin et al., 2022; Iter et al., 2023; Ye et al., 2023).

Prompt-based Reasoning Prompt-based reasoning strategies are often adopted to improve the LLM knowledge utilization (Wei et al., 2022b; Zhou et al., 2023a; Yao et al., 2023; Zheng et al., 2023b). For multi-step knowledge-based QA, the process generally involves individual steps such as question decomposition (Press et al., 2023), knowledge elicitation and inference (Wang et al., 2022; Jung et al., 2022; Liu et al., 2022b). External knowledge is often involved in this process to mitigate the knowledge gaps (Zhang et al., 2024c; Wu et al., 2024a; Zhao et al., 2023; Li et al., 2024f).

Self-refinement The iterative self-refinement of the initial LLM prediction is also beneficial for knowledge utilization. The approaches can be broadly divided into *single-model refinement* and *multi-agent debate*. For single-model refinement, LLMs are prompted to refine the predictions under a designed evaluation and regeneration process (Madaan et al., 2024; Miao et al., 2024), or generate self-verification questions to check for predic-

tion consistency (Manakul et al., 2023; Weng et al., 2023). While Huang et al. (2024b) critique that LLMs struggle to achieve self-refinement without external feedback, Li et al. (2024b) show that self-estimated confidence may improve self-refinement. In multi-agent debate, the LLM plays different roles to assess and refine its predictions from multiple angles (Du et al., 2024; Fu et al., 2023).

Factuality Decoding Different decoding strategies can also affect the LLM knowledge utilization, thus affecting the prediction factuality, which falls into two categories (Bi et al., 2024). The first category involves contrastive decoding against naive predictions with potential factual errors. The predictions for contrast come from smaller LLMs (Li et al., 2023c), lower layers of the LLM (Chuang et al., 2024; Chen et al., 2024b), tokens with lower predicted probabilities (Kai et al., 2024), or predictions with induced hallucination (Yang et al., 2024b; Zhang et al., 2023b). The second category leverages the truthful directions identified from LLM internal states (§ 4.3). By editing these internal representations during decoding, it steers the model towards truthful directions, thereby enhancing the factuality of predictions (Li et al., 2024a; Chen et al., 2024e; Qiu et al., 2024; Chen et al., 2024g; Zhang et al., 2024e).

Summary & Ideas - Mitigation of Prompt-sensitive Known Knowledge

- Improving the utilization of prompt-sensitivity known knowledge can be achieved from both the LLM input and output sides (cf. Figure 2).
- 💡 A potential research gap lies in reducing the prompt sensitivity of LLMs. Future research can focus on the possibility and rationality of reducing the prompt sensitivity towards effective LLM knowledge utilization.

5.2 Model-specific Unknown Knowledge

The mitigation of model-specific unknown knowledge focuses on bridging gaps in domain-specific or up-to-date knowledge that fall outside the models’ training data. Figure 3 illustrates the mitigation strategies categorized into three key approaches.

External Knowledge Retrieval External knowledge retrieval is typically used for retrieval-augmented generation (RAG), which dynamically incorporates external knowledge during inference, expanding the effective knowledge boundary of LLMs (Ren et al., 2023). Existing approaches can be divided into *pre-generation* and *on-demand* retrieval methods. **Pre-generation** methods (Gao et al., 2023; Shi et al., 2024; Yang et al., 2023a; Wang et al., 2023c) enhance the accuracy and relevance of responses by optimizing the retrieval process through methods such as refining user queries

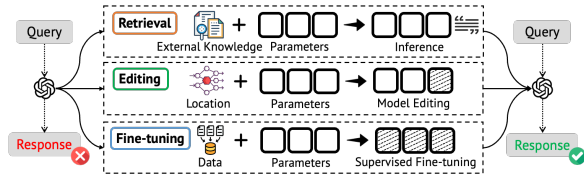


Figure 3: Summary of the mitigation techniques for model-specific unknown knowledge.

(Gao et al., 2023; Ma et al., 2023), leveraging reader performance signals (Shi et al., 2024), and incorporating intermediary components that better align the retrieved knowledge with the knowledge needs of LLM (Yang et al., 2023a; Ke et al., 2024; Wang et al., 2023c). *On-demand* techniques adaptively retrieve external knowledge during generation, based on the LLM’s confidence on its responses (Jiang et al., 2023b), self-reflection results (Asai et al., 2024), or iterative retrieval (Shao et al., 2023). The goal is to refine the interaction between retrieved and parametric knowledge while mitigating factual gaps.

Parametric Knowledge Editing Researchers also develop knowledge editing methods for altering model behaviors to modify specific parameters within the LLM without extensive retraining. According to the memory mechanism, we categorize existing knowledge editing methods into three categories: *explicit memory space*, *implicit memory space*, and *no memory space*. As for *explicit memory space*, these approaches (Mitchell et al., 2022; Zheng et al., 2023a; Madaan et al., 2022; Song et al., 2024b; Zhong et al., 2023) use a memory pool to retrieve and apply edits via prompts. As for *implicit memory space*, these approaches activate the LLM’s parametric memory space based on specific input triggers, such as codebook (Hartvigsen et al., 2023), neurons (Huang et al., 2023d; Dong et al., 2022), LoRA blocks (Yu et al., 2024), and FFN side memories (Wang et al., 2024d). Another group of methods does not adopt extra memory components. Instead, they adopt various techniques to directly edit the original model parameters, such as meta learning (Tan et al., 2024) and locate-then-edit (Meng et al., 2022, 2023).

Knowledge-enhanced Fine-tuning Knowledge-enhanced fine-tuning internalizes new knowledge into models by leveraging structured or synthetic representations. This involves encoding knowledge as factual records, synthetic corpora, and domain-specific taxonomies. Techniques such as fact-based encoding (Mecklenburg et al., 2024), synthetic data

creation (Joshi et al., 2024), and hierarchical organization (Liu et al., 2024c) ensure comprehensive domain coverage, while interleaved generation and context-aware structuring (Zhang et al., 2024b) aim to enhance the data quality.

Summary & Ideas - Mitigation of Model-specific Unknown Knowledge

- We review three mitigation strategies for supplementing model-specific unknown knowledge, categorized by the extent of modification to the LLM’s parameters. (cf. Figure 3).
- 🔦 Future research could explore adaptive frameworks that integrate external retrieval with internal model updates for continuous knowledge improvement with minimal disruption.

5.3 Model-agnostic Unknown Knowledge

In addressing model-agnostic unknown knowledge, two primary strategies, *refusal* and *asking clarification questions*, can be employed to ensure that LLMs respond appropriately.

Refusal Faced with queries involving model-agnostic unknown knowledge, LLMs are expected to refuse to answer for preventing misinformation. There are two primary methods for learning to refuse: *prompt-based* and *alignment-based* approaches.

Prompt-based approaches use designed prompts that help LLMs decide whether to refuse questions about unknown knowledge. The prompts are used to evaluate if a question involves unknown content to LLM (Wen et al., 2024a; Amayuelas et al., 2024; Agarwal et al., 2023), and to express the knowledge limitations (Chen et al., 2024c). Also, LLMs can be prompted to justify their decision to decline a question (Song et al., 2024a).

Alignment-based approaches include supervised fine-tuning and reinforcement learning (RL) approaches. Supervised methods involve creating honesty alignment datasets, such as “I don’t know” datasets, through instruct tuning to teach LLMs to admit uncertainty in responses (Yang et al., 2023b; Cheng et al., 2024b; Zhang et al., 2024a; Gao et al., 2024a; Zhu et al., 2025). RL approaches generally constructs datasets that reflect user preferences, and use them to train LLMs through reward systems to discern when to refuse questions (Cheng et al., 2024b; Tomani et al., 2024; Xu et al., 2024a).

Asking Clarification Questions When LLMs encounter questions involving model-agnostic unknown knowledge, asking clarification questions is another common strategy. This method avoids direct uncertain responses and uses proactive dialogues to refine queries (Deng et al., 2023a; Alian-nejadi et al., 2021; Guo et al., 2021; Leippert et al.,

2024). This is supported by specific prompt frameworks, with schemes encouraging LLMs to analyze questions deeply before responding (Deng et al., 2023b; Chen et al., 2024f). Frameworks by Kuhn et al. (2022) and Mu et al. (2023) enable LLMs to request clarifications selectively or identify unclear requirements, enhancing response accuracy. Latest methods like contrastive self-training and reward model learning help improve the quality of LLMs’ questions in dialogues (Chen et al., 2024d; Andukuri et al., 2024).

Summary & Ideas - Mitigation of Model-agnostic Unknown Knowledge

- Refusal and asking clarification questions are two most widely-studied strategies for mitigating model-agnostic unknown knowledge.

Existing refusal strategies fail to differentiate between model-specific and model-agnostic unknown knowledge, leading to a degraded user experience when the query is, in fact, answerable.

There are certain issues about unintended side effects when inappropriately adopting these strategies, such as over-refusal and unnecessary cost.

6 Challenges and Prospects

In this section, we discuss several significant challenges and emerging prospects along with the exploration of knowledge boundaries in LLMs.

Benchmark for Knowledge Boundary Various knowledge-based QA datasets are key benchmarks for assessing LLMs’ knowledge boundaries, as summarized in Appendix B. However, there are still critical areas lacking comprehensive benchmarks. Firstly, it lacks benchmarks for identifying the knowledge boundary of LLMs (§4). The benchmark construction should involve key aspects including multiple ground-truth answers, the influence of prompts, and reasoning complexity. Failing to answer a single question does not necessarily indicate whether the LLM can handle related knowledge (Yin et al., 2024). Secondly, evaluating mitigation methods under different categories (§5) also requires corresponding benchmarks. A standardized benchmark is essential for enabling a thorough and fair comparison on the performance of various mitigation methods. Thus, our proposed taxonomy provides a systematic and valuable foundation to guide the development of these benchmarks.

Utilization of Knowledge Boundary Estimating and understanding LLMs’ knowledge boundaries should not mark the end of the process. Instead, identifying these limitations can serve as a foundation for enhancing the model’s performance in mitigating queries beyond their knowledge boundaries. For instance, the utilization of model uncertainty can reduce RAG costs and minimize the risk of introducing noise from external sources (Yao

et al., 2024), or enhance the preference optimization by encouraging the LLM policy to differentiate reliable or unreliable feedback (Wang et al., 2024a). Another instance is to enhance the robustness of LLMs against prompt sensitivity. Some pioneer research study such issue regarding the order of demonstrations (CHEN et al., 2025; Lu et al., 2022). Further studies could investigate the role of outward knowledge boundary of LLMs in the overall prompt robustness of LLMs, enabling them to express more knowledge they already possess.

Understanding Knowledge Boundary through Knowledge Mechanisms

Existing research on knowledge mechanisms, including memorization, comprehension, creation, and evolution, investigates how LLMs acquire, store, and utilize knowledge (Wang et al., 2024c). It is worth studying different phenomena of LLM knowledge boundaries under these mechanism views. For example, the outward knowledge boundary showcases how mechanisms like memorization and comprehension manifest in the explicit behaviors and outputs of the model, while the parametric boundary reflects a deeper, less visible level of how knowledge is embedded and structured through these mechanisms. The universal boundary can help measure the creative and evolutionary capabilities of LLMs.

In addition to the primary challenges outlined above, we also explore several critical prospects for the real-world applications of LLMs in relation to their knowledge boundaries in Appendix E. These include the *generalization of knowledge boundary*, *unintended side effects* of mitigation strategies, and issues related to the *knowledge boundary in long-form language modeling*.

7 Conclusions

This survey present a comprehensive overview of the knowledge boundary of LLMs, offering a formalized taxonomy and addressing key questions in the field. By exploring undesirable behaviors, identification techniques, and mitigation strategies, we emphasize the critical role of understanding and managing these boundaries to improve the reliability and utility of LLMs. Despite significant progress, challenges persist, including lack of comprehensive benchmarks, potential uses of knowledge boundary, and the role of knowledge boundary under various mechanisms. We hope this survey inspires continued exploration and innovation toward more trustworthy and reliable LLMs.

Limitations

We identify several limitations of our work.

Formal Definition of Knowledge This survey does not give a formal definition of the knowledge k , which is a critical problem in the scope of NLP research on knowledge. In this survey, we define the abstracted concept of knowledge as k , which is represented by a set of textual expressions of input and output Q_k . This definition can facilitate practical NLP experiments and efficient validation. In fact, the formal definition of knowledge is still a debatable topic, calling for future exploration. For example, Fierro et al. (2024) try to bridge the philosophical definition to the knowledge of LLMs, though significant disagreements persist among various philosophical schools of thought.

Various Forms of Textual Expressions regarding Different Knowledge Types Different types of knowledge may correspond to various forms of textual input-output Q_k , while we aim to provide a universal definition without the loss of generality. For instance, outputs for complex concepts may be open-ended and long-form, while simpler concepts might be expressed in a multiple-choice format. Some knowledge can be explicitly stated in the input, whereas others, such as commonsense knowledge, may need to be inferred from the input. Additionally, a single input may have multiple valid outputs. Some knowledge types, like mathematical knowledge, may inevitably involve multiple pieces of knowledge in a single input-output instance. As research progresses, a more nuanced definition for Q_k may be necessary to accommodate different knowledge types effectively.

(Un)Known to Human or Models Besides, in our definition, LLMs operate within the universal knowledge boundary, typically limited to human-known knowledge. We generally believe that LLMs do not possess knowledge beyond this boundary. However, there may be outliers that LLMs have knowledge that is unknown for human, which is not clearly studied in existing research. Wang et al. (2024b) hypothesize that LLMs may create new knowledge, but its reliability remains uncertain. While such outputs could reflect meaningful discoveries, they may also stem from implicit correlations in training data. Since existing research has not systematically examined LLM-generated unknown knowledge, its nature and implications remain unclear.

Missing Latest Studies Finally, we try to include all the related research, but it is possible that our survey miss some related work. Currently this is still an active research area, while our content has limited pages. We will maintain a Github repository to keep track on the research progress¹.

References

- Vaibhav Adlakha, Shehzaad Dhuliawala, Kaheer Suleman, Harm de Vries, and Siva Reddy. 2022. [Top-iOCQA: Open-domain conversational question answering with topic switching](#). *Transactions of the Association for Computational Linguistics*, 10:468–483.
- Ayushi Agarwal, Nisarg Patel, Neeraj Varshney, Mihir Parmar, Pavan Mallina, Aryan Bhavin Shah, Srihari Raju Sangaraju, Tirth Patel, Nihar Thakkar, and Chitta Baral. 2023. [Can NLP models ‘identify’, ‘distinguish’, and ‘justify’ questions that don’t have a definitive answer?](#) In *TrustNLP Workshop at ACL 2023*.
- Gustaf Ahlritz, Tian Qin, Nikhil Vyas, Boaz Barak, and Benjamin L. Edelman. 2024. [Distinguishing the knowable from the unknowable with language models](#). In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Lukas Aichberger, Kajetan Schweighofer, Mykyta Ielanskyi, and Sepp Hochreiter. 2024. [How many opinions does your LLM have? improving uncertainty estimation in NLG](#). In *ICLR 2024 Workshop on Secure and Trustworthy Large Language Models*.
- Mohammad Aliannejadi, Julia Kiseleva, Aleksandr Chuklin, Jeff Dalton, and Mikhail Burtsev. 2021. [Building and evaluating open-domain dialogue corpora with clarifying questions](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 4473–4484. Association for Computational Linguistics.
- Zeyuan Allen-Zhu and Yuanzhi Li. 2024. [Physics of language models: Part 3.1, knowledge storage and extraction](#). In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Alfonso Amayuelas, Kyle Wong, Liangming Pan, Wenhui Chen, and William Yang Wang. 2024. [Knowledge of knowledge: Exploring known-unknowns uncertainty with large language models](#). In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 6416–6432. Association for Computational Linguistics.

¹https://anonymous.4open.science/r/knowledge_boundary_survey-6E15.

- Chinmaya Andukuri, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah D. Goodman. 2024. [Star-gate: Teaching language models to ask clarifying questions](#). *CoRR*, abs/2403.19154.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. [Self-rag: Learning to retrieve, generate, and critique through self-reflection](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Amos Azaria and Tom M. Mitchell. 2023. [The internal state of an LLM knows when it’s lying](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 967–976. Association for Computational Linguistics.
- Baolong Bi, Shenghua Liu, Yiwei Wang, Lingrui Mei, and Xueqi Cheng. 2024. [Is factuality decoding a free lunch for llms? evaluation on knowledge editing benchmark](#). *CoRR*, abs/2404.00216.
- Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. 2023. [Discovering latent knowledge in language models without supervision](#). In *The Eleventh International Conference on Learning Representations*.
- Arslan Chaudhry, Sridhar Thiagarajan, and Dilan Gorur. 2024. [Finetuning language models to emit linguistic expressions of uncertainty](#). *arXiv preprint arXiv:2409.12180*.
- Canyu Chen, Baixiang Huang, Zekun Li, Zhaorun Chen, Shiyang Lai, Xiong Xiao Xu, Jia-Chen Gu, Jindong Gu, Huaxiu Yao, Chaowei Xiao, Xifeng Yan, William Yang Wang, Philip Torr, Dawn Song, and Kai Shu. 2024a. [Can editing llms inject harm?](#) *Preprint*, arXiv:2407.20224.
- Dingwei Chen, Feiteng Fang, Shiwen Ni, Feng Liang, Ruifeng Xu, Min Yang, and Chengming Li. 2024b. [Lower layer matters: Alleviating hallucination via multi-layer fusion contrastive decoding with truthfulness refocused](#). *arXiv preprint arXiv:2408.08769*.
- Jiuhai Chen and Jonas Mueller. 2024. [Quantifying uncertainty in answers from any language model and enhancing their trustworthiness](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5186–5200.
- Liang CHEN, Li Shen, Yang Deng, Xiaoyan Zhao, Bin Liang, and Kam-Fai Wong. 2025. [PEARL: Towards permutation-resilient LLMs](#). In *The Thirteenth International Conference on Learning Representations*.
- Lida Chen, Zujie Liang, Xintao Wang, Jiaqing Liang, Yanghua Xiao, Feng Wei, Jinglei Chen, Zhenghong Hao, Bing Han, and Wei Wang. 2024c. [Teaching large language models to express knowledge boundary from their own signals](#). *CoRR*, abs/2406.10881.
- Maximillian Chen, Ruoxi Sun, Sercan Ö. Arik, and Tomas Pfister. 2024d. [Learning to clarify: Multi-turn conversations with action-based contrastive self-training](#). *ArXiv*, abs/2406.00222.
- Shiqi Chen, Miao Xiong, Junteng Liu, Zhengxuan Wu, Teng Xiao, Siyang Gao, and Junxian He. 2024e. [In-context sharpness as alerts: An inner representation perspective for hallucination mitigation](#). In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Yue Chen, Chen Huang, Yang Deng, Wenqiang Lei, Dingnan Jin, Jia Liu, and Tat-Seng Chua. 2024f. [Style: Improving domain transferability of asking clarification questions in large language model powered conversational agents](#). In *Annual Meeting of the Association for Computational Linguistics*.
- Zhiyu Chen, Wenhui Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema N Moussa, Matthew I. Beane, Ting-Hao ‘Kenneth’ Huang, Bryan R. Routledge, and William Yang Wang. 2021. [Finqa: A dataset of numerical reasoning over financial data](#). *ArXiv*, abs/2109.00122.
- Zhongzhi Chen, Xingwu Sun, Xianfeng Jiao, Fengzong Lian, Zhanhui Kang, Di Wang, and Chengzhong Xu. 2024g. [Truth forest: Toward multi-scale truthfulness in large language models through intervention without tuning](#). In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024*, pages 20967–20974. AAAI Press.
- Jeffrey Cheng, Marc Marone, Orion Weller, Dawn Lawrie, Daniel Khazabi, and Benjamin Van Durme. 2024a. [Dated data: Tracing knowledge cutoffs in large language models](#). *ArXiv*, abs/2403.12958.
- Qinyuan Cheng, Tianxiang Sun, Xiangyang Liu, Wenwei Zhang, Zhangyue Yin, Shimin Li, Linyang Li, Zhengfu He, Kai Chen, and Xipeng Qiu. 2024b. [Can AI assistants know what they don’t know?](#) In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- John J. Cherian, Isaac Gibbs, and Emmanuel J. Candès. 2024. [Large language model validity via enhanced conformal prediction methods](#). *CoRR*, abs/2406.09714.
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R. Glass, and Pengcheng He. 2024. [Dola: Decoding by contrasting layers improves factuality in large language models](#). In *The Twelfth International Conference on Learning Representations*.
- Matthew Dahl, Varun Magesh, Mirac Suzgun, and Daniel E Ho. 2024. [Large Legal Fictions: Profiling Legal Hallucinations in Large Language Models](#). *Journal of Legal Analysis*, 16(1):64–93.
- Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yi-han Wang, Han Guo, Tianmin Shu, Meng Song, Eric

921	Xing, and Zhiting Hu. 2022. Rlprompt: Optimizing discrete text prompts with reinforcement learning. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3369–3391.	979
922		980
923		981
924		982
925		983
926	Yang Deng, Wenqiang Lei, Wai Lam, and Tat-Seng Chua. 2023a. A survey on proactive dialogue systems: Problems, methods, and prospects. In <i>Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI 2023</i> , pages 6583–6591.	984
927		985
928		986
929		987
930		988
931		989
932	Yang Deng, Lizi Liao, Liang Chen, Hongru Wang, Wenqiang Lei, and Tat-Seng Chua. 2023b. Prompting and evaluating large language models for proactive dialogues: Clarification, target-guided, and non-collaboration. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 10602–10621.	990
933		991
934		992
935		993
936		994
937		995
938		
939	Yang Deng, Yong Zhao, Moxin Li, See-Kiong Ng, and Tat-Seng Chua. 2024. Don't just say "i don't know"! self-aligning large language models for responding to unknown questions with explanations. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 13652–13673. Association for Computational Linguistics.	996
940		997
941		998
942		999
943		1000
944		1001
945		1002
946		1003
947	Shizhe Diao, Zhichao Huang, Ruijia Xu, Xuechun Li, LIN Yong, Xiao Zhou, and Tong Zhang. 2023. Black-box prompt learning for pre-trained language models. <i>Transactions on Machine Learning Research</i> .	1004
948		1005
949		1006
950		1007
951	Qingxiu Dong, Damai Dai, Yifan Song, Jingjing Xu, Zhifang Sui, and Lei Li. 2022. Calibrating factual knowledge in pretrained language models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022</i> , pages 5937–5947. Association for Computational Linguistics.	1008
952		1009
953		1010
954		1011
955		
956		
957		
958	Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Zhifang Sui, and Lei Li. 2024. Statistical knowledge assessment for large language models. <i>Advances in Neural Information Processing Systems</i> , 36.	1012
959		1013
960		1014
961		1015
962	Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2024. Improving factuality and reasoning in language models through multiagent debate. In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	1016
963		1017
964		1018
965		1019
966		1020
967		
968	Jinhao Duan, Hao Cheng, Shiqi Wang, Chenan Wang, Alex Zavalny, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. 2023. Shifting attention to relevance: Towards the uncertainty estimation of large language models. <i>CoRR</i> , abs/2307.01379.	1021
969		1022
970		1023
971		1024
972		1025
973	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.	1026
974		1027
975		
976		
977		
978		
	Ekaterina Fadeeva, Aleksandr Rubashevskii, Artem Shelmanov, Sergey Petrakov, Haonan Li, Hamdy Mubarak, Evgenii Tsymbalov, Gleb Kuzmin, Alexander Panchenko, Timothy Baldwin, Preslav Nakov, and Maxim Panov. 2024. Fact-checking the output of large language models via token-level uncertainty quantification. In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 9367–9385. Association for Computational Linguistics.	1028
		1029
		1030
		1031
		1032
		1033
	Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, Zhuo Han, Alan Huang, Songyang Zhang, Kai Chen, Zhixin Yin, Zongwen Shen, Jidong Ge, and Vincent Ng. 2024. Lawbench: Benchmarking legal knowledge of large language models. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	1034
		1035
	Shangbin Feng, Weijia Shi, Yike Wang, Wenxuan Ding, Vidhisha Balachandran, and Yulia Tsvetkov. 2024. Don't hallucinate, abstain: Identifying LLM knowledge gaps via multi-LLM collaboration. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14664–14690, Bangkok, Thailand. Association for Computational Linguistics.	
	Constanza Fierro, Ruchira Dhar, Filippos Stamatiou, Nicolas Garneau, and Anders Søgaard. 2024. Defining knowledge: Bridging epistemology and large language models. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 16096–16111. Association for Computational Linguistics.	
	Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata. 2023. Improving language model negotiation with self-play and in-context learning from ai feedback. <i>arXiv preprint arXiv:2305.10142</i> .	
	Chujie Gao, Siyuan Wu, Yue Huang, Dongping Chen, Qihui Zhang, Zhengyan Fu, Yao Wan, Lichao Sun, and Xiangliang Zhang. 2024a. Honestllm: Toward an honest and helpful large language model. <i>Preprint, arXiv:2406.00380</i> .	
	Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2023. Precise zero-shot dense retrieval without relevance labels. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1762–1777, Toronto, Canada. Association for Computational Linguistics.	
	Xiang Gao, Jiaxin Zhang, Lalla Mouatadid, and Kamalika Das. 2024b. Spuq: Perturbation-based uncertainty quantification for large language models. In <i>Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2336–2346.	
	Meiqi Guo, Mingda Zhang, Siva Reddy, and Malihe Alikhani. 2021. Abg-coqa: Clarifying ambiguity in	

1036	conversational question answering . In <i>3rd Conference on Automated Knowledge Base Construction, AKBC 2021, Virtual, October 4-8, 2021</i> .	1092
1037		1093
1038		1094
1039	Tom Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. 2023. Aging with GRACE: lifelong model editing with discrete key-value adaptors . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023</i> .	1095
1040		1096
1041		1097
1042		1098
1043		1099
1044		1100
1045		1101
1046	Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 6609–6625.	1102
1047		
1048		1103
1049		1104
1050		1105
1051		1106
1052	Bairu Hou, Yujian Liu, Kaizhi Qian, Jacob Andreas, Shiyu Chang, and Yang Zhang. 2024. Decomposing uncertainty for large language models through input clarification ensembling . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	1107
1053		1108
1054		
1055		1109
1056		1110
1057		1111
1058	Mengting Hu, Zhen Zhang, Shiwan Zhao, Minlie Huang, and Bingzhe Wu. 2023. Uncertainty in natural language processing: Sources, quantification, and applications. <i>arXiv preprint arXiv:2306.04459</i> .	1112
1059		1113
1060		1114
1061		1115
1062	Chao-Wei Huang and Yun-Nung Chen. 2024a. Fac-talign: Long-form factuality alignment of large language models . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024</i> , pages 16363–16375. Association for Computational Linguistics.	1116
1063		1117
1064		1118
1065		1119
1066		1120
1067		
1068	Chao-Wei Huang and Yun-Nung Chen. 2024b. Fac-talign: Long-form factuality alignment of large language models . <i>ArXiv</i> , abs/2410.01691.	1121
1069		1122
1070		1123
1071	Hsiu-Yuan Huang, Yutong Yang, Zhaoxi Zhang, San-woo Lee, and Yunfang Wu. 2024a. A survey of uncertainty estimation in llms: Theory meets practice . <i>CoRR</i> , abs/2410.15326.	1124
1072		1125
1073		
1074		1126
1075	Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2023a. Large language models can self-improve. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 1051–1068.	1127
1076		1128
1077		1129
1078		1130
1079		1131
1080	Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. 2024b. Large language models cannot self-correct reasoning yet . In <i>The Twelfth International Conference on Learning Representations</i> .	1132
1081		1133
1082		1134
1083		1135
1084		1136
1085		1137
1086	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023b. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions . <i>ArXiv</i> , abs/2311.05232.	1138
1087		1139
1088		1140
1089		1141
1090		1142
1091		1143
		1144
		1145
		1146
		1147
	Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming Zhao, Huaming Chen, Felix Juefei-Xu, and Lei Ma. 2023c. Look before you leap: An exploratory study of uncertainty measurement for large language models. <i>arXiv preprint arXiv:2307.10236</i> .	
	Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. 2023d. Transformer-patcher: One mistake worth one neuron . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
	Dan Iter, Reid Pryzant, Ruochen Xu, Shuohang Wang, Yang Liu, Yichong Xu, and Chenguang Zhu. 2023. In-context demonstration selection with cross entropy difference. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 1150–1162.	
	Ziwei Ji, Delong Chen, Etsuko Ishii, Samuel Cahyawijaya, Yejin Bang, Bryan Wilie, and Pascale Fung. 2024. LLM internal states reveal hallucination risk faced with a query . In <i>Proceedings of the 7th Black-boxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP</i> , pages 88–104, Miami, Florida, US. Association for Computational Linguistics.	
	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation . <i>ACM Comput. Surv.</i> , 55(12):248:1–248:38.	
	Mingjian Jiang, Yangjun Ruan, Sicong Huang, Saifei Liao, Silviu Pitis, Roger Baker Grosse, and Jimmy Ba. 2023a. Calibrating language models via augmented prompt ensembles. <i>ICML 2023 Workshop DeployableGenerativeAI</i> .	
	Zhengbao Jiang, Frank Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023b. Active retrieval augmented generation . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 7969–7992, Singapore. Association for Computational Linguistics.	
	Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W. Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	
	Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. <i>arXiv preprint arXiv:1705.03551</i> .	
	Raviraj Joshi, Kanishk Singla, Anusha Kamath, Raulnak Kalani, Rakesh Paul, Utkarsh Vaidya, Sanjay Singh Chauhan, Niranjana Wartikar, and Eileen Long. 2024. Adapting multilingual llms to low-resource languages using continued pre-training and synthetic corpus . <i>ArXiv</i> , abs/2410.14815.	

1148	Jaehun Jung, Lianhui Qin, Sean Welleck, Faeze Brah-	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	1205
1149	man, Chandra Bhagavatula, Ronan Le Bras, and	field, Michael Collins, Ankur Parikh, Chris Alberti,	1206
1150	Yejin Choi. 2022. Maieutic prompting: Logically	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-	1207
1151	consistent reasoning with recursive explanations. In	ton Lee, et al. 2019. Natural questions: a benchmark	1208
1152	<i>Proceedings of the 2022 Conference on Empirical</i>	for question answering research. <i>Transactions of the</i>	1209
1153	<i>Methods in Natural Language Processing</i> , pages	<i>Association for Computational Linguistics</i> , 7:453–	1210
1154	1266–1279.	466.	1211
1155	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom	James Law. 2006. Review of "algorithmic learning in	1212
1156	Henighan, Dawn Drain, Ethan Perez, Nicholas	a random world by vovk, gammerman and shafer" ,	1213
1157	Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli	springer , 2005, ISBN: 0-387-00152-2. <i>SIGACT</i>	1214
1158	Tran-Johnson, et al. 2022. Language models	<i>News</i> , 37(4):38–40.	1215
1159	(mostly) know what they know. <i>arXiv preprint</i>		
1160	<i>arXiv:2207.05221</i> .	Alina Leippert, Tatiana Anikina, Bernd Kiefer, and Josef	1216
1161	Jushi Kai, Tianhang Zhang, Hai Hu, and Zhouhan Lin.	Genabith. 2024. To clarify or not to clarify: A compar-	1217
1162	2024. SH2: self-highlighted hesitation helps you	ative analysis of clarification classification with	1218
1163	decode more truthfully . In <i>Findings of the Association</i>	fine-tuning, prompt tuning, and prompt engineering .	1219
1164	<i>for Computational Linguistics: EMNLP 2024,</i>	<i>Proceedings of the 2024 Conference of the North</i>	1220
1165	<i>Miami, Florida, USA, November 12-16, 2024</i> , pages	<i>American Chapter of the Association for Computa-</i>	1221
1166	4514–4530. Association for Computational Linguis-	<i>tional Linguistics: Human Language Technologies</i>	1222
1167	tics.	(<i>Volume 4: Student Research Workshop</i>).	1223
1168	Haoqiang Kang and Xiao-Yang Liu. 2024. Deficiency	Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter	1224
1169	of large language models in finance: An empirical	Pfister, and Martin Wattenberg. 2023a. Inference-	1225
1170	examination of hallucination . In <i>I Can't Believe It's</i>	time intervention: Eliciting truthful answers from a	1226
1171	<i>Not Better Workshop: Failure Modes in the Age of</i>	language model . In <i>Thirty-seventh Conference on</i>	1227
1172	<i>Foundation Models</i> .	<i>Neural Information Processing Systems</i> .	1228
1173	Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi,	Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter	1229
1174	Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir	Pfister, and Martin Wattenberg. 2024a. Inference-	1230
1175	Radev, Noah A. Smith, Yejin Choi, and Kentaro Inui.	time intervention: Eliciting truthful answers from a	1231
1176	2023. Realtime QA: what's the answer right now?	language model . <i>Advances in Neural Information</i>	1232
1177	In <i>Advances in Neural Information Processing Sys-</i>	<i>Processing Systems</i> , 36.	1233
1178	<i>tems 36: Annual Conference on Neural Information</i>		
1179	<i>Processing Systems 2023, NeurIPS 2023</i> .	Loka Li, Guangyi Chen, Yusheng Su, Zhenhao	1234
1180	Zixuan Ke, Weize Kong, Cheng Li, Mingyang Zhang,	Chen, Yixuan Zhang, Eric Xing, and Kun Zhang.	1235
1181	Qiaozhu Mei, and Michael Bendersky. 2024. Bridg-	2024b. Confidence matters: Revisiting intrinsic	1236
1182	ing the preference gap between retrievers and LLMs .	self-correction capabilities of large language mod-	1237
1183	In <i>Proceedings of the 62nd Annual Meeting of the</i>	els. <i>arXiv preprint arXiv:2402.12563</i> .	1238
1184	<i>Association for Computational Linguistics (Volume 1:</i>	Moxin Li, Wenjie Wang, Fuli Feng, Yixin Cao, Jizhi	1239
1185	<i>Long Papers)</i> , pages 10438–10451, Bangkok, Thai-	Zhang, and Tat-Seng Chua. 2023b. Robust prompt	1240
1186	land. Association for Computational Linguistics.	optimization for large language models against distri-	1241
1187	Jannik Kossen, Jiatong Han, Muhammed Razzak, Lisa	bution shifts. In <i>Proceedings of the 2023 Conference</i>	1242
1188	Schut, Shreshth Malik, and Yarin Gal. 2024. Sema-	<i>on Empirical Methods in Natural Language Process-</i>	1243
1189	ntic entropy probes: Robust and cheap hallucination	<i>ing</i> , pages 1539–1554.	1244
1190	detection in llms. <i>arXiv preprint arXiv:2406.15927</i> .	Moxin Li, Wenjie Wang, Fuli Feng, Fengbin Zhu, Qi-	1245
1191	Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2022.	fan Wang, and Tat-Seng Chua. 2024c. Think twice	1246
1192	Clam: Selective clarification for ambiguous ques-	before trusting: Self-detection for large language	1247
1193	tions with large language models. <i>arXiv preprint</i>	models through comprehensive answer reflection .	1248
1194	<i>arXiv:2212.07769</i> .	In <i>Findings of the Association for Computational</i>	1249
1195	Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023.	<i>Linguistics: EMNLP 2024, Miami, Florida, USA,</i>	1250
1196	Semantic uncertainty: Linguistic invariances for un-	<i>November 12-16, 2024</i> , pages 11858–11875. Associ-	1251
1197	certainty estimation in natural language generation .	ation for Computational Linguistics.	1252
1198	In <i>The Eleventh International Conference on Learn-</i>	Shuo Li, Sangdon Park, Insup Lee, and Osbert Bastani.	1253
1199	<i>ing Representations</i> .	2024d. TRAQ: trustworthy retrieval augmented ques-	1254
1200	Bhawesh Kumar, Charlie Lu, Gauri Gupta, Anil Palepu,	tion answering via conformal prediction . In <i>Proceed-</i>	1255
1201	David R. Bellamy, Ramesh Raskar, and Andrew	<i>ings of the 2024 Conference of the North American</i>	1256
1202	Beam. 2023. Conformal prediction with large lan-	<i>Chapter of the Association for Computational Lin-</i>	1257
1203	guage models for multi-choice question answering .	<i>guistics: Human Language Technologies (Volume 1:</i>	1258
1204	<i>CoRR</i> , abs/2305.18404.	<i>Long Papers)</i> , NAACL 2024, Mexico City, Mexico,	1259
		<i>June 16-21, 2024</i> , pages 3799–3821. Association for	1260
		Computational Linguistics.	1261

1262	Siheng Li, Cheng Yang, Taiqiang Wu, Chufan Shi,	Chen Ling, Xujiang Zhao, Xuchao Zhang, Wei Cheng,	1317
1263	Yuji Zhang, Xinyu Zhu, Zesen Cheng, Deng Cai,	Yanchi Liu, Yiyao Sun, Mika Oishi, Takao Osaki,	1318
1264	Mo Yu, Lemao Liu, Jie Zhou, Yujiu Yang, Ngai	Katsushi Matsuda, Jie Ji, et al. 2024. Uncertainty	1319
1265	Wong, Xixin Wu, and Wai Lam. 2024e. A survey on	quantification for in-context learning of large lan-	1320
1266	the honesty of large language models. <i>arXiv preprint</i>	guage models. In <i>Proceedings of the 2024 Confer-</i>	1321
1267	<i>arXiv:2409.18786</i> .	<i>ence of the North American Chapter of the Associ-</i>	1322
		<i>ation for Computational Linguistics: Human Lan-</i>	1323
1268	Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang,	<i>guage Technologies (Volume 1: Long Papers)</i> , pages	1324
1269	Jason Eisner, Tatsunori B Hashimoto, Luke Zettle-	3357–3370.	1325
1270	moyer, and Mike Lewis. 2023c. Contrastive decod-		
1271	ing: Open-ended text generation as optimization. In	Alisa Liu, Zhaofeng Wu, Julian Michael, Alane Suhr,	1326
1272	<i>Proceedings of the 61st Annual Meeting of the As-</i>	Peter West, Alexander Koller, Swabha Swayamdipta,	1327
1273	<i>sociation for Computational Linguistics (Volume 1:</i>	Noah A. Smith, and Yejin Choi. 2023. We’re afraid	1328
1274	<i>Long Papers)</i> , pages 12286–12312.	language models aren’t modeling ambiguity . In <i>Pro-</i>	1329
		<i>ceedings of the 2023 Conference on Empirical Meth-</i>	1330
1275	Xiaonan Li, Kai Lv, Hang Yan, Tianyang Lin, Wei	<i>ods in Natural Language Processing, EMNLP 2023,</i>	1331
1276	Zhu, Yuan Ni, Guotong Xie, Xiaoling Wang, and	<i>Singapore, December 6-10, 2023</i> , pages 790–807.	1332
1277	Xipeng Qiu. 2023d. Unified demonstration retriever	Association for Computational Linguistics.	1333
1278	for in-context learning. In <i>Proceedings of the 61st</i>		
1279	<i>Annual Meeting of the Association for Computational</i>	Genglin Liu, Xingyao Wang, Lifan Yuan, Yangyi Chen,	1334
1280	<i>Linguistics (Volume 1: Long Papers)</i> , pages 4644–	and Hao Peng. 2024a. Examining llms’ uncer-	1335
1281	4668.	tainty expression towards questions outside paramet-	1336
		ric knowledge . <i>Preprint</i> , arXiv:2311.09731.	1337
1282	Xiaonan Li and Xipeng Qiu. 2023. Finding support		
1283	examples for in-context learning. In <i>Findings of the</i>	Jiachang Liu, Dinghan Shen, Yizhe Zhang, William B	1338
1284	<i>Association for Computational Linguistics: EMNLP</i>	Dolan, Lawrence Carin, and Weizhu Chen. 2022a.	1339
1285	<i>2023</i> , pages 6219–6235.	What makes good in-context examples for gpt-3?	1340
		In <i>Proceedings of Deep Learning Inside Out (Dee-</i>	1341
1286	Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng	<i>LIO 2022): The 3rd Workshop on Knowledge Extrac-</i>	1342
1287	Ding, Shafiq Joty, Soujanya Poria, and Lidong Bing.	<i>tion and Integration for Deep Learning Architectures,</i>	1343
1288	2024f. Chain-of-knowledge: Grounding large lan-	pages 100–114.	1344
1289	guage models via dynamic knowledge adapting over		
1290	heterogeneous sources . In <i>The Twelfth International</i>	Jiacheng Liu, Alisa Liu, Ximing Lu, Sean Welleck, Pe-	1345
1291	<i>Conference on Learning Representations</i> .	ter West, Ronan Le Bras, Yejin Choi, and Hannaneh	1346
		Hajishirzi. 2022b. Generated knowledge prompting	1347
1292	Yinghui Li, Haojing Huang, Jiayi Kuang, Yangning Li,	for commonsense reasoning. In <i>Proceedings of the</i>	1348
1293	Shu-Yu Guo, Chao Qu, Xiaoyu Tan, Hai-Tao Zheng,	<i>60th Annual Meeting of the Association for Compu-</i>	1349
1294	Ying Shen, and Philip S. Yu. 2025. Refine knowledge	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	1350
1295	of large language models via adaptive contrastive	3154–3169.	1351
1296	learning .		
		Junteng Liu, Shiqi Chen, Yu Cheng, and Junxian He.	1352
1297	Sheng-Chieh Lin, Luyu Gao, Barlas Oguz, Wenhan	2024b. On the universal truthfulness hyperplane in-	1353
1298	Xiong, Jimmy Lin, Scott Yih, and Xilun Chen. 2024a.	side llms . In <i>Proceedings of the 2024 Conference on</i>	1354
1299	FLAME : Factuality-aware alignment for large lan-	<i>Empirical Methods in Natural Language Processing,</i>	1355
1300	guage models . In <i>Advances in Neural Information</i>	<i>EMNLP 2024, Miami, FL, USA, November 12-16,</i>	1356
1301	<i>Processing Systems 38: Annual Conference on Neu-</i>	2024, pages 18199–18224. Association for Computa-	1357
1302	<i>ral Information Processing Systems 2024, NeurIPS</i>	tional Linguistics.	1358
1303	<i>2024, Vancouver, BC, Canada, December 10 - 15,</i>		
1304	<i>2024</i> .	Kai Liu, Ze Chen, Zhihang Fu, Rongxin Jiang, Fan	1359
		Zhou, Yao-Shen Chen, Yue Wu, and Jieping Ye.	1360
1305	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022a.	2024c. Structure-aware domain knowledge injection	1361
1306	Teaching models to express their uncertainty in	for large language models .	1362
1307	words . <i>Trans. Mach. Learn. Res.</i> , 2022.		
		Xin Liu, Muhammad Khalifa, and Lu Wang. 2024d.	1363
1308	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022b.	Litcab: Lightweight language model calibration over	1364
1309	Truthfulqa: Measuring how models mimic human	short-and long-form responses. In <i>The Twelfth Inter-</i>	1365
1310	falsehoods. In <i>Proceedings of the 60th Annual Meet-</i>	<i>national Conference on Learning Representations</i> .	1366
1311	<i>ing of the Association for Computational Linguistics</i>		
1312	<i>(Volume 1: Long Papers)</i> , pages 3214–3252.	Do Xuan Long, Yiran Zhao, Hannah Brown, Yuxi Xie,	1367
		James Xu Zhao, Nancy F. Chen, Kenji Kawaguchi,	1368
1313	Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2024b.	Michael Shieh, and Junxian He. 2024. Prompt op-	1369
1314	Generating with confidence: Uncertainty quantifica-	timization via adversarial in-context learning . In	1370
1315	tion for black-box large language models . <i>Transac-</i>	<i>Proceedings of the 62nd Annual Meeting of the As-</i>	1371
1316	<i>tions on Machine Learning Research</i> .	<i>sociation for Computational Linguistics (Volume 1:</i>	1372

1373		<i>Long Papers</i>), <i>ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 7308–7327. Association for Computational Linguistics.	
1374			
1375			
1376	Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel,		
1377	and Pontus Stenetorp. 2022. Fantastically ordered		
1378	prompts and where to find them: Overcoming few-		
1379	shot prompt order sensitivity . In <i>Proceedings of the</i>		
1380	<i>60th Annual Meeting of the Association for Computa-</i>		
1381	<i>tional Linguistics (Volume 1: Long Papers), ACL</i>		
1382	<i>2022, Dublin, Ireland, May 22-27, 2022</i> , pages 8086–		
1383	8098. Association for Computational Linguistics.		
1384	Man Luo, Xin Xu, Zhuyun Dai, Panupong Pasu-		
1385	pat, Mehran Kazemi, Chitta Baral, Vaiva Im-		
1386	brasaitė, and Vincent Y Zhao. 2023. Dr.ICL:		
1387	Demonstration-retrieved in-context learning . In <i>R0-</i>		
1388	<i>FoMo: Robustness of Few-shot and Zero-shot Learn-</i>		
1389	<i>ing in Large Foundation Models</i> .		
1390	Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao,		
1391	and Nan Duan. 2023. Query rewriting in retrieval-		
1392	augmented large language models . In <i>Proceedings of</i>		
1393	<i>the 2023 Conference on Empirical Methods in Natural</i>		
1394	<i>Language Processing</i> , pages 5303–5315, Singa-		
1395	pore. Association for Computational Linguistics.		
1396	Aman Madaan, Niket Tandon, Peter Clark, and Yim-		
1397	ing Yang. 2022. Memory-assisted prompt editing to		
1398	improve GPT-3 after deployment . In <i>Proceedings of</i>		
1399	<i>the 2022 Conference on Empirical Methods in Natural</i>		
1400	<i>Language Processing, EMNLP 2022, Abu Dhabi,</i>		
1401	<i>United Arab Emirates, December 7-11, 2022</i> , pages		
1402	2833–2861. Association for Computational Linguis-		
1403	tics.		
1404	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler		
1405	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,		
1406	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,		
1407	et al. 2024. Self-refine: Iterative refinement with		
1408	self-feedback. <i>Advances in Neural Information Pro-</i>		
1409	<i>cessing Systems</i> , 36.		
1410	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das,		
1411	Daniel Khashabi, and Hannaneh Hajishirzi. 2022.		
1412	When not to trust language models: Investigating		
1413	effectiveness of parametric and non-parametric mem-		
1414	ories. <i>arXiv preprint arXiv:2212.10511</i> .		
1415	Potsawee Manakul, Adian Liusie, and Mark Gales. 2023.		
1416	Selfcheckgpt: Zero-resource black-box hallucina-		
1417	tion detection for generative large language models.		
1418	In <i>Proceedings of the 2023 Conference on Empiri-</i>		
1419	<i>cal Methods in Natural Language Processing</i> , pages		
1420	9004–9017.		
1421	Samuel Marks and Max Tegmark. 2023. The geometry		
1422	of truth: Emergent linear structure in large language		
1423	model representations of true/false datasets. <i>arXiv</i>		
1424	<i>preprint arXiv:2310.06824</i> .		
1425	Costas Mavromatis, Balasubramaniam Srinivasan,		
1426	Zhengyuan Shen, Jiani Zhang, Huzefa Rangwala,		
1427	Christos Faloutsos, and George Karypis. 2023.		
1428	Which examples to annotate for in-context learn-		
1429	ing? towards effective and efficient selection. <i>arXiv</i>		
1430	<i>preprint arXiv:2310.20046</i> .		
	Nick Mecklenburg, Yiyu Lin, Xiaoxiao Li, Daniel		1431
	Holstein, Leonardo Nunes, Sara Malvar, Bruno		1432
	Leonardo Barros Silva, Ranveer Chandra, Vijay Aski,		1433
	Pavan Kumar Reddy Yannam, Tolga Aktas, and Todd		1434
	Hendry. 2024. Injecting new knowledge into large		1435
	language models via supervised fine-tuning . <i>ArXiv</i> ,		1436
	abs/2404.00213.		1437
	Kevin Meng, David Bau, Alex Andonian, and Yonatan		1438
	Belinkov. 2022. Locating and editing factual associ-		1439
	ations in GPT. In <i>Advances in Neural Information</i>		1440
	<i>Processing Systems 35: Annual Conference on Neu-</i>		1441
	<i>ral Information Processing Systems 2022, NeurIPS</i>		1442
	<i>2022</i> .		1443
	Kevin Meng, Arnab Sen Sharma, Alex J. Andonian,		1444
	Yonatan Belinkov, and David Bau. 2023. Mass-		1445
	editing memory in a transformer . In <i>The Eleventh</i>		1446
	<i>International Conference on Learning Representa-</i>		1447
	<i>tions, ICLR 2023</i> .		1448
	Ning Miao, Yee Whye Teh, and Tom Rainforth. 2024.		1449
	Selfcheck: Using LLMs to zero-shot check their own		1450
	step-by-step reasoning . In <i>The Twelfth International</i>		1451
	<i>Conference on Learning Representations</i> .		1452
	Sabrina J Mielke, Arthur Szlam, Emily Dinan, and Y-		1453
	Lan Boureau. 2022. Reducing conversational agents’		1454
	overconfidence through linguistic calibration. <i>Trans-</i>		1455
	<i>actions of the Association for Computational Linguis-</i>		1456
	<i>tics</i> , 10:857–872.		1457
	Sewon Min, Kalpesh Krishna, Xinxin Lyu, Mike		1458
	Lewis, Wen tau Yih, Pang Wei Koh, Mohit Iyyer,		1459
	Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023.		1460
	Factscore: Fine-grained atomic evaluation of fac-		1461
	tual precision in long form text generation . <i>ArXiv</i> ,		1462
	abs/2305.14251.		1463
	Eric Mitchell, Charles Lin, Antoine Bosselut, Christo-		1464
	pher D. Manning, and Chelsea Finn. 2022. Memory-		1465
	based model editing at scale . In <i>International Con-</i>		1466
	<i>ference on Machine Learning, ICML 2022</i> , volume		1467
	162 of <i>Proceedings of Machine Learning Research</i> ,		1468
	pages 15817–15831.		1469
	Christopher Mohri and Tatsunori Hashimoto. 2024.		1470
	Language models with conformal factuality guar-		1471
	antees . In <i>Forty-first International Conference on</i>		1472
	<i>Machine Learning, ICML 2024, Vienna, Austria, July</i>		1473
	<i>21-27, 2024</i> . OpenReview.net.		1474
	Fangwen Mu, Lin Shi, Song Wang, Zhuohao Yu, Bin-		1475
	quan Zhang, ChenXue Wang, Shichao Liu, and Qing		1476
	Wang. 2023. Clarifygpt: Empowering llm-based		1477
	code generation with intention clarification . <i>ArXiv</i> ,		1478
	abs/2310.10996.		1479
	Tarek Naous, Michael J. Ryan, Alan Ritter, and Wei		1480
	Xu. 2024. Having beer after prayer? measuring cul-		1481
	tural bias in large language models . In <i>Proceedings</i>		1482
	<i>of the 62nd Annual Meeting of the Association for</i>		1483
	<i>Computational Linguistics (Volume 1: Long Papers),</i>		1484
	<i>ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> ,		1485
	pages 16366–16393. Association for Computational		1486
	Linguistics.		1487

1488	Alexander Nikitin, Jannik Kossen, Yarin Gal, and Pekka Marttinen. 2024. Kernel language entropy: Fine-grained uncertainty quantification for llms from semantic similarities. <i>arXiv preprint arXiv:2405.20003</i> .	1548
1489		1547
1490		1548
1491		1549
1492		1550
1493	Yasumasa Onoe, Michael J. Q. Zhang, Eunsol Choi, and Greg Durrett. 2022. Entity cloze by date: What llms know about unseen entities . In <i>Findings of the Association for Computational Linguistics: NAACL 2022, Seattle, WA, United States, July 10-15, 2022</i> , pages 693–702. Association for Computational Linguistics.	1551
1494		1552
1495		1553
1496		1554
1497		1555
1498		1556
1499	Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2023. Med-halt: Medical domain hallucination test for large language models . In <i>Proceedings of the 27th Conference on Computational Natural Language Learning, CoNLL 2023, Singapore, December 6-7, 2023</i> , pages 314–334. Association for Computational Linguistics.	1557
1500		1558
1501		1559
1502		1560
1503		1561
1504		1562
1505		
1506	Soumen Pal, Manojit Bhattacharya, Sang-Soo Lee, and Chiranjib Chakraborty. 2024. A domain-specific next-generation large language model (llm) or chatgpt is required for biomedical engineering and research. <i>Annals of biomedical engineering</i> , 52(3):451–454.	1563
1507		1564
1508		1565
1509		1566
1510		
1511		
1512	Yikang Pan, Liangming Pan, Wenhu Chen, Preslav Nakov, Min-Yen Kan, and William Yang Wang. 2023. On the risk of misinformation pollution with large language models . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 1389–1403. Association for Computational Linguistics.	1567
1513		1568
1514		1569
1515		1570
1516		1571
1517		1572
1518		1573
1519	Archiki Prasad, Peter Hase, Xiang Zhou, and Mohit Bansal. 2023. Grips: Gradient-free, edit-based instruction search for prompting large language models. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 3845–3864.	1574
1520		1575
1521		1576
1522		1577
1523		1578
1524		
1525	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2023. Measuring and narrowing the compositionality gap in language models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 5687–5711.	1579
1526		1580
1527		1581
1528		1582
1529		1583
1530	Reid Pryzant, Dan Iter, Jerry Li, Yin Lee, Chenguang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with “gradient descent” and beam search. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 7957–7968.	1584
1531		1585
1532		1586
1533		1587
1534		1588
1535		1589
1536	Yifu Qiu, Zheng Zhao, Yftah Ziser, Anna Korhonen, Edoardo M Ponti, and Shay B Cohen. 2024. Spectral editing of activations for large language model alignment. <i>arXiv preprint arXiv:2405.09719</i> .	1590
1537		1591
1538		1592
1539		1593
1540	Victor Quach, Adam Fisch, Tal Schuster, Adam Yala, Jae Ho Sohn, Tommi S. Jaakkola, and Regina Barzilay. 2024. Conformal language modeling . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	1594
1541		1595
1542		1596
1543		
1544		
1545		
	Shauli Ravfogel, Yoav Goldberg, and Jacob Goldberger. 2023. Conformal nucleus sampling . In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 27–34. Association for Computational Linguistics.	1597
		1598
		1599
		1600
		1601
	Ruiyang Ren, Yuhao Wang, Yingqi Qu, Wayne Xin Zhao, Jing Liu, Hao Tian, Hua Wu, Ji-Rong Wen, and Haifeng Wang. 2023. Investigating the factual knowledge boundary of large language models with retrieval augmentation. <i>arXiv preprint arXiv:2307.11019</i> .	
	Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. Learning to retrieve prompts for in-context learning. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2655–2671.	
	Abulhair Saparov and He He. 2023. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought . In <i>The Eleventh International Conference on Learning Representations</i> .	
	Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9248–9274, Singapore. Association for Computational Linguistics.	
	Maohao Shen, Subhro Das, Kristjan Greenewald, Prasanna Sattigeri, Gregory Wornell, and Soumya Ghosh. 2024. Thermometer: Towards universal calibration for large language models. <i>arXiv preprint arXiv:2403.08819</i> .	
	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In <i>International Conference on Machine Learning</i> , pages 31210–31227. PMLR.	
	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024. REPLUG: Retrieval-augmented black-box language models . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 8371–8384, Mexico City, Mexico. Association for Computational Linguistics.	
	Vaishnavi Shrivastava, Percy Liang, and Ananya Kumar. 2024. Llamas know what GPTs don’t show: Surrogate models for selective classification .	
	Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Lee Boyd-Graber, and Lijuan Wang. 2023. Prompting GPT-3 to be reliable . In <i>The Eleventh International Conference on Learning Representations</i> .	

- Smriti Singh, Shuvam Keshari, Vinija Jain, and Aman Chadha. 2024. [Born with a silver spoon? investigating socioeconomic bias in large language models](#). *ArXiv*, abs/2403.14633.
- Maojia Song, Shang Hong Sim, Rishabh Bhardwaj, Hai Leong Chieu, Navonil Majumder, and Soujanya Poria. 2024a. [Measuring and enhancing trustworthiness of llms in rag through grounded attributions and learning to refuse](#). *ArXiv*, abs/2409.11242.
- Xiaoshuai Song, Zhengyang Wang, Keqing He, Guanting Dong, Yutao Mou, Jinxu Zhao, and Weiran Xu. 2024b. [Knowledge editing on black-box large language models](#). *CoRR*, abs/2402.08631.
- Elias Stengel-Eskin, Peter Hase, and Mohit Bansal. 2024. [Lacie: Listener-aware finetuning for confidence calibration in large language models](#). *arXiv preprint arXiv:2405.21028*.
- Jiayuan Su, Jing Luo, Hongwei Wang, and Lu Cheng. 2024. [API is enough: Conformal prediction for large language models without logit-access](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 979–995. Association for Computational Linguistics.
- Chenmien Tan, Ge Zhang, and Jie Fu. 2024. [Massive editing for large language models via meta learning](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Shuchang Tao, Liuyi Yao, Hanxing Ding, Yuexiang Xie, Qi Cao, Fei Sun, Jinyang Gao, Huawei Shen, and Bolin Ding. 2024. [When to trust LLMs: Aligning confidence with response quality](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 5984–5996, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D. Manning, and Chelsea Finn. 2024. [Fine-tuning language models for factuality](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442.
- Christian Tomani, Kamalika Chaudhuri, I. Evtimov, Daniel Cremers, and Mark Ibrahim. 2024. [Uncertainty-based abstention in llms improves safety and reduces hallucinations](#). *ArXiv*, abs/2404.10960.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *CoRR*, abs/2307.09288.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. [MuSiQue: Multi-hop questions via single-hop question composition](#). *Transactions of the Association for Computational Linguistics*, 10:539–554.
- Neeraj Varshney, Pavel Dolin, Agastya Seth, and Chitta Baral. 2023. [The art of defending: A systematic evaluation and analysis of llm defense strategies on safety and over-defensiveness](#). *ArXiv*, abs/2401.00287.
- Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and Thang Luong. 2023. [Freshllms: Refreshing large language models with search engine augmentation](#). In *Annual Meeting of the Association for Computational Linguistics*.
- Boshi Wang, Xiang Deng, and Huan Sun. 2022. Iteratively prompt pre-trained language models for chain of thought. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2714–2730.
- Boshi Wang, Xiang Yue, and Huan Sun. 2023a. [Can chatgpt defend its belief in truth? evaluating LLM reasoning via debate](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11865–11881.
- Jianing Wang, Yang Zhou, Xiaocheng Zhang, Mengjiao Bao, and Peng Yan. 2024a. Self-evolutionary large language models through uncertainty-enhanced preference optimization. *arXiv preprint arXiv:2409.11212*.
- Mengru Wang, Yunzhi Yao, Ziwen Xu, Shuofei Qiao, Shumin Deng, Peng Wang, Xiang Chen, Jia-Chen Gu, Yong Jiang, Pengjun Xie, Fei Huang, Huajun Chen, and Ningyu Zhang. 2024b. [Knowledge mechanisms in large language models: A survey and perspective](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA*,

1716	November 12-16, 2024, pages 7097–7135. Association	2024, Miami, Florida, USA, November 12-16, 2024,	1772
1717	tion for Computational Linguistics.	pages 3437–3450. Association for Computational	1773
1718	Mengru Wang, Yunzhi Yao, Ziwen Xu, Shuofei Qiao,		1774
1719	Shumin Deng, Peng Wang, Xiang Chen, Jia-Chen Gu,	Bingbing Wen, Jihan Yao, Shangbin Feng, Chenjun	1775
1720	Yong Jiang, Pengjun Xie, et al. 2024c. Knowledge	Xu, Yulia Tsvetkov, Bill Howe, and Lucy Lu Wang.	1776
1721	mechanisms in large language models: A survey and	2024b. Know your limits: A survey of abstention	1777
1722	perspective. <i>arXiv preprint arXiv:2407.15017</i> .	in large language models. <i>arXiv preprint</i>	1778
1723		<i>arXiv:2407.18418</i> .	1779
1724	Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi	Zhihua Wen, Zhiliang Tian, Zexin Jian, Zhen Huang,	1780
1725	Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Hua-	Pei Ke, Yifu Gao, Minlie Huang, and Dongsheng Li.	1781
1726	jun Chen. 2024d. WISE: rethinking the knowledge	2024c. Perception of knowledge boundary for large	1782
1727	memory for lifelong model editing of large language	language models through semi-open-ended question	1783
	models . <i>CoRR</i> , abs/2405.14768.	answering . <i>CoRR</i> , abs/2405.14383.	1784
1728	Xinyuan Wang, Chenxi Li, Zhen Wang, Fan Bai, Hao-		
1729	tian Luo, Jiayou Zhang, Nebojsa Jojic, Eric Xing, and	Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He,	1785
1730	Zhiting Hu. 2024e. Promptagent: Strategic planning	Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao.	1786
1731	with language models enables expert-level prompt op-	2023. Large language models are better reasoners	1787
1732	timization . In <i>The Twelfth International Conference</i>	with self-verification. In <i>Findings of the Association</i>	1788
1733	<i>on Learning Representations</i> .	<i>for Computational Linguistics: EMNLP 2023</i> , pages	1789
		2550–2575.	1790
1734	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le,	Jian Wu, Linyi Yang, Yuliang Ji, Wenhao Huang,	1791
1735	Ed H. Chi, Sharan Narang, Aakanksha Chowdhery,	Börje F Karlsson, and Manabu Okumura. 2024a.	1792
1736	and Denny Zhou. 2023b. Self-consistency improves	Gendec: A robust generative question-decomposition	1793
1737	chain of thought reasoning in language models . In	method for multi-hop reasoning . <i>arXiv preprint</i>	1794
1738	<i>The Eleventh International Conference on Learning</i>	<i>arXiv:2402.11166</i> .	1795
1739	<i>Representations</i> .		
1740	Zhen Wang, Peide Zhu, and Jie Yang. 2024f. Controversialqa: Exploring controversy in question answering .	Siye Wu, Jian Xie, Jiangjie Chen, Tinghui Zhu, Kai	1796
1741	In <i>Proceedings of the 2024 Joint International Con-</i>	Zhang, and Yanghua Xiao. 2024b. How easily do	1797
1742	<i>ference on Computational Linguistics, Language Re-</i>	irrelevant inputs skew the responses of large language	1798
1743	<i>sources and Evaluation, LREC/COLING 2024, 20-25</i>	models? <i>ArXiv</i> , abs/2404.03302.	1799
1744	<i>May, 2024, Torino, Italy</i> , pages 3962–3966. ELRA		
1745	and ICCL.	Johnathan Xie, Annie S Chen, Yoonho Lee, Eric	1800
1746		Mitchell, and Chelsea Finn. 2024. Calibrating lan-	1801
1747	Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md Rizwan	guage models with adaptive temperature scaling . In	1802
1748	Parvez, and Graham Neubig. 2023c. Learning to fil-	<i>Proceedings of the 2024 Conference on Empirical</i>	1803
1749	ter context for retrieval-augmented generation. <i>arXiv</i>	<i>Methods in Natural Language Processing</i> , pages	1804
1750	<i>preprint arXiv:2311.08377</i> .	18128–18138, Miami, Florida, USA. Association for	1805
		Computational Linguistics.	1806
1751	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong	1807
1752	Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le,	Zhang. 2024a. Benchmarking retrieval-augmented	1808
1753	and Denny Zhou. 2022a. Chain-of-thought prompt-	generation for medicine . <i>ArXiv</i> , abs/2402.13178.	1809
1754	ing elicits reasoning in large language models . In		
1755	<i>Advances in Neural Information Processing Systems</i> ,	Miao Xiong, Zhiyuan Hu, Xinyang Lu, YIFEI LI, Jie	1810
1756	volume 35, pages 24824–24837. Curran Associates,	Fu, Junxian He, and Bryan Hooi. 2024b. Can LLMs	1811
1757	Inc.	express their uncertainty? an empirical evaluation of	1812
		confidence elicitation in LLMs . In <i>The Twelfth Inter-</i>	1813
1758	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	<i>national Conference on Learning Representations</i> .	1814
1759	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,		
1760	et al. 2022b. Chain-of-thought prompting elicits rea-	Hongshen Xu, Zichen Zhu, Da Ma, Situo Zhang, Shuai	1815
1761	soning in large language models. <i>Advances in neural</i>	Fan, Lu Chen, and Kai Yu. 2024a. Rejection im-	1816
1762	<i>information processing systems</i> , 35:24824–24837.	proves reliability: Training llms to refuse unknown	1817
		questions using rl from knowledge feedback . <i>ArXiv</i> ,	1818
1763	Jerry Wei, Chengrun Yang, Xinying Song, Yifeng Lu,	abs/2403.18349.	1819
1764	Nathan Hu, Dustin Tran, Daiyi Peng, Ruibo Liu,		
1765	Da Huang, Cosmo Du, et al. 2024. Long-form fac-	Rongwu Xu, Brian S. Lin, Shujian Yang, Tianqi Zhang,	1820
1766	tuality in large language models. <i>arXiv preprint</i>	Weiyang Shi, Tianwei Zhang, Zhixuan Fang, Wei Xu,	1821
1767	<i>arXiv:2403.18802</i> .	and Han Qiu. 2024b. The earth is flat because...:	1822
		Investigating llms’ belief towards misinformation	1823
1768	Bingbing Wen, Bill Howe, and Lucy Lu Wang. 2024a.	via persuasive conversation . In <i>Proceedings of the</i>	1824
1769	Characterizing LLM abstention behavior in science	<i>62nd Annual Meeting of the Association for Compu-</i>	1825
1770	QA with context perturbations . In <i>Findings of the</i>	<i>tational Linguistics (Volume 1: Long Papers)</i> , <i>ACL</i>	1826
1771	<i>Association for Computational Linguistics: EMNLP</i>	2024, pages 16259–16303.	1827

1828	Xin Xu, Yue Liu, Panupong Pasupat, Mehran Kazemi, et al. 2024c. In-context learning with retrieved demonstrations for language models: A survey. <i>arXiv preprint arXiv:2401.11624</i> .	1884
1829		1885
1830		1886
1831		1887
1832	Yasin Abbasi Yadkori, Ilja Kuzborskij, András György, and Csaba Szepesvári. 2024. To believe or not to believe your llm. <i>arXiv preprint arXiv:2406.02543</i> .	1888
1833		
1834		
1835	Yuzi Yan, Xingzhou Lou, Jialian Li, Yiping Zhang, Jian Xie, Chao Yu, Yu Wang, Dong Yan, and Yuan Shen. 2024. Reward-robust rlhf in llms. <i>arXiv preprint arXiv:2409.15360</i> .	1889
1836		1890
1837		1891
1838		1892
1839	Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. 2024a. Large language models as optimizers . In <i>The Twelfth International Conference on Learning Representations</i> .	1893
1840		1894
1841		1895
1842		1896
1843		
1844	Dingkang Yang, Dongling Xiao, Jinjie Wei, Mingcheng Li, Zhaoyu Chen, Ke Li, and Lihua Zhang. 2024b. Improving factuality in large language models via decoding-time hallucinatory and truthful comparators. <i>arXiv preprint arXiv:2408.12325</i> .	1897
1845		1898
1846		1899
1847		1900
1848		
1849	Haoyan Yang, Zhitao Li, Yong Zhang, Jianzong Wang, Ning Cheng, Ming Li, and Jing Xiao. 2023a. PRCA: Fitting black-box large language models for retrieval question answering via pluggable reward-driven contextual adapter . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 5364–5375, Singapore. Association for Computational Linguistics.	1901
1850		1902
1851		1903
1852		1904
1853		1905
1854		1906
1855		1907
1856		1908
1857	Muchen Yang, Moxin Li, Yongle Li, Zijun Chen, Chongming Gao, Junqi Zhang, Yangyang Li, and Fuli Feng. 2024c. Dual-phase accelerated prompt optimization . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 12163–12173, Miami, Florida, USA. Association for Computational Linguistics.	1909
1858		1910
1859		1911
1860		1912
1861		1913
1862		1914
1863		1915
1864	Yuqing Yang, Ethan Chern, Xipeng Qiu, Graham Neubig, and Pengfei Liu. 2023b. Alignment for honesty . <i>ArXiv</i> , abs/2312.07000.	1916
1865		1917
1866		1918
1867	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018</i> , pages 2369–2380. Association for Computational Linguistics.	1919
1868		1920
1869		
1870		
1871		
1872		
1873		
1874		
1875		
1876	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1921
1877		1922
1878		1923
1879		1924
1880		1925
1881		1926
1882		1927
1883		1928
	Zijun Yao, Weijian Qi, Liangming Pan, Shulin Cao, Linmei Hu, Weichuan Liu, Lei Hou, and Juanzi Li. 2024. Seakr: Self-aware knowledge retrieval for adaptive retrieval augmented generation. <i>arXiv preprint arXiv:2406.19215</i> .	1929
		1930
	Fanghua Ye, Mingming Yang, Jianhui Pang, Longyue Wang, Derek F. Wong, Emine Yilmaz, Shuming Shi, and Zhaopeng Tu. 2024. Benchmarking llms via uncertainty quantification . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	1931
		1932
		1933
		1934
	Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. 2023. Compositional exemplars for in-context learning. In <i>International Conference on Machine Learning</i> , pages 39818–39833. PMLR.	1935
		1936
	Xunjian Yin, Xu Zhang, Jie Ruan, and Xiaojun Wan. 2024. Benchmarking knowledge boundary for large language models: A different perspective on model evaluation . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 2270–2286. Association for Computational Linguistics.	1937
		1938
		1939
		1940
		1941
	Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. 2023. Do large language models know what they don't know? In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 8653–8665. Association for Computational Linguistics.	
	Lang Yu, Qin Chen, Jie Zhou, and Liang He. 2024. MELO: enhancing model editing with neuron-indexed dynamic lora . In <i>Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024</i> , pages 19449–19457. AAAI Press.	
	Hanning Zhang, Shizhe Diao, Yong Lin, Yi R. Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng Ji, and Tong Zhang. 2024a. R-tuning: Instructing large language models to say 'i don't know' . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , NAACL 2024, pages 7113–7139.	
	Jiaxin Zhang, Wendi Cui, Yiran Huang, Kamalika Das, and Kumar Sricharan. 2024b. Synthetic knowledge ingestion: Towards knowledge refinement and injection for enhancing large language models . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	
	Kun Zhang, Jiali Zeng, Fandong Meng, Yuanzhuo Wang, Shiqi Sun, Long Bai, Huawei Shen, and Jie Zhou. 2024c. Tree-of-reasoning question decomposition for complex question answering with large language models . In <i>Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024</i> , pages 19560–19568. AAAI Press.	

1942	Mozhi Zhang, Mianqiu Huang, Rundong Shi, Linsen Guo, Chong Peng, Peng Yan, Yaqian Zhou, and Xipeng Qiu. 2024d. Calibrating the confidence of large language models by eliciting fidelity . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 2959–2979. Association for Computational Linguistics.	
1943		
1944		
1945		
1946		
1947		
1948		
1949		
1950	Shaolei Zhang, Tian Yu, and Yang Feng. 2024e. Truthx: Alleviating hallucinations by editing large language models in truthful space . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 8908–8949. Association for Computational Linguistics.	
1951		
1952		
1953		
1954		
1955		
1956		
1957	Tianjun Zhang, Xuezhi Wang, Denny Zhou, Dale Schuurmans, and Joseph E. Gonzalez. 2023a. TEMPERA: Test-time prompt editing via reinforcement learning . In <i>The Eleventh International Conference on Learning Representations</i> .	
1958		
1959		
1960		
1961		
1962	Tong Zhang, Peixin Qin, Yang Deng, Chen Huang, Wenqiang Lei, Junhong Liu, Dingnan Jin, Hongru Liang, and Tat-Seng Chua. 2024f. CLAMBER: A benchmark of identifying and clarifying ambiguous information needs in large language models . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024</i> , pages 10746–10766.	
1963		
1964		
1965		
1966		
1967		
1968		
1969		
1970	Yue Zhang, Leyang Cui, Wei Bi, and Shuming Shi. 2023b. Alleviating hallucinations of large language models through induced hallucinations . <i>CoRR</i> , abs/2312.15710.	
1971		
1972		
1973		
1974	Bowen Zhao, Zander Brumbaugh, Yizhong Wang, Hananeh Hajishirzi, and Noah A. Smith. 2024a. Set the clock: Temporal alignment of pretrained language models . In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 15015–15040. Association for Computational Linguistics.	
1975		
1976		
1977		
1978		
1979		
1980		
1981	Ruochen Zhao, Xingxuan Li, Shafiq Joty, Chengwei Qin, and Lidong Bing. 2023. Verify-and-edit: A knowledge-enhanced chain-of-thought framework . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 5823–5840. Association for Computational Linguistics.	
1982		
1983		
1984		
1985		
1986		
1987		
1988		
1989	Xinran Zhao, Hongming Zhang, Xiaoman Pan, Wenlin Yao, Dong Yu, Tongshuang Wu, and Jianshu Chen. 2024b. Fact-and-reflection (FaR) improves confidence calibration of large language models . In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pages 8702–8718, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.	
1990		
1991		
1992		
1993		
1994		
1995		
1996		
1997	Yukun Zhao, Lingyong Yan, Weiwei Sun, Guoliang Xing, Chong Meng, Shuaiqiang Wang, Zhicong	
1998		
	Cheng, Zhaochun Ren, and Dawei Yin. 2024c. Knowing what llms DO NOT know: A simple yet effective self-detection method . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico City, Mexico, June 16-21, 2024</i> , pages 7051–7063. Association for Computational Linguistics.	1999
		2000
		2001
		2002
		2003
		2004
		2005
		2006
		2007
	Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023a. Can we edit factual knowledge by in-context learning? In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 4862–4876. Association for Computational Linguistics.	2008
		2009
		2010
		2011
		2012
		2013
		2014
	Chuanyang Zheng, Zhengying Liu, Enze Xie, Zhenguo Li, and Yu Li. 2023b. Progressive-hint prompting improves reasoning in large language models . <i>CoRR</i> , abs/2304.09797.	2015
		2016
		2017
		2018
	Haoxiang Zhong, Chaojun Xiao, Cunchao Tu, T. Zhang, Zhiyuan Liu, and Maosong Sun. 2019. Jec-qa: A legal-domain question answering dataset . <i>ArXiv</i> , abs/1911.12011.	2019
		2020
		2021
		2022
	Zexuan Zhong, Zhengxuan Wu, Christopher D. Manning, Christopher Potts, and Danqi Chen. 2023. Mquake: Assessing knowledge editing in language models via multi-hop questions . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 15686–15702. Association for Computational Linguistics.	2023
		2024
		2025
		2026
		2027
		2028
		2029
		2030
	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V. Le, and Ed H. Chi. 2023a. Least-to-most prompting enables complex reasoning in large language models . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	2031
		2032
		2033
		2034
		2035
		2036
		2037
		2038
	Kaitlyn Zhou, Jena D. Hwang, Xiang Ren, and Maarten Sap. 2024. Relying on the unreliable: The impact of language models’ reluctance to express uncertainty . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 3623–3643. Association for Computational Linguistics.	2039
		2040
		2041
		2042
		2043
		2044
		2045
		2046
	Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2023b. Large language models are human-level prompt engineers . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	2047
		2048
		2049
		2050
		2051
		2052
	Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat seng Chua. 2021. Tat-qa: A question answering	2053
		2054
		2055

2056	benchmark on a hybrid of tabular and textual content	2107
2057	in finance. <i>ArXiv</i> , abs/2105.07624.	2108
2058	Runchuan Zhu, Zinco Jiang, Jiang Wu, Zhipeng Ma,	2109
2059	Jiahe Song, Fengshuo Bai, Dahua Lin, Lijun Wu, and	2110
2060	Conghui He. 2025. <i>Grait: Gradient-driven refusal-</i>	2111
2061	<i>aware instruction tuning for effective hallucination</i>	2112
2062	<i>mitigation. Preprint</i> , arXiv:2502.05911.	2113
2063	Runchuan Zhu, Zhipeng Ma, Jiang Wu, Junyuan Gao, Ji-	
2064	aqi Wang, Dahua Lin, and Conghui He. 2024. <i>Utilize</i>	
2065	<i>the flow before stepping into the same river twice:</i>	2114
2066	<i>Certainty represented knowledge flow for refusal-</i>	
2067	<i>aware instruction tuning. CoRR</i> , abs/2410.06913.	
2068	Thomas P. Zollo, Todd Morrill, Zhun Deng, Jake Snell,	2115
2069	Toniann Pitassi, and Richard S. Zemel. 2024. <i>Prompt</i>	2116
2070	<i>risk control: A rigorous framework for responsible</i>	2117
2071	<i>deployment of large language models. In The Twelfth</i>	2118
2072	<i>International Conference on Learning Representations,</i>	2119
2073	<i>ICLR 2024, Vienna, Austria, May 7-11, 2024.</i>	2120
2074	OpenReview.net.	
2075	A Overview	2121
2076	We begin by introducing the definition of knowl-	
2077	edge boundary, outlining three types of knowl-	
2078	edge boundaries and a four-type knowledge tax-	
2079	onomy. Following this, we describe the typical	
2080	undesired behaviors that arise from knowledge lim-	
2081	itations, emphasizing the importance of addressing	
2082	such issues. These challenges highlight the criti-	
2083	cal need for methods that can detect when LLMs	
2084	operate beyond their knowledge capabilities. To	
2085	this end, we present three distinct identification	
2086	techniques that help delineate where knowledge	
2087	gaps exist. Once these gaps are identified, various	
2088	mitigation strategies can be employed to address	
2089	the issues caused by the knowledge boundaries. Fi-	
2090	nally, we explored several significant challenges	
2091	and emerging prospects in understanding and man-	
2092	aging knowledge boundaries in LLMs. Figure 4	
2093	illustrates a comprehensive framework for manag-	
2094	ing the knowledge boundaries of LLMs, focusing	
2095	on three key components: Undesired Behaviors,	
2096	Identification of Knowledge Boundaries, and Mit-	
2097	igation Strategies.	
2098	Summary of Contribution As a survey paper,	2142
2099	our primary goal is to synthesize and analyze	2143
2100	existing research while providing new insights	2144
2101	and frameworks for understanding the knowledge	2145
2102	boundaries of LLMs. We believe our work offers	2146
2103	significant novelty in the following aspects:	2147
2104	1) Scope and Coverage	2148
2105	• Novelty in Scope: This survey covers a topic	2149
2106	or area that has not been thoroughly reviewed	2150
	before. We address an emerging field and under-	
	explored topics.	
	• Comprehensiveness: This survey provides a	
	comprehensive or up-to-date overview. Nov-	
	elty lies in including recent advancements, over-	
	looked studies, or a broader range of perspec-	
	tives.	
	2) Organization and Structure	
	• Unique Frameworks or Taxonomies: This sur-	
	vey introduces a novel taxonomy, as a new way	
	of categorizing, organizing, or analyzing the liter-	
	ature of LLM knowledge boundary studies. This	
	new taxonomy provides fresh insights into the	
	field.	
	3) Insights and Critical Analysis	
	• Original Insights: This survey provides orig-	
	inal interpretations and thought-provoking per-	
	spectives on the existing literature. For instance,	
	we provide a unique categorization of mitiga-	
	tion strategies for prompt-sensitive known knowl-	
	edge according to the process on the LLM input	
	and output sides, and also for model-specific un-	
	known knowledge based on the extent of modifi-	
	cation to LLM’s parameters.	
	• Identification of Gaps: This survey identifies	
	several underexplored areas or open problems in	
	the field.	
	4) Timeliness and Relevance	
	• As the very first survey paper on the topic of	
	LLM knowledge boundaries, our work serves as	
	a foundational resource for researchers and prac-	
	titioners. By consolidating and organizing the	
	existing literature, we provide a starting point for	
	further exploration and innovation in this critical	
	area.	
	In summary, our work contributes novelty through	
	comprehensive coverage, innovative taxonomy, and	
	original insights. We believe these contributions	
	significantly advance the understanding of LLM	
	knowledge boundaries and provide a valuable re-	
	source for the research community.	
	B Dataset	
	In the pursuit of advancing LLM capabilities	
	and understanding their boundaries in knowledge	

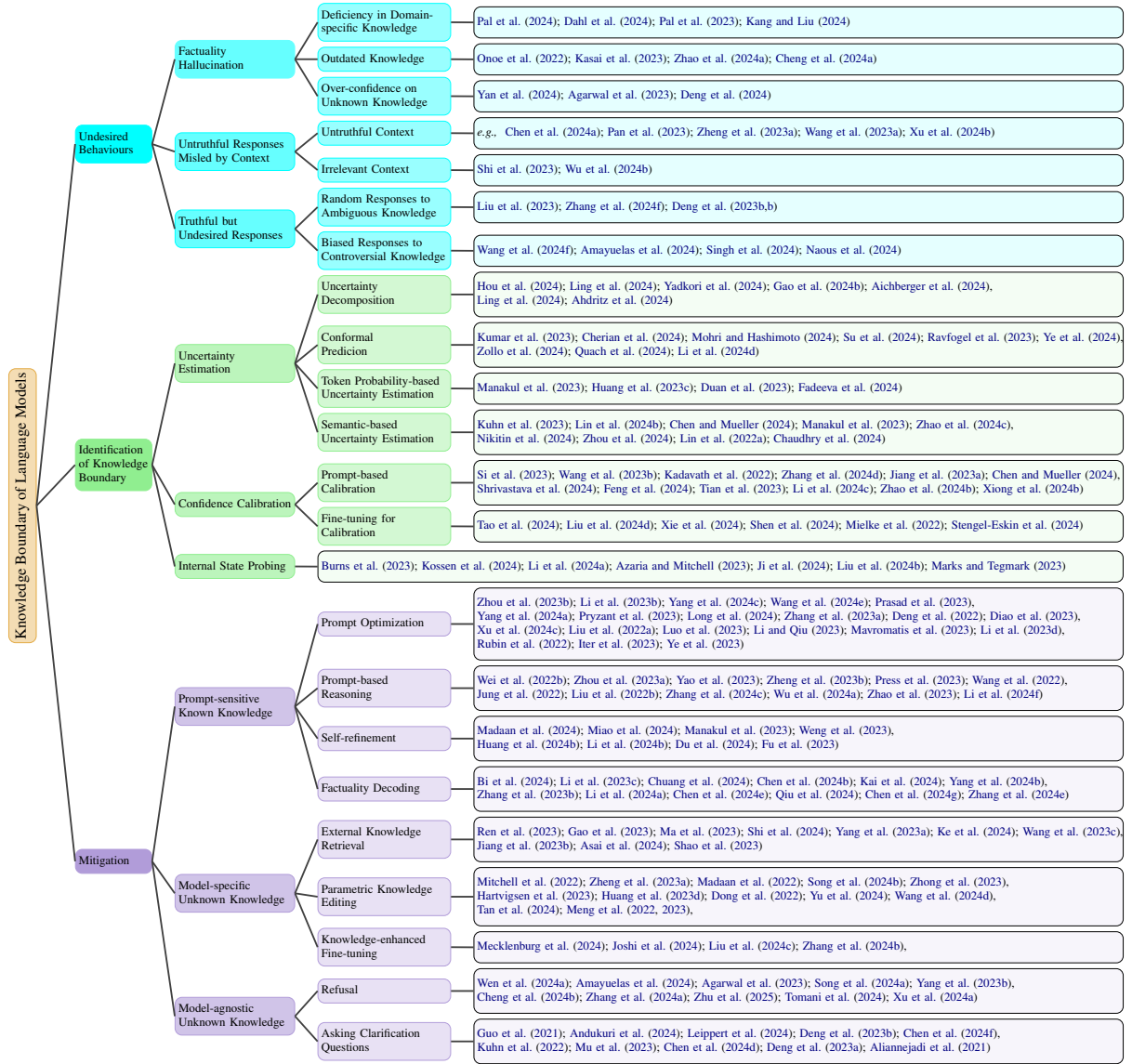


Figure 4: The main content flow and categorization of this survey.

Knowledge Category	Dataset	Reference	Size	Description
Prompt-Sensitive Known Knowledge	ProntoQA	Saparov and He (2023)	9.7k	A question-answering dataset which generates examples with chains-of-thought that describe the reasoning required to answer the questions correctly.
	2WikiMultiHopQA	Ho et al. (2020)	192,606	A multi-hop QA benchmark combining structured and unstructured data.
	MuSiQue	Trivedi et al. (2022)	25k	A multi-hop QA benchmark with 2-4 hop questions.
	HotpotQA	Yang et al. (2018)	113k	A multi-hop QA dataset requiring reasoning over two Wikipedia paragraphs, with supporting facts provided for explainability and evaluation.
	TruthfulQA	Lin et al. (2022b)	817	A benchmark across 38 categories, designed to evaluate whether language models generate truthful answers, particularly in cases prone to false beliefs.
	PARAREL	Elazar et al. (2021)	328	A dataset of English cloze-style query paraphrases for 38 relations, designed to evaluate the consistency of PLMs in handling factual knowledge across meaning-preserving input variations.
	KAAssess	Dong et al. (2024)	139k	A comprehensive assessment suite with 994,123 entities and 600 relations, designed to evaluate the factual knowledge of LLMs by estimating their ability to generate correct answers across diverse prompts compared to random chance.
	FARM	Xu et al. (2024b)	1,952	A dataset of factual questions paired with systematically generated persuasive misinformation, designed to evaluate the susceptibility of LLMs to belief manipulation through multi-turn persuasive conversations.
	Misinfo-QA	Pan et al. (2023)	3,034	A dataset designed to study the impact of misinformation on open-domain question answering (ODQA) systems by injecting synthetic misinformation passages to evaluate how QA models respond under such conditions.
Model-Specific Unknown Knowledge	Natural Questions	Kwiatkowski et al. (2019)	7,842	A large-scale dataset of real anonymized Google queries, annotated with long and short answers from Wikipedia or marked null if no answer is present.
	TopiOCQA	Adlakha et al. (2022)	3,920	An open-domain conversational dataset with information-seeking conversations featuring topic switches.
	PopQA	Mallen et al. (2022)	14k	Long-tail relation triples from WikiData are converted into QA pairs; no explicit unanswerable questions but questions are about long-tail entities.
	TriviaQA	Joshi et al. (2017)	950k	A realistic text-based question answering dataset which includes question-answer pairs from documents collected from Wikipedia and the web.
	RealtimeQA	Kasai et al. (2023)	4,356	A dynamic open-domain question-answering dataset that evaluates models based on real-time, time-sensitive questions sourced weekly from news articles.
	FreshQA	Vu et al. (2023)	600	A dynamic QA benchmark designed to evaluate LLMs on fast-changing world knowledge and debunking false premises.
	PubMedQA	Jin et al. (2019)	273.5k	A biomedical research question-answering dataset, which features questions derived from research article titles in PubMed, requiring complex reasoning and interpretation of quantitative biomedical content.
	MIRAGE	Xiong et al. (2024a)	7,663	A benchmark dataset for medical question answering, focusing on retrieving information from medical literature to answer multiple-choice medical questions, with an emphasis on zero-shot reasoning and systematic evaluation of retrieval performance.
	TAT-QA	Zhu et al. (2021)	16,552	A question-answering dataset for the financial domain, combining tabular and textual content from real financial reports.
	FinQA	Chen et al. (2021)	8,281	A question-answering dataset for the financial domain, with questions and answers crafted by financial experts, involving complex numerical reasoning over tables and text from financial reports.
	JEC-QA	Zhong et al. (2019)	26,365	A legal-domain question-answering dataset with questions sourced from the National Judicial Examination of China, covering legal concept understanding and case analysis.
	LawBench	Fei et al. (2024)	20,000	A legal reasoning evaluation benchmark designed for the Chinese legal environment, covering tasks such as legal knowledge memorization, document proofreading, case analysis, charge prediction, and legal consultation.
Model-Agnostic Unknown Knowledge	KUQ	Amayuelas et al. (2024)	6,884	A dataset designed to explore uncertainty in question-answering by focusing on questions without definitive answers.
	UnknownBench	Liu et al. (2024a)	13,319	A benchmark consisting of answerable and unanswerable questions, designed to evaluate LLMs' ability to express uncertainty and handle knowledge gaps while maintaining honesty and helpfulness.
	SelfAware	Yin et al. (2023)	2,337	A dataset containing unanswerable questions across five categories, designed to evaluate LLMs' self-knowledge by detecting uncertainty and their ability to identify limitations in their knowledge.
	QnotA	Agarwal et al. (2023)	400	A dataset featuring questions without definitive answers across five categories, paired with corresponding answerable alternatives.
	KUQP	Deng et al. (2024)	320	A dataset of known and unknown question pairs, designed to evaluate language models' ability to handle unanswerable, ambiguous, or incorrect queries.

Table 1: Representative datasets for studying the knowledge boundary of language models.

processing, various datasets have been meticulously designed and utilized. The following sections categorize these datasets into three distinct groups based on the type of knowledge they aim to verify: Prompt-Sensitive Known Knowledge, Model-Specific Unknown Knowledge, and Model-Agnostic Unknown Knowledge. A summary of these datasets can be viewed in Table 1.

Datasets for Prompt-Sensitive Known Knowledge This type of datasets mainly aim to assess the prompt-sensitive known knowledge of LLMs, requiring specific prompting strategies and decoding strategies for the LLM to fully recall and utilize such knowledge.

The first type of datasets focuses on *multi-step reasoning*, such as multi-step knowledge-based question answering datasets (e.g., 2WikiMultiHopQA (Ho et al., 2020), MuSiQue (Trivedi et al., 2022), and HotpotQA (Yang et al., 2018)) and logical reasoning datasets like ProntoQA (Saparov and He, 2023). These tasks require the LLM to achieve a step-by-step reasoning process or benefit from prompting strategies that focus on question decomposition and explicit knowledge recall.

The second type is *fact-based question answering* datasets that evaluate the LLM’s factuality, e.g., TruthfulQA (Lin et al., 2022b). In these datasets, the decoding strategy can influence how accurately knowledge is expressed (Li et al., 2024a).

The third type of datasets explicitly study the influence of *varied prompt phrasing* in LLM knowledge, including PARAREL (Elazar et al., 2021) and KAssess (Dong et al., 2024).

The fourth type involves datasets with *misleading contexts*. Wang et al. (2023a) curate queries with misleading user opinion to test LLM’s ability to defend its response. FARM (Xu et al., 2024b) contains persuasive misinformation in the dialog context to evaluate LLM’s belief change. MisinfoQA (Pan et al., 2023) includes model-generated misinformation to perturb open-domain QA.

Dataset for Model-Specific Unknown Knowledge This type of datasets can be used for assessing the model-specific unknown knowledge of LLMs, which challenges LLMs by probing their ability to handle highly specialized and temporally-sensitive information, testing their adaptive knowledge boundaries. These datasets are specifically designed to evaluate knowledge that lies outside the parametric scope of LLMs, requiring external

knowledge retrieval or new knowledge injection to generate accurate responses.

Open-domain question answering datasets form an important category. These datasets evaluate the ability of language models to answer questions across a broad range of domains, leveraging both retrieval and parametric knowledge. Representative examples include Natural Questions (Kwiatkowski et al., 2019), TopiOCQA (Adlakha et al. 2022), PopQA (Mallen et al. 2022), and TriviaQA-unfiltered (Joshi et al. 2017). These datasets often focus on queries that require world knowledge or niche details, testing the model’s capacity to combine retrieval and internalized knowledge effectively. Meanwhile, various domain-specific QA datasets can be adopted to evaluate the model-specific unknown knowledge for each specialized applications, such as medical domain (e.g., PubMedQA (Jin et al., 2019) and MIRAGE (Xiong et al., 2024a)), finance domain (e.g., TAT-QA (Zhu et al., 2021) and FinQA (Chen et al., 2021)), and legal domain (e.g., JEC-QA (Zhong et al., 2019) and LawBench (Fei et al., 2024)).

Another crucial subdomain focuses on time-sensitive datasets that test a model’s ability to generalize to out-of-distribution data. Datasets such as RealtimeQA (Kasai et al. 2023) and FreshQA (Vu et al. 2023) require language models to stay current with global events and provide accurate, up-to-date responses. These datasets evaluate the model’s capacity to adapt to evolving information and address queries that rely on recent developments.

This diverse set of datasets for studying model-sensitive unknown knowledge systematically evaluates the gaps in parametric knowledge of language models, testing their ability to retrieve, adapt, and reason with external information under various constraints.

Dataset for Model-Agnostic Unknown Knowledge As for the model-agnostic unknown knowledge, datasets such as Known-Unknown Questions (KUQ) (Amayuelas et al., 2024) and Unknown-Bench (Liu et al., 2024a) are specifically crafted to probe questions that remain unresolved or are based on uncertain future developments and incorrect assumptions. These datasets encapsulate complex scenarios including counterfactuals and ambiguities, which emphasize the current boundaries of our knowledge and the unpredictable nature of future inquiries.

Further pushing these boundaries, the SelfAware

dataset (Yin et al., 2023) explores questions that defy scientific consensus, are subjective, or philosophical, often requiring responses that extend beyond factual representation and into personal belief or theoretical speculation. Similarly, resources like QnotA (Agarwal et al., 2023) and Known-Unknown Question Pairs (KUQP) (Deng et al., 2024) challenge models with incomplete or erroneous information and speculative predictions about the future. These datasets collectively serve to test LLM’s capability in navigating the complexities of human inquiry where the answers are unknown.

C Details in Mitigation Approaches

C.1 Prompt-Sensitive Known Knowledge

Prompt Optimization. For instruction optimization, APE (Zhou et al., 2023b) leverages LLMs to automatically forward-generate and perform Monte Carlo search on the prompts, and evaluate the performance of the candidate prompts via reverse generation, which consists of n rounds. For demonstration optimization, KATE (Liu et al., 2022a) retrieve the K nearest in-context examples by the semantic similarity to the test example, measured by the embedding from an encoder model.

Prompt-based Reasoning. Chain-of-thoughts (Wei et al., 2022b) generates the step-by-step rationales followed by the answer. Tree-of-thoughts (Yao et al., 2023) improves the linear chain-of-thoughts reasoning into tree structure, each node representing a piece of thoughts, and branches represents alternative thoughts. It allows LLMs to perform various forms of reasoning steps. Progressive-hint-prompting (Zheng et al., 2023b) appends the LLM-generated answers to the prompt as hints to iteratively arrive at the correct answers.

Self-refinement. Self-refine (Madaan et al., 2024) prompts LLMs to generate feedback on its previous answer for iterative answer refinement. Self-verification (Weng et al., 2023) transforms the generated answer into abductive reasoning questions to examine the consistency with the given context. Self-correction (Huang et al., 2023a) employs an iterative initial CoT, review, and answer improvement process. MAD (Du et al., 2024) utilize multiple LLM agents to evaluate other LLMs’ answers and update their own answers until they reach a consensus.

Factuality Decoding. DoLA (Chuang et al., 2024) contrasts the logits obtained from the later layers with that obtained from the earlier layers to reduce generating factual errors. ITI (Li et al., 2024a) changes the direction of the activations towards a factuality-improving direction obtained via probing to enhance factuality during inference.

C.2 Model-specific Unknown Knowledge

External Knowledge Retrieval For pre-generation methods, HyDE (Gao et al., 2023) enhance retrieval by rewriting or expanding the user’s input to obtain more comprehensive and accurate relevant information required by the model. This approach focuses on adapting the query to improve retrieval performance. For on-demand methods, FLARE (Jiang et al., 2023b) evaluates the confidence levels in the model’s generated content and actively retrieves pertinent documents to regenerate low-confidence segments, enhancing factual accuracy.

Parametric Knowledge Editing PostEdit (Song et al., 2024b) edits the outputs of black-box LLMs while preserving data privacy and maintaining the original text style through fine-grained modifications. MELO (Yu et al., 2024) dynamically activates LoRA blocks using a neuron-indexed vector database, enabling efficient and precise updates to LLMs with minimal computational cost.

Knowledge-enhanced Fine-tuning Fact-based SFT (Mecklenburg et al., 2024) constructs a systematically covered fact-level question-answer dataset by extracting key facts from documents and generating diverse training examples, then enhances LLMs through SFT to improve their accuracy and adaptability to out-of-domain knowledge. StructTuning (Liu et al., 2024c) constructs structured domain knowledge by automatically extracting knowledge taxonomies from corpora, linking text segments to specific knowledge points for efficient model fine-tuning. Factuality alignment methods (Lin et al., 2024a; Huang and Chen, 2024a; Tian et al., 2024) is also a category of approach under this type, enhancing LLM knowledge via alignment approaches such as DPO.

C.3 Model-agnostic Unknown Knowledge

Refusal Amayuelas et al. (2024) guides LLMs to recognize “known-unknown” questions and express uncertainty in high-uncertainty scenarios, enabling them to refrain from answering questions

lacking definitive answers. R-tuning (Zhang et al., 2024a) identifies the gap between the knowledge contained in the dataset and the knowledge encapsulated in the pre-trained parameters, thereby constructing a refusal-aware dataset and training the model based on it.

Asking Clarification Questions Deng et al. (2023b) constructed a proactive prompting scheme for dialogue between users and LLMs, requiring LLMs to carefully analyze and think through the question before posing clarification questions. ACT (Chen et al., 2024d) guides the model to optimize dialogue strategies through contrastive learning in multi-turn conversations, especially when facing ambiguous user requests, enabling it to automatically recognize and ask clarification questions instead of guessing user intent or providing incorrect answers.

D Cost-effective Summarization of Representative Mitigation Techniques

We present a cost-effective comparison of representative mitigation techniques in Section 5, aiming to compare their usefulness and provide recommendations, as summarized in Table 2. This table offers a clearer and more structured comparison of these methods, helping readers better understand their relative strengths and limitations. However, directly and fairly comparing the exact performance of these methods remains challenging due to the current lack of a general and comprehensive benchmark for evaluating different mitigation approaches. Further discussions on this challenge can be found in Section 6.

From the table, we can make the following observations: (1) Prompt optimization, prompt-based reasoning, and self-refinement typically follow two main patterns: step-by-step reasoning and multi-round refinement. These fundamental approaches enhance performance, though their specific design and cost vary depending on the method used. (2) In factuality decoding, DoLA operates purely as a decoding method, whereas ITI includes a probing stage with parameter updates. This distinction can guide the choice between the two methods. (3) The main frameworks of external knowledge retrieval and parametric knowledge editing focus on integrating retrieval and inference while minimizing the cost of both components. (4) The cost of refusal and asking clarification questions methods mainly depends on whether fine-tuning on a constructed

dataset is required.

E More Challenges and Prospects

Apart from the main and general challenges and prospects discussed in §6, we further elaborate more challenges that are of great importance in real-world applications.

Generalization of Knowledge Boundary While knowledge boundary studies are often conducted in specific domains, understanding the general knowledge boundary in LLMs is vital. The internal state probing approach has been validated with a certain generalization ability (Liu et al., 2024b), but it is still an open challenge whether trained probes can generalize well across domains as a general knowledge boundary detector, fostering refusal and input clarification in open domains. Further theoretical analysis and studies are needed to identify the existence and utility of general knowledge boundaries, which may be related to fundamental theories of LLM knowledge mechanism (Wang et al., 2024b; Allen-Zhu and Li, 2024).

Unintended Side Effects Although the mitigation strategies mentioned above aim to improve the performance of LLMs, they can also introduce a range of unintended side effects that may compromise the utility and effectiveness of the model. In the following, we detail several of these effects, highlighting the challenges and potential trade-offs.

- **Over-refusal** occurs when models excessively avoid responding, even to valid queries within their knowledge boundaries. Studies like Varshney et al. (2023) show that techniques like “self-check” can make LLMs overly cautious, reducing their utility. Zhu et al. (2024) further explores this issue, identifying static and dynamic conflicts in training as key contributors.
- **Unnecessary Cost** arises when LLMs use strategies (e.g., clarifications, RAG, or self-correction) to manage queries beyond their knowledge boundaries. Although effective in avoiding undesired behaviors, these methods often consume additional time or effort, delaying responses. For instance, clarifications increase the round of interactions (Chen et al., 2024f), while RAG can introduce noise if LLMs already possess the necessary knowledge (Asai et al., 2024).

Type	Method	Training Cost	Inference Cost
Prompt Optimization	APE (Zhou et al., 2023b)	N/A	n round (prompt forward generation/monte carlo search + prompt reverse generation)
Prompt-based Reasoning	CoT (Wei et al., 2022a) PHP (Zheng et al., 2023b)	N/A N/A	step-by-step reasoning + answer n round * (step-by-step reasoning + answer)
Self-refinement	Self-correction (Huang et al., 2024b) MAD (Du et al., 2024)	N/A N/A	initial generation + review + revise n round * m agent
Factuality Decoding	DoLA (Chuang et al., 2024) ITI (Li et al., 2023a)	N/A probing the truthful direction	initial decoding + contrastive decoding step perturbed attention decoding step
External Knowledge Retrieval	HyDE (Gao et al., 2023) FLARE (Jiang et al., 2023b)	N/A N/A	hypothetical document generation + retrieval + generation n * (retrieval + generation)
Parametric Knowledge Editing	postEdit (Song et al., 2024b) MELO (Yu et al., 2024)	retrieval + SFT retrieval + SFT	generation generation
Knowledge-enhanced Fine-tuning	Fact-based SFT (Mecklenburg et al., 2024) StructTuning (Liu et al., 2024c)	fact extraction + SFT structure-aware continual pre-training + SFT	generation generation
Refusal	KUQ (Amayuelas et al., 2024) R-tuning (Zhang et al., 2024a)	SFT SFT	generation generation
Asking Clarification Questions	ProCoT (Deng et al., 2023b) ACT (Chen et al., 2024d)	N/A preference data construction + direct preference optimization	step-by-step reasoning + generation n * generation

Table 2: Cost-effective comparison of representative mitigation techniques.

Knowledge Boundary in Long-Form Language Modeling Knowledge boundaries critically impact long-form factuality, defining how well LLMs generate coherent and accurate extended responses. Unlike short-form factuality, which depends on individual fact retrieval, long-form factuality is affected by cumulative knowledge gaps, where minor errors propagate over extended discourse. Existing research (Wei et al., 2024; Min et al., 2023; Huang and Chen, 2024b) has explored evaluation and mitigation strategies, providing a lens to examine how LLMs navigate and extend their knowledge boundaries, but the interaction between knowledge boundaries and factuality degradation remains an open research area.