

# UniKnow: A Unified Framework for Reliable Language Model Behavior across Parametric and External Knowledge

Anonymous ACL submission

## Abstract

Language models often benefit from external knowledge beyond parametric knowledge. While this combination enhances performance, achieving reliable knowledge utilization remains challenging, as it requires assessing the state of each knowledge source based on the presence of relevant information. Yet, prior work on knowledge integration often overlooks this challenge by assuming ideal conditions and provides limited coverage of knowledge scenarios. To address this gap, we introduce **UniKnow**, a **Unified** framework for reliable LM behavior across parametric and external **Knowledge**. UniKnow enables controlled evaluation across knowledge scenarios such as knowledge conflict, distraction, and absence conditions that are rarely addressed together. Beyond evaluating existing methods under this setting, we extend our work by introducing UniKnow-Aware methods to support comprehensive evaluation. Experiments on UniKnow reveal that existing methods struggle to generalize across a broader range of knowledge configurations and exhibit scenario-specific biases. UniKnow thus provides a foundation for systematically exploring and improving reliability under knowledge scenarios.

## 1 Introduction

Language models (LMs), trained on large-scale corpora, exhibit the capacity to address a broad range of tasks by leveraging their pre-trained parametric knowledge (Grattafiori et al., 2024; Yang et al., 2024). However, LMs are confined to the static pre-trained knowledge and therefore struggle to handle tasks requiring information beyond this boundary, such as long-tail (Kandpal et al., 2023; Mallen et al., 2023) or time-sensitive information (Liska et al., 2022). To overcome these limitations, LMs often benefit from dynamically incorporating external knowledge, commonly through retrieval-augmented generation (RAG), thereby granting ac-



Figure 1: Four knowledge scenarios in UniKnow are defined by the boundaries of parametric and external knowledge sources. Each region illustrates the expected LM behavior for each scenario.

cess to up-to-date, task-relevant information at inference time (Chen et al., 2017; Asai et al., 2023).

The integration of parametric and contextual knowledge has broadened the capabilities of LMs, driving their application in knowledge-intensive and sensitive domains (Tsatsaronis et al., 2015; Jin et al., 2019; Dasigi et al., 2021). Consequently, the reliability of LMs has become a vital consideration (Wen et al., 2024a), with models expected to not only recognize the boundaries of their possessed knowledge but also identify when relevant information is missing. While prior work has tackled various dimensions of knowledge integration (Su et al., 2024; Yoran et al., 2023), these studies have typically remained fragmented, providing an incomplete assessment of reliability (Li et al., 2023; Cheng et al., 2024). Moreover, knowledge utilization methods developed under such narrow environments still lack validation in more realistic and compositional knowledge scenarios.

To this end, we introduce **UniKnow**, a unified framework for reliable LM behavior across parametric and external knowledge. While reliability may encompass a broader range of factors, this work focuses on the presence of *relevant* information in each parametric and external knowledge

source. Central to UniKnow is the notion of *relevance*, which we define as whether a knowledge source provides sufficient and contextually supporting information to answer a query. For example, when asked “Who is the president of the United States?”, an LM might answer “Biden” based on its parametric knowledge, while the context might refer to “Trump”—both are considered relevant.

UniKnow is designed to categorize and assess four distinct scenarios as illustrated in Figure 1: (1) Conflict, (2) Parametric-Only, (3) External-Only, and (4) Unknown. When only a single relevant source is available, the model is expected to ground its output solely in that source. Furthermore, if both sources are relevant but conflicting (1), the model should prioritize the external knowledge, as it generally offers more up-to-date and task-specific information. If neither source provides relevant knowledge (4), the model should recognize its limitations and abstain from generating hallucinations (Zhang et al., 2024a; Feng et al., 2024).

To examine how existing methods developed under partial scenario coverage generalize to UniKnow, we evaluate two naïve baselines and three existing methods with different scenario coverage. We further complement this evaluation by introducing UniKnow-Aware methods—inference-based and training-based—covering all scenarios in UniKnow by explicitly incorporating relevance-based knowledge conditions into their formulation.

Our in-depth analysis under UniKnow reveals that methods appearing reliable in individual scenarios often fail in composite scenarios requiring simultaneous consideration of both knowledge sources. We further uncover how LM behavior shifts across scenarios, highlighting biases specific to scenario types. Notably, training with UniKnow-aligned supervision significantly improves reliability. Together, these findings enable a more comprehensive understanding of LM alignment potential under UniKnow and mark a substantial step toward bridging the gap between narrow knowledge settings and a unified framework.

## 2 Related Works

**Knowledge Conflict** Parametric knowledge is inherently static, whereas external knowledge can be delivered in response to diverse circumstances. This dynamic provision often results in discrepancies between the parametric memory and the external context. Studies have examined the conflict

through the lens of external knowledge features, such as temporal shifts (Kasai et al., 2023; Dhingra et al., 2022), synthetically updated facts (Longpre et al., 2021), and contextual plausibility (Xie et al., 2023; Tan et al., 2024). Building on these findings, several approaches aim to improve external knowledge incorporation, primarily through contrastive decoding (Shi et al., 2024; Jin et al., 2024b; Yuan et al., 2024). Yet many existing approaches (Liu et al., 2024; Wang et al., 2024a; Jin et al., 2024a) still treat any mismatch between model output and context as a conflict, often neglecting whether the model had prior access to that information.

**Robustness against Irrelevance** Although external knowledge is intended to supply LM’s knowledge, in real-world scenarios (i.e. RAG), it may not always be relevant. LMs face challenges in handling irrelevant context, which often leads to performance degradation (Shen et al., 2024). RAG is particularly susceptible, as retrieval errors can introduce a misleading but plausible context (Wu et al., 2024). To mitigate this, researchers have explored methods to encourage LMs to rely on parametric knowledge when external information is irrelevant—either at inference time (Yu et al., 2024b; Park et al., 2024; Baek et al., 2023) or through training (Yoran et al., 2023; Asai et al., 2024; Xia et al., 2024; Luo et al., 2023). Despite their effectiveness at mitigating the influence of irrelevant context, these approaches entirely overlook the presence of relevant information in their parametric knowledge.

**Abstention** A growing line of work focuses on aligning LMs to abstain when appropriate—specifically when the model does not possess the relevant knowledge—to prevent hallucination and ensure reliable model behavior (Wen et al., 2025). Some approaches quantify uncertainty (Huang et al., 2025; Kuhn et al., 2023; Kadavath et al., 2022) in parametric knowledge and relabel training data accordingly to guide abstention behavior (Feng et al., 2024; Zhang et al., 2024a; Wen et al., 2024b). Recently, studies have begun to explore abstention based on the relevance of external knowledge (Wen et al., 2024a; Kim et al., 2025).

**Knowledge Frameworks** There have been efforts to unify various aspects of knowledge utilization to understand LM behaviors. Li et al. (2023) trains LMs to generate either parametric- or context-grounded responses depending on the context type, whereas Neeman et al. (2023) trains

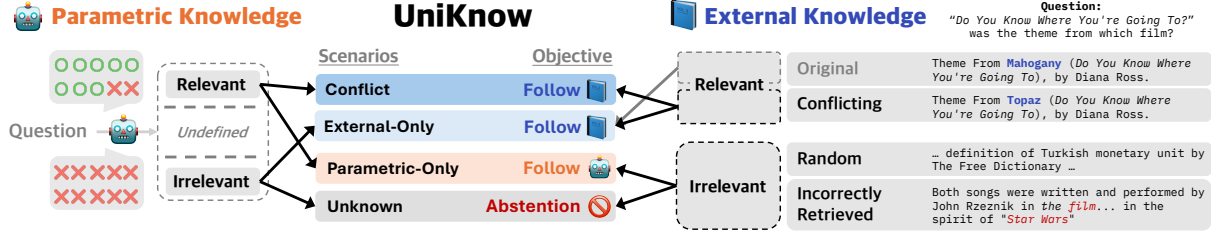


Figure 2: Overview of UniKnow.

LMs to generate both in parallel. Similar to our work, Cheng et al. (2024) proposes a benchmark to investigate whether LMs can express possessed parametric knowledge when exposed to various context types. While prior approaches have provided valuable insights into how LMs utilize knowledge, we extend this perspective with a framework centered on reliability. We offer a more comprehensive view by jointly considering the relevance of information within both parametric and external knowledge, and by addressing conflict, irrelevance, and abstention within a unified framework.

### 3 UniKnow

This work focuses on context-augmented generation in open-domain question-answering, facilitating LMs to leverage their **parametric** knowledge while simultaneously utilizing **external** knowledge to answer a given query  $q$ . This section first defines each knowledge source based on the availability of relevant information. Guided by this taxonomy, we introduce **UniKnow**, a **Unified** framework for reliable LM behavior across parametric and external **Knowledge**, covering four distinct scenarios as illustrated in Figure 2. We then describe the construction process of estimating the parametric knowledge and designing diverse context types.

#### 3.1 Definition of Knowledge Sources

**Parametric knowledge (PK)** refers to information encoded in an LM during pretraining. Since this knowledge is bound by its pretraining data, we define that relevant information resides in PK ( $\exists_{PK}$ ) if  $LM(\hat{a} | q) = a_{PK}^*$ , where  $a_{PK}^*$  denotes the answer grounded in the LM’s pretraining data (Bang et al., 2025). Still, PK remains inherently static and may not align with the most recent world knowledge.

**External knowledge (EK)** indicates any information provided at inference time as the input context. To solely evaluate the LM’s ability to utilize relevant knowledge, we assume that all provided EK is factually aligned with world knowledge. Un-

der this assumption, we assess LM behavior in both relevant ( $\exists_{EK}$ ) and irrelevant ( $\emptyset_{EK}$ ) contexts.

#### 3.2 Scenarios in UniKnow

UniKnow is designed to cover all possible scenarios regarding the presence of relevant PK and EK. This gives rise to four distinct scenarios, each reflecting real-world challenges such as conflict resolution, over-reliance, and hallucination risk. Since each challenge has its own expected behavior, we define scenario-specific expectations as follows.

- **Conflict (C):** ( $\exists_{PK}, \exists_{EK}$ ) and  $a_{PK}^* \neq a_{EK}^*$   
The conflict between knowledge sources arises when EK presents relevant information contradicting what LM knows (Xu et al., 2024b). While PK and EK may either align or conflict, we focus on the latter, allowing us to evaluate whether LMs can correctly prioritize EK.
- **External-Only (E-Only):** ( $\emptyset_{PK}, \exists_{EK}$ )  
The model lacks PK with relevant information and is expected to rely on relevant EK.
- **Parametric-Only (P-Only):** ( $\exists_{PK}, \emptyset_{EK}$ )  
The model is required to rely on its PK with relevant information and ignore irrelevant EK.
- **Unknown (U):** ( $\emptyset_{PK}, \emptyset_{EK}$ )  
Neither knowledge source is sufficient, and the model is expected to abstain from answering.

#### 3.3 Parametric Knowledge Estimation

We estimate the presence of relevant PK by assessing whether the LM is capable of generating a correct answer to a given  $q$  without access to external context. Following prior works, we assess the factual *correctness* (Zhang et al., 2024a,b; Wang et al., 2024b) and *consistency* (Kuhn et al., 2023; Huang et al., 2025; Amayuelas et al., 2024b) of the prediction utilizing its PK. We classify  $q$  as  $\exists_{EK}$  if both conditions are satisfied, and as  $\emptyset_{PK}$  otherwise.

For each  $q$ , we sample  $n$  responses using  $q$  alone:  $a_i \sim LM(a | q)$  for  $i = 1, \dots, n$ . If the proportion

of correct responses is greater than or equal to the threshold  $\tau$ , we classify  $q$  as  $\exists_{EK}$ :

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}[a_i = a_{PK}^*] \geq \tau \Rightarrow q \in \exists_{PK} \quad (1)$$

If none of the responses are correct, we assign  $q \in \emptyset_{PK}$ . Questions falling between these thresholds are considered *undefined* and excluded from scenario construction. We set  $n = 10$  and  $\tau = 0.7$  in our implementation.

### 3.4 External Knowledge Construction

To operationalize each scenario, we construct context types tailored to diverse conditions. In addition to the original context, we construct conflicting and two types of irrelevant contexts: (1) topically unrelated random contexts, and (2) incorrectly retrieved contexts with high retriever score. This allows fine-grained control over the degree of relevance, capturing challenges ranging from knowledge conflicts to misleading but plausible distractors.

**Relevant contexts** The *original* context refers to the context paired with the question-answer pair in the dataset. We derive a *conflicting* context by providing LLAMA 3 70B INSTRUCT (Grattafiori et al., 2024) with the original context and the corresponding answer to generate an alternative answer while preserving its part of speech. The original answer span is then replaced with the conflicting answer, introducing an intended conflict with the model’s PK. Note that in C scenario, we use only conflicting contexts, whereas E-Only scenario includes both original and conflicting contexts.

**Irrelevant contexts** We consider two key aspects for irrelevant context selection: the absence of the answer span (i.e., uninformative) and the potential semantic relevance (Wu et al., 2024) that may mislead the model (i.e., misleading). To capture both uninformative and misleading cases, we include two types of contexts. A *randomly* sampled context from the same dataset, topically unrelated to the question, and not containing the original answer. The *incorrectly retrieved contexts* also lack the answer but may appear topically relevant, thereby creating a false sense of relevance. We obtain these incorrectly retrieved contexts by querying a Wikipedia corpus using the CONTRIEVER-MSMARCO retriever (Izacard et al., 2022), and then select the highest-ranked context that does not contain the answer. This setting captures challenges

Methods	Conflict	E-Only	P-Only	Unknown
COIECD	✓	✓	✗	✗
RetRobust	✗	✓	✓	✗
KAFT	✓	✓	✓	✗
COIECD <sub>Prompt</sub>	✓	✓	✓	✓
LM <sub>UniKnow</sub>	✓	✓	✓	✓

Table 1: Comparison of methods based on their consideration of each UniKnow scenarios. ‘✓’ indicates the method explicitly accounts for the corresponding scenario, while ‘✗’ denotes that it does not.

in real-world RAG, where retrieval often returns plausible but irrelevant information.

## 4 Knowledge Utilization Methods

This section describes the methods used to evaluate model behavior under UniKnow. As an initial baseline, we take a **prompting** approach, instructing the model to consider the presence of knowledge sources for a reliable generation. We also perform **naïve** greedy generation with a simple QA template. We further assess three existing methods alongside two UniKnow-aware approaches. Table 1 summarizes the scenario coverage of each method. Implementation details of the approaches are presented in Appendix B.1.

### 4.1 Existing Methods

We adapt three existing methods, each designed to handle only partial scenarios of UniKnow.

**Conflict** Methods for resolving knowledge conflict aim to overwrite the model’s PK with EK. To this end, context-aware contrastive decoding approaches have been widely explored (Shi et al., 2024; Zhao et al., 2024). Among them, we utilize **COIECD**<sup>1</sup> (Yuan et al., 2024), a state-of-the-art method that amplifies the context-informed distribution when conflict arises.

**Parametric-Only RetRobust** (Yoran et al., 2023) fine-tunes the LM with augmented training data, incorporating irrelevant context alongside the original context. The goal of RetRobust is to improve its robustness against irrelevant contexts.

**Conflict and Parametric-Only** Li et al., 2023 also adopts a fine-tuning approach, **KAFT**, a knowledge-aware fine-tuning that addresses both knowledge conflict and irrelevance. Their training data includes original, conflicting, and irrelevant contexts, aiming to improve the LM’s overall ability to utilize external knowledge effectively.

<sup>1</sup>Contextual Information-Entropy Constraint Decoding



## 4.2 UniKnow-Aware Methods

To evaluate the impact of covering all UniKnow scenarios on model behavior, we construct two UniKnow-Aware methods.

**UniKnow-Aware Inference** To explicitly account for UniKnow’s scenarios during inference, we introduce **COIECD<sub>Prompt</sub>**, an extension of COIECD that additionally incorporates prompting into the decoding process. By explicitly considering all the possible scenarios, we expect COIECD<sub>Prompt</sub> to cover a broader range of cases.

**UniKnow-Aware Training** We investigate whether reliability can be improved by training LMs with supervision aligned to knowledge scenarios defined in UniKnow. We design scenario-aware training data that explicitly reflects the presence or absence of relevant information in both knowledge sources. The key lies in the scenario-aware construction of the training data.

To prepare training data, we sample a balanced set of  $q \in \exists_{PK}$  and  $q \in \emptyset_{PK}$ , as determined by the criteria in Section 3.3. As illustrated in Figure 3, each  $q$  is then paired with four types of external contexts described in Section 3.4 to cover knowledge scenarios. To maintain the LM’s ability to answer when the context contains information that matches with its PK ( $a_{PK}^* = a_{EK}^*$ ), we include the original context paired with  $q \in \exists_{PK}$  during training, although it is excluded from the C scenario analysis. For scenarios where relevant information is available—C, E-Only, and P-Only—LM<sub>UniKnow</sub> is optimized to produce the expected answer corresponding to each scenario. In the U scenario, LM<sub>UniKnow</sub> is trained to abstain by generating "unknown".

## 5 Experimental Setting

### 5.1 Implementation Details

**Datasets** UniKnow employs seven QA datasets from diverse knowledge domains: NaturalQuestions (NQ), TriviaQA, HotpotQA, SQuAD, BioASQ, TextbookQA, and RelationExtraction (RE) (Kwiatkowski et al., 2019; Joshi et al., 2017; Yang et al., 2018; Rajpurkar et al., 2016; Tsatsaronis et al., 2015; Kembhavi et al., 2017; Levy et al., 2017). We use the dataset versions curated by the Machine Reading for Question Answering (MRQA) benchmark (Fisch et al., 2019). As the impact of context length is beyond the scope of our

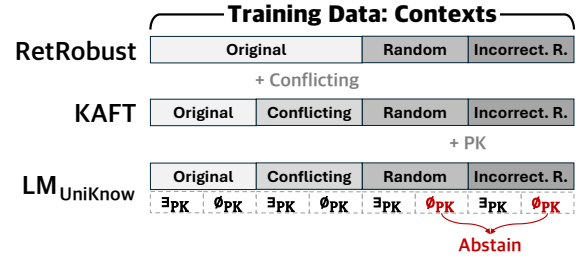


Figure 3: Comparison of training-based methods on their training data.

study, we limit the context to approximately 100 words to ensure experimental control.

**Models** We use open-source auto-regressive language models, including LLAMA2 (7B & 13B, Touvron et al., 2023), LLAMA3-8B (Grattafiori et al., 2024), MISTRAL-7B v0.3 (Jiang et al., 2023), and QWEN 2.5 (1.5B & 3B & 7B & 14B, Yang et al., 2024). Training-based methods are evaluated in a zero-shot setting, whereas inference-only methods utilize two-shot demonstrations. More details on datasets and templates are in Appendix A.

**Training Details** For a fair comparison, all training-based methods share the same settings. Utilizing the training set of NQ and TriviaQA, we randomly sample 250 questions from each of  $\exists_{PK}$  and  $\emptyset_{PK}$ , resulting in a total of 1,000 samples. As illustrated in Figure 3, we pair each  $q$  with four context types, resulting in 4,000 question-context pairs. In case of RetRobust, since it does not use conflicting contexts, we additionally sample 1,000 questions and pair them with the original context. QLoRA (Detrmers et al., 2023) is applied for efficient training. Appendix B.2 provides additional training details.

### 5.2 Evaluation Metrics

We use Exact Match (EM) to assess whether the model’s prediction aligns with the expected answer. Note that the expected answer for each scenario differs, as defined in Section 3.2. Still, it is equally important to evaluate LM behavior on samples that are *undefined* with respect to the presence of relevant PK. To reflect more realistic usage settings, we also evaluate the full samples within UniKnow and report the accuracy (Acc) and reliability (Rel<sub>y</sub>) score (Xu et al., 2024a), which captures both correctness and appropriate abstention. These are computed based on the number of correct ( $N_c$ ), incorrect ( $N_i$ ),

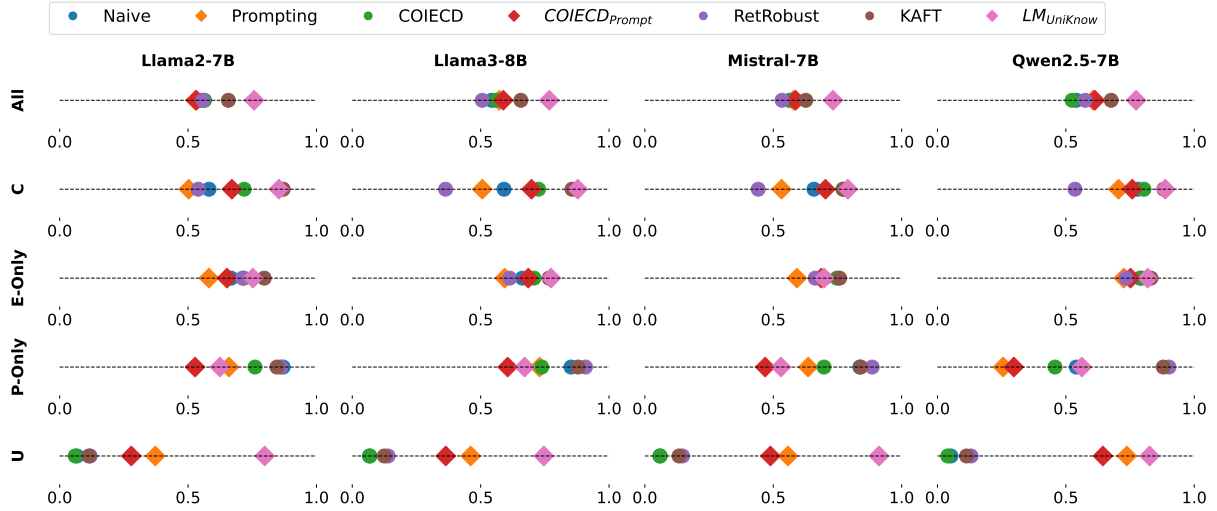


Figure 4: EM scores by scenario and model. All indicates scores averaged across all scenarios. Methods marked with diamonds incorporate abstention, while those with circle markers do not.

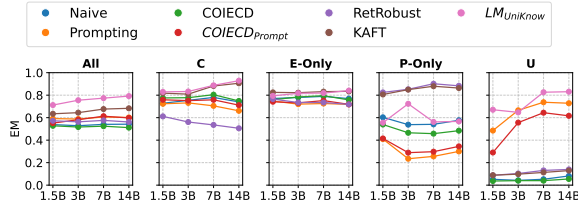


Figure 5: EM scores of QWEN models across different sizes, shown by scenario.

and abstained ( $N_a$ ) responses with EM.<sup>2</sup>

Rely score captures the balance of two components: Acc ( $\frac{N_c}{N}$ ) and truthfulness (Truth). Truth measures the proportion of responses that are either correct or abstained ( $\frac{N_c + N_a}{N}$ ), thereby ensuring that the model avoids generating incorrect outputs. To discourage excessive abstention, the answer rate ( $\text{Ans} = \frac{N_c + N_i}{N}$ ) is used as a weighting factor. Rely is computed as  $\text{Ans} \times \text{Truth} + (1 - \text{Ans}) \times \text{Acc}$ . Rely is high when LM provides correct answers and abstains appropriately, while penalizing both incorrect outputs and excessive abstention.

## 6 Results on UniKnow

### 6.1 Main Results

Figure 4 illustrates the performance across the four UniKnow scenarios, C, E-Only, P-Only, and U, and the overall averaged performance (All). To assess generalization across knowledge domains, we report EM scores averaged over all datasets, comprising two in-domain and five out-of-domain sets for training-based methods.

<sup>2</sup> $N$ : The total number of responses.

**Broader scenario coverage leads to better overall results.** LM<sub>UniKnow</sub>, which covers all scenarios, achieves the best overall performance, followed by KAFT. Other methods, designed with a subset of scenarios, lead to limited performance gains, often falling below or only marginally above naïve. Meanwhile, COIECD<sub>Prompt</sub> consistently outperforms both COIECD and Prompting in three out of four models, demonstrating the extensibility potential of existing methods.

**Resolving conflicts with known knowledge poses a greater challenge than simply incorporating new, unknown information.** Compared to C scenario, the performance points in E-Only are more tightly clustered with less variance. It demonstrates that LM behavior is influenced not only by context type itself, but also by its interaction with PK. Still, a similar trend is observed across methods in both C and E-Only scenarios. Notably, the performance drop of RetRobust is more pronounced in the C scenario than in E-Only, reflecting its limited ability to handle contradictory information effectively.

**A trade-off between answering and abstention arises under irrelevant contexts.** Methods that prioritize answerability without accounting for the presence of PK, such as COIECD, RetRobust, and KAFT, achieve strong performance in P-Only scenario. However, in U scenario, they are more likely to generate hallucinations. In contrast, methods that incorporate abstention ability, including Prompting, COIECD<sub>Prompt</sub>, and LM<sub>UniKnow</sub>, handle U with abstention behavior, but suffer in a trade-

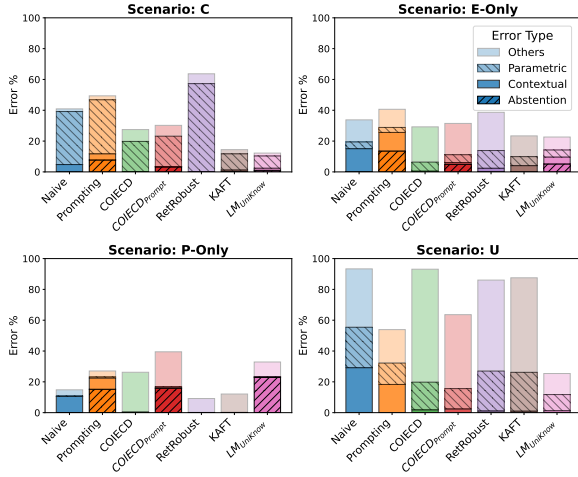


Figure 6: Stacked error type distributions across methods for each knowledge scenario. Transparency reflects error type. Evaluated using LLAMA3-8B.

off of exhibiting lower performance in P-Only. Among these,  $LM_{UniKnow}$  demonstrates the largest performance gain in U scenario, driven by its consideration of the model’s knowledge state.

**Larger LMs generally improve reliability, with distinct trends across scenarios.** Based on Figure 5, the performance in E-Only scenario remains relatively unaffected by scale, suggesting that EK utilization does not strongly benefit from larger LMs. In C and P-Only scenarios, gains depend on whether the method is explicitly trained for those conditions. By contrast, in U scenario, abstention performance improves consistently with scale, indicating that larger LMs are better at recognizing knowledge limitations and abstaining accordingly.

## 6.2 Error Analysis

Since LMs may exhibit scenario-specific biases, we analyze output errors to examine such patterns in detail. Incorrect responses are categorized into four types: contextual, parametric, false abstention, and others. *Contextual errors* occur when the model generates an incorrect response grounded on the given context. In case of relevant context, this involves extracting incorrect information; in the case of irrelevant content, the model is misled by unrelated content. *Parametric errors* refer to errors generated based on the model’s PK. In the C scenario, this reflects the model’s failure to follow the given context, exhibiting a parametric bias. *False abstention* is counted as an error in three scenarios where the model possesses at least one relevant knowledge, except U. *Other* includes incorrect re-

sponses that do not fall into the above categories. Figure 6 shows the error distribution for Llama 3 8B across the four knowledge-handling scenarios.

**Over-reliance on PK depends on the presence of PK.** In the C scenario, where the model possesses the relevant information, all methods exhibit the highest rate of parametric errors compared to other error types. In contrast, such error is much less common in E-Only scenario. Even with COIECD, which explicitly targets knowledge conflict, the rate of parametric error remains significantly higher in C than in E-Only. Unlike prior works that focus solely on controlling EK via conflicting contexts, our findings highlight that over-reliance becomes more evident when scenarios are further distinguished by the presence of PK.

**Contextual errors are rare across most methods, except for naïve approaches.** In naïve approaches, contextual errors are observed in all scenarios, particularly in E-Only and U. This indicates that when the required knowledge is absent from the model’s parametric memory, it tends to rely on the provided context but often fails to utilize it correctly (E-Only) or is misled by irrelevant information (U). In contrast, most other methods effectively mitigate context misinterpretation, as evidenced by the near absence of contextual errors.

**Abstention error occurs most frequently in P-Only scenario, while it is rare under relevant contexts.** Methods guided to abstain appropriately tend to exhibit relatively high abstention bias in P-Only. This again highlights the importance of the trade-off mitigation. Interestingly, the abstention error rate of COIECD<sub>Prompt</sub> remains comparable to that of KAFT in P-Only, but is significantly reduced in E-Only. This indicates that combining the strengths of COIECD and Prompting leads to more proper abstention behavior across scenarios.

## 7 Additional Analysis on Reliability

Figure 7 visualizes the Acc and Rely scores for each method. Despite including *undefined* samples in the evaluation, the overall trend in Rely scores remains consistent with the scenario-averaged results in UniKnow (All in Figure 2). Note that methods on the dotted line, where Acc equals Rely, limit their performance in terms of answerability.  $LM_{UniKnow}$  achieves the highest Rely, and its Acc remains comparable to methods which primarily focus on answerability. This suggests that, through

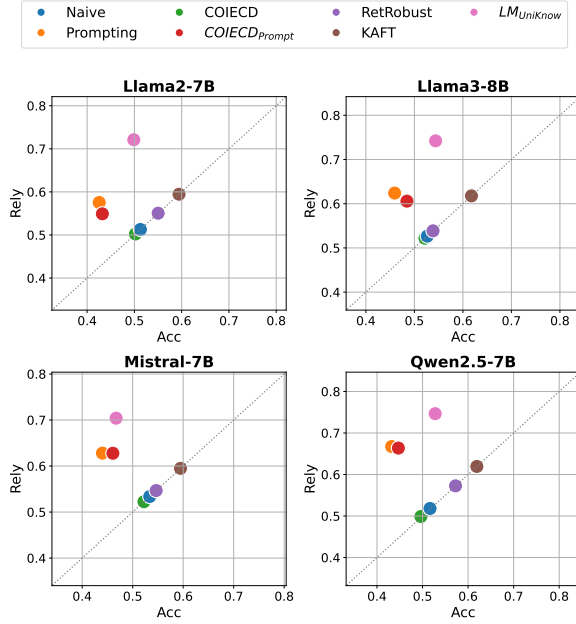


Figure 7: Acc and Rely scores across models. Each point represents a method averaged over all datasets. The dotted line indicates equal values of Acc and Rely.

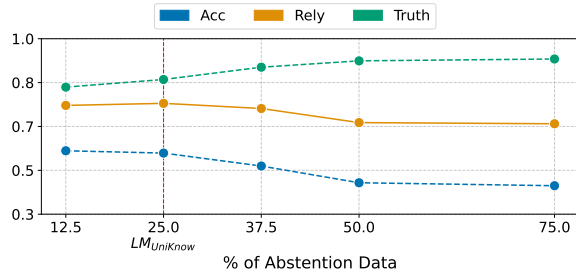


Figure 8: Effect of varying the proportion of abstention data on model performance for LLAMA3-8B. The red dashed line indicates the proportion used in  $LM_{UniKnow}$ .

alignment with UniKnow,  $LM_{UniKnow}$  effectively minimizes incorrect responses via abstention while maintaining adaptability to various scenarios.

## 7.1 Impact of Abstention Data

$LM_{UniKnow}$  allocates an equal proportion (25%) to each of the four scenarios within UniKnow. To investigate the effect of abstention supervision, we conduct an ablation study using LLAMA3-8B by varying the proportion of samples from U scenario. With a fixed number of training samples, we adjust the proportions of the remaining three scenarios equally. From Figure 8, we observe a trade-off between Acc and Truth as the proportion of abstention data increases. The lower proportions of abstention data lead to higher Acc, while higher proportions improve Truth. This reflects the inher-

Dataset Metric	TriviaQA			NQ		
	Acc	Truth	Rely	Acc	Truth	Rely
$LM_{UniKnow}$	<b>0.6915</b>	<b>0.8762</b>	<b>0.8421</b>	<b>0.5396</b>	<b>0.8161</b>	<b>0.7396</b>
−C	0.6695	0.7040	0.7028	0.4987	0.6430	0.6222
−IR	0.6872	0.7352	0.7329	0.5056	0.6410	0.6227
−C, IR	0.6836	0.7084	0.7078	0.4987	0.6406	0.6205

Table 2: Ablation study on context types in the training data for LLAMA3-8B, measuring the impact of excluding conflicting contexts (−C), incorrectly retrieved contexts (−IR), or both (−C, IR). **Bold** indicates the best.

ent trade-off between maximizing correct answer generation (Acc) and minimizing incorrect outputs through abstention (Truth). Notably, the equal allocation across the four scenarios—25% abstention data ( $LM_{UniKnow}$ )—achieves the highest Rely score in both datasets, indicating a balanced performance between answering correctly and abstaining appropriately.

## 7.2 Impact of Context Type Diversity

We conduct an ablation study in which specific types of contexts are selectively removed, while maintaining the total number of training data. We consider three ablation settings: (1) −C, which excludes conflicting contexts and replaces them with original contexts; (2) −IR, which removes incorrectly retrieved contexts and retains only randomly sampled irrelevant contexts; and (3) −C, IR, which excludes both conflicting and incorrectly retrieved contexts. These settings allow us to isolate the contribution of each context type to overall reliability. As shown in Table 2, excluding conflicting or incorrectly retrieved contexts results in a noticeable drop in Truth and Rely, while having minimal impact on Acc. These findings underscore the importance of incorporating diverse context types, reflecting those encountered in practical settings, to enhance the reliability of knowledge-handling.

## 8 Conclusion

We present UniKnow, a unified framework for evaluating LM reliability across PK and EK. By systematically defining scenarios based on knowledge relevance, UniKnow enables fine-grained analysis of LM behavior. Our experiments reveal that existing methods often struggle to jointly handle scenarios and exhibit scenario-specific biases. We show that training with UniKnow-aligned supervision improves reliability, particularly evident in U scenario. UniKnow provides a foundation for building reliable LMs in knowledge utilization.



## Limitations

**Scope of Knowledge Tasks** We primarily focus on the QA task, which provides a clear view of knowledge requirements and serves as a representative of knowledge-intensive tasks. Nevertheless, extending the scope to other tasks—such as reasoning (Xiong et al., 2024) or claim verification (Hagström et al., 2024)—is crucial, since the influence of knowledge sources may vary depending on the task. Additionally, we adopt a simplified RAG setting in which a single context is provided per query, allowing fine-grained control over context relevance and supporting targeted analysis of LM behavior. However, in real-world applications, LMs often receive multiple retrieved contexts simultaneously. This introduces new challenges, such as conflicts between external contexts (Xu et al., 2024b). Incorporating diverse tasks and extending UniKnow to support multi-context would be a valuable step toward modeling more complex and realistic RAG scenarios.

**Factuality of External knowledge** This study assumes that external knowledge is factually accurate, considering scenarios involving changed or newly emerging facts (Longpre et al., 2021; Xie et al., 2023). While this assumption enables controlled analysis, it may be strong in practice, as the quality of external knowledge depends heavily on the underlying database and retrieval system. The research area of factuality verification in external contexts using LLMs (Yu et al., 2024a; Fatahi Bayat et al., 2023) is closely related to this limitation. Exploring this aspect in conjunction with our framework could further strengthen the setting of the framework.

**Limited Strategies for UniKnow-Aware Training** Our study focuses on demonstrating the potential to UniKnow-aware supervised fine-tuning to equip LMs with a comprehensive knowledge utilization capability. Still, future work could explore alternative training techniques such as direct preference optimization or reward-based fine-tuning (Rafailov et al., 2023; Tian et al., 2024). Broadening the scope of training strategies may offer deeper insights into optimizing LM behavior across scenarios and improve the reliability. We also leave out trends beyond 14B model scale (e.g. 32B, 70B, 72B), which may further impact behavior in knowledge-intensive tasks.

## References

- Alfonso Amayuelas, Kyle Wong, Liangming Pan, Wenhu Chen, and William Yang Wang. 2024a. [Knowledge of knowledge: Exploring known-unknowns uncertainty with large language models](#). *Preprint*, arXiv:2305.13712.
- Alfonso Amayuelas, Kyle Wong, Liangming Pan, Wenhu Chen, and William Yang Wang. 2024b. [Knowledge of knowledge: Exploring known-unknowns uncertainty with large language models](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 6416–6432, Bangkok, Thailand. Association for Computational Linguistics.
- Akari Asai, Sewon Min, Zexuan Zhong, and Danqi Chen. 2023. [Retrieval-based language models and applications](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 6: Tutorial Abstracts)*, pages 41–46, Toronto, Canada. Association for Computational Linguistics.
- 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*.
- Jinheon Baek, Soyeong Jeong, Minki Kang, Jong Park, and Sung Hwang. 2023. [Knowledge-augmented language model verification](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1720–1736, Singapore. Association for Computational Linguistics.
- Yejin Bang, Ziwei Ji, Alan Schelten, Anthony Hartshorn, Tara Fowler, Cheng Zhang, Nicola Cancedda, and Pascale Fung. 2025. [Hallulens: Llm hallucination benchmark](#). *Preprint*, arXiv:2504.17550.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. [Reading Wikipedia to answer open-domain questions](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879, Vancouver, Canada. Association for Computational Linguistics.
- Sitao Cheng, Liangming Pan, Xunjian Yin, Xinyi Wang, and William Yang Wang. 2024. [Understanding the interplay between parametric and contextual knowledge for large language models](#). *Preprint*, arXiv:2410.08414.
- Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A. Smith, and Matt Gardner. 2021. [A dataset of information-seeking questions and answers anchored in research papers](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4599–4610, Online. Association for Computational Linguistics.

707	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	764
708	Luke Zettlemoyer. 2023. <a href="#">QLoRA: Efficient finetun-</a>	laume Lample, Lucile Saulnier, L��lio Renard Lavaud,	765
709	<a href="#">ing of quantized LLMs</a> . In <i>Thirty-seventh Confer-</i>	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	766
710	<i>ence on Neural Information Processing Systems</i> .	Thibaut Lavril, Thomas Wang, Timoth��e Lacroix,	767
		and William El Sayed. 2023. <i>Mistral 7b</i> . <i>Preprint</i> ,	768
		arXiv:2310.06825.	769
711	Bhuwan Dhingra, Jeremy R. Cole, Julian Martin	Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William	770
712	Eisenschlos, Daniel Gillick, Jacob Eisenstein, and	Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset	771
713	William W. Cohen. 2022. <a href="#">Time-aware language mod-</a>	for biomedical research question answering. In <i>Pro-</i>	772
714	<a href="#">els as temporal knowledge bases</a> . <i>Transactions of the</i>	<i>ceedings of the 2019 Conference on Empirical Meth-</i>	773
715	<i>Association for Computational Linguistics</i> , 10:257–	<i>ods in Natural Language Processing and the 9th In-</i>	774
716	273.	<i>ternational Joint Conference on Natural Language</i>	775
		<i>Processing (EMNLP-IJCNLP)</i> , pages 2567–2577.	776
717	Farima Fatahi Bayat, Kun Qian, Benjamin Han, Yisi	Zhuoran Jin, Pengfei Cao, Yubo Chen, Kang Liu, Xiao-	777
718	Sang, Anton Belyy, Samira Khorshidi, Fei Wu, Ihab	jian Jiang, Jiexin Xu, Li Qiuxia, and Jun Zhao. 2024a.	778
719	Ilyas, and Yunyao Li. 2023. <a href="#">FLEEK: Factual error</a>	<a href="#">Tug-of-war between knowledge: Exploring and re-</a>	779
720	<a href="#">detection and correction with evidence retrieved from</a>	<a href="#">solving knowledge conflicts in retrieval-augmented</a>	780
721	<a href="#">external knowledge</a> . In <i>Proceedings of the 2023 Con-</i>	<a href="#">language models</a> . In <i>Proceedings of the 2024 Joint</i>	781
722	<i>ference on Empirical Methods in Natural Language</i>	<i>International Conference on Computational Linguis-</i>	782
723	<i>Processing: System Demonstrations</i> , pages 124–130,	<i>tics, Language Resources and Evaluation (LREC-</i>	783
724	Singapore. Association for Computational Linguis-	<i>COLING 2024)</i> , pages 16867–16878, Torino, Italia.	784
725	tics.	ELRA and ICCL.	785
726	Shangbin Feng, Weijia Shi, Yike Wang, Wenxuan Ding,	Zhuoran Jin, Pengfei Cao, Hongbang Yuan, Yubo Chen,	786
727	Vidhisha Balachandran, and Yulia Tsvetkov. 2024.	Jiexin Xu, Huaijun Li, Xiaojian Jiang, Kang Liu,	787
728	<a href="#">Don’t hallucinate, abstain: Identifying LLM knowl-</a>	and Jun Zhao. 2024b. <a href="#">Cutting off the head ends</a>	788
729	<a href="#">edge gaps via multi-LLM collaboration</a> . In <i>Proceeed-</i>	<a href="#">the conflict: A mechanism for interpreting and mit-</a>	789
730	<i>ings of the 62nd Annual Meeting of the Association</i>	<a href="#">igating knowledge conflicts in language models</a> . In	790
731	<i>for Computational Linguistics (Volume 1: Long Pa-</i>	<i>Findings of the Association for Computational Lin-</i>	791
732	<i>pers)</i> , pages 14664–14690, Bangkok, Thailand. As-	<i>guistics: ACL 2024</i> , pages 1193–1215, Bangkok,	792
733	sociation for Computational Linguistics.	Thailand. Association for Computational Linguistics.	793
734	Adam Fisch, Alon Talmor, Robin Jia, Minjoon Seo,	Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke	794
735	Eunsol Choi, and Danqi Chen. 2019. <a href="#">Mrqa 2019</a>	Zettlemoyer. 2017. <a href="#">Triviaqa: A large scale distantly</a>	795
736	<a href="#">shared task: Evaluating generalization in reading</a>	<a href="#">supervised challenge dataset for reading comprehen-</a>	796
737	<a href="#">comprehension</a> . <i>Preprint</i> , arXiv:1910.09753.	<i>sion</i> . <i>Preprint</i> , arXiv:1705.03551.	797
738	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom	798
739	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	Henighan, Dawn Drain, Ethan Perez, Nicholas	799
740	Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-	Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli	800
741	ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh	Tran-Johnson, Scott Johnston, Sheer El-Showk,	801
742	Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-	Andy Jones, Nelson Elhage, Tristan Hume, Anna	802
743	tra, Archie Sravankumar, Artem Korenev, Arthur	Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, and	803
744	Hinsvark, and 542 others. 2024. <a href="#">The llama 3 herd of</a>	17 others. 2022. <a href="#">Language models (mostly) know</a>	804
745	<a href="#">models</a> . <i>Preprint</i> , arXiv:2407.21783.	<a href="#">what they know</a> . <i>Preprint</i> , arXiv:2207.05221.	805
746	Lovisa Hagstr��m, Sara Vera Marjanovi��, Haeun Yu,	Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric	806
747	Arnav Arora, Christina Lioma, Maria Maistro, Pepa	Wallace, and Colin Raffel. 2023. <a href="#">Large language</a>	807
748	Atanasova, and Isabelle Augenstein. 2024. <a href="#">A reality</a>	<a href="#">models struggle to learn long-tail knowledge</a> . In	808
749	<a href="#">check on context utilisation for retrieval-augmented</a>	<i>Proceedings of the 40th International Conference</i>	809
750	<a href="#">generation</a> . <i>Preprint</i> , arXiv:2412.17031.	<i>on Machine Learning</i> , volume 202 of <i>Proceedings</i>	810
751	Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming	<i>of Machine Learning Research</i> , pages 15696–15707.	811
752	Zhao, Huaming Chen, Felix Juefei-Xu, and Lei Ma.	PMLR.	812
753	2025. <a href="#">Look before you leap: An exploratory study of</a>	Jungo Kasai, Keisuke Sakaguchi, yoichi takahashi, Ro-	813
754	<a href="#">uncertainty analysis for large language models</a> . <i>IEEE</i>	nan Le Bras, Akari Asai, Xinyan Yu, Dragomir	814
755	<i>Transactions on Software Engineering</i> , 51(2):413–	Radev, Noah A Smith, Yejin Choi, and Kentaro Inui.	815
756	429.	2023. <a href="#">Realtime qa: What’s the answer right now?</a> In	816
757	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Se-	<i>Advances in Neural Information Processing Systems</i> ,	817
758	bastian Riedel, Piotr Bojanowski, Armand Joulin,	volume 36, pages 49025–49043. Curran Associates,	818
759	and Edouard Grave. 2022. <a href="#">Unsupervised dense infor-</a>	Inc.	819
760	<a href="#">mation retrieval with contrastive learning</a> . <i>Preprint</i> ,		
761	arXiv:2112.09118.		
762	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-		
763	sch, Chris Bamford, Devendra Singh Chaplot, Diego		

820	Aniruddha Kembhavi, Minjoon Seo, Dustin Schwenk,	the KNOT: Interweaving conflicting knowledge and	877
821	Jonghyun Choi, Ali Farhadi, and Hannaneh Ha-	reasoning skills in large language models. In <i>Pro-</i>	878
822	jishirzi. 2017. Are you smarter than a sixth grader?	<i>ceedings of the 2024 Joint International Conference</i>	879
823	textbook question answering for multimodal machine	<i>on Computational Linguistics, Language Resources</i>	880
824	comprehension. In <i>Proceedings of the IEEE Confer-</i>	<i>and Evaluation (LREC-COLING 2024)</i> , pages 17186–	881
825	<i>ence on Computer Vision and Pattern Recognition</i>	17204, Torino, Italia. ELRA and ICCL.	882
826	(CVPR).		
827	Hyuhng Joon Kim, Youna Kim, Sang goo Lee, and	Shayne Longpre, Kartik Perisetla, Anthony Chen,	883
828	Taeuk Kim. 2025. <a href="#">When to speak, when to abstain:</a>	Nikhil Ramesh, Chris DuBois, and Sameer Singh.	884
829	<a href="#">Contrastive decoding with abstention.</a> <i>Preprint,</i>	2021. <a href="#">Entity-based knowledge conflicts in question</a>	885
830	arXiv:2412.12527.	<a href="#">answering.</a> In <i>Proceedings of the 2021 Conference</i>	886
831	Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023.	<i>on Empirical Methods in Natural Language Process-</i>	887
832	<a href="#">Semantic uncertainty: Linguistic invariances for un-</a>	<i>ing</i> , pages 7052–7063, Online and Punta Cana, Do-	888
833	<a href="#">certainty estimation in natural language generation.</a>	minican Republic. Association for Computational	889
834	In <i>The Eleventh International Conference on Learn-</i>	Linguistics.	890
835	<i>ing Representations.</i>	I. Loshchilov and F. Hutter. 2017. Decoupled weight	891
836	Tom Kwiakowski, Jennimaria Palomaki, Olivia Red-	decay regularization. <i>International Conference on</i>	892
837	field, Michael Collins, Ankur Parikh, Chris Alberti,	<i>Learning Representations.</i>	893
838	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-	Hongyin Luo, Tianhua Zhang, Yung-Sung Chuang,	894
839	ton Lee, Kristina Toutanova, Llion Jones, Matthew	Yuan Gong, Yoon Kim, Xixin Wu, Helen Meng, and	895
840	Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob	James Glass. 2023. <a href="#">Search augmented instruction</a>	896
841	Uszkoreit, Quoc Le, and Slav Petrov. 2019. <a href="#">Natu-</a>	<a href="#">learning.</a> In <i>Findings of the Association for Computa-</i>	897
842	<a href="#">ral questions: A benchmark for question answering</a>	<i>tional Linguistics: EMNLP 2023</i> , pages 3717–3729,	898
843	<a href="#">research.</a> <i>Transactions of the Association for Compu-</i>	Singapore. Association for Computational Linguis-	899
844	<i>tational Linguistics</i> , 7:452–466.	tics.	900
845	Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettle-	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das,	901
846	moyer. 2017. <a href="#">Zero-shot relation extraction via read-</a>	Daniel Khashabi, and Hannaneh Hajishirzi. 2023.	902
847	<a href="#">ing comprehension.</a> <i>Preprint</i> , arXiv:1706.04115.	<a href="#">When not to trust language models: Investigating</a>	903
848	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan	<a href="#">effectiveness of parametric and non-parametric mem-</a>	904
849	Ghazvininejad, Abdelrahman Mohamed, Omer Levy,	<a href="#">ories.</a> In <i>Proceedings of the 61st Annual Meeting of</i>	905
850	Veselin Stoyanov, and Luke Zettlemoyer. 2020.	<i>the Association for Computational Linguistics (Vol-</i>	906
851	<a href="#">BART: Denoising sequence-to-sequence pre-training</a>	<i>ume 1: Long Papers)</i> , pages 9802–9822, Toronto,	907
852	<a href="#">for natural language generation, translation, and com-</a>	Canada. Association for Computational Linguistics.	908
853	<a href="#">prehension.</a> In <i>Proceedings of the 58th Annual Meet-</i>	Ella Neeman, Roei Aharoni, Or Honovich, Leshem	909
854	<i>ing of the Association for Computational Linguistics,</i>	Choshen, Idan Szpektor, and Omri Abend. 2023.	910
855	pages 7871–7880, Online. Association for Computa-	<a href="#">DisentQA: Disentangling parametric and contextual</a>	911
856	tional Linguistics.	<a href="#">knowledge with counterfactual question answering.</a>	912
857	Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin	In <i>Proceedings of the 61st Annual Meeting of the</i>	913
858	Wang, Michal Lukasik, Andreas Veit, Felix Yu, and	<i>Association for Computational Linguistics (Volume 1:</i>	914
859	Sanjiv Kumar. 2023. <a href="#">Large language models with</a>	<i>Long Papers)</i> , pages 10056–10070, Toronto, Canada.	915
860	<a href="#">controllable working memory.</a> In <i>Findings of the As-</i>	Association for Computational Linguistics.	916
861	<i>sociation for Computational Linguistics: ACL 2023,</i>	Seong-II Park, Seung-Woo Choi, Na-Hyun Kim, and	917
862	pages 1774–1793, Toronto, Canada. Association for	Jay-Yoon Lee. 2024. <a href="#">Enhancing robustness of</a>	918
863	Computational Linguistics.	<a href="#">retrieval-augmented language models with in-context</a>	919
864	Adam Liska, Tomas Kocisky, Elena Gribovskaya, Tay-	<a href="#">learning.</a> <i>KNOWLEDGENLP.</i>	920
865	fun Terzi, Eren Sezener, Devang Agrawal, Cyprien	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	921
866	De Masson D’Autume, Tim Scholtes, Manzil Zaheer,	pher D Manning, Stefano Ermon, and Chelsea Finn.	922
867	Susannah Young, Ellen Gilsenan-Mcmahon, Sophia	2023. <a href="#">Direct preference optimization: Your language</a>	923
868	Austin, Phil Blunsom, and Angeliki Lazaridou. 2022.	<a href="#">model is secretly a reward model.</a> In <i>Thirty-seventh</i>	924
869	<a href="#">StreamingQA: A benchmark for adaptation to new</a>	<i>Conference on Neural Information Processing Sys-</i>	925
870	<a href="#">knowledge over time in question answering models.</a>	<i>tems.</i>	926
871	In <i>Proceedings of the 39th International Conference</i>	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev,	927
872	<i>on Machine Learning</i> , volume 162 of <i>Proceedings</i>	and Percy Liang. 2016. <a href="#">Squad: 100,000+ ques-</a>	928
873	<i>of Machine Learning Research</i> , pages 13604–13622.	<a href="#">tions for machine comprehension of text.</a> <i>Preprint,</i>	929
874	PMLR.	arXiv:1606.05250.	930
875	Yantao Liu, Zijun Yao, Xin Lv, Yuchen Fan, Shulin Cao,	Xiaoyu Shen, Rexhina Blloshmi, Dawei Zhu, Jiahuan	931
876	Jifan Yu, Lei Hou, and Juanzi Li. 2024. <a href="#">Untangle</a>	Pei, and Wei Zhang. 2024. <a href="#">Assessing “implicit” re-</a>	932
		<a href="#">trieval robustness of large language models.</a> In <i>Pro-</i>	933



934	<i>ceedings of the 2024 Conference on Empirical Meth-</i>		
935	<i>ods in Natural Language Processing</i> , pages 8988–		
936	9003, Miami, Florida, USA. Association for Compu-		
937	tational Linguistics.		
938	Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia		
939	Tsvetkov, Luke Zettlemoyer, and Wen-tau Yih. 2024.		
940	<a href="#">Trusting your evidence: Hallucinate less with context-</a>		
941	<a href="#">aware decoding</a> . In <i>Proceedings of the 2024 Confer-</i>		
942	<i>ence of the North American Chapter of the Associ-</i>		
943	<i>ation for Computational Linguistics: Human Lan-</i>		
944	<i>guage Technologies (Volume 2: Short Papers)</i> , pages		
945	783–791, Mexico City, Mexico. Association for Com-		
946	putational Linguistics.		
947	Zhaochen Su, Jun Zhang, Xiaoye Qu, Tong Zhu, Yanshu		
948	Li, Jiashuo Sun, Juntao Li, Min Zhang, and Yu Cheng.		
949	2024. <a href="#">\$texttt{ConflictBank}\$: A benchmark for</a>		
950	<a href="#">evaluating the influence of knowledge conflicts in</a>		
951	<a href="#">LLMs</a> . In <i>The Thirty-eight Conference on Neural</i>		
952	<i>Information Processing Systems Datasets and Bench-</i>		
953	<i>marks Track</i> .		
954	Hexiang Tan, Fei Sun, Wanli Yang, Yuanzhuo Wang,		
955	Qi Cao, and Xueqi Cheng. 2024. <a href="#">Blinded by gen-</a>		
956	<a href="#">erated contexts: How language models merge gen-</a>		
957	<a href="#">erated and retrieved contexts when knowledge con-</a>		
958	<a href="#">flicts?</a> In <i>Proceedings of the 62nd Annual Meeting of</i>		
959	<i>the Association for Computational Linguistics (Vol-</i>		
960	<i>ume 1: Long Papers)</i> , pages 6207–6227, Bangkok,		
961	Thailand. Association for Computational Linguistics.		
962	Katherine Tian, Eric Mitchell, Huaxiu Yao, Christo-		
963	pher D Manning, and Chelsea Finn. 2024. <a href="#">Fine-</a>		
964	<a href="#">tuning language models for factuality</a> . In <i>The Twelfth</i>		
965	<i>International Conference on Learning Representa-</i>		
966	<i>tions</i> .		
967	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-		
968	bert, Amjad Almahairi, Yasmine Babaei, Nikolay		
969	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti		
970	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton		
971	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,		
972	Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 oth-		
973	ers. 2023. <a href="#">Llama 2: Open foundation and fine-tuned</a>		
974	<a href="#">chat models</a> . <i>Preprint</i> , arXiv:2307.09288.		
975	George Tsatsaronis, Georgios Balikas, Prodromos		
976	Malakasiotis, Ioannis Partalas, Matthias Zschunke,		
977	Michael R Alvers, Dirk Weissenborn, Anastasia		
978	Krithara, Sergios Petridis, Dimitris Polychronopou-		
979	los, and 1 others. 2015. An overview of the bioasq		
980	large-scale biomedical semantic indexing and ques-		
981	tion answering competition. <i>BMC bioinformatics</i> ,		
982	16:1–28.		
983	Yike Wang, Shangbin Feng, Heng Wang, Weijia Shi,		
984	Vidhisha Balachandran, Tianxing He, and Yulia		
985	Tsvetkov. 2024a. <a href="#">Resolving knowledge conflicts in</a>		
986	<a href="#">large language models</a> . In <i>First Conference on Lan-</i>		
987	<i>guage Modeling</i> .		
988	Yuxia Wang, Minghan Wang, Muhammad Arslan Man-		
989	zoor, Fei Liu, Georgi Nenkov Georgiev, Rocktim Jy-		
990	oti Das, and Preslav Nakov. 2024b. <a href="#">Factuality of</a>		
	<a href="#">large language models: A survey</a> . In <i>Proceedings</i>		991
	<i>of the 2024 Conference on Empirical Methods in</i>		992
	<i>Natural Language Processing</i> , pages 19519–19529,		993
	Miami, Florida, USA. Association for Computational		994
	Linguistics.		995
	Bingbing Wen, Bill Howe, and Lucy Lu Wang. 2024a.		996
	<a href="#">Characterizing LLM abstention behavior in science</a>		997
	<a href="#">QA with context perturbations</a> . In <i>Findings of the</i>		998
	<i>Association for Computational Linguistics: EMNLP</i>		999
	2024, pages 3437–3450, Miami, Florida, USA. Asso-		1000
	ciation for Computational Linguistics.		1001
	Bingbing Wen, Jihan Yao, Shangbin Feng, Chenjun		1002
	Xu, Yulia Tsvetkov, Bill Howe, and Lucy Lu Wang.		1003
	2024b. <a href="#">Know your limits: A survey of abstention</a>		1004
	<a href="#">in large language models</a> . <i>arXiv preprint arXiv:</i>		1005
	<i>2407.18418</i> .		1006
	Bingbing Wen, Jihan Yao, Shangbin Feng, Chenjun Xu,		1007
	Yulia Tsvetkov, Bill Howe, and Lucy Lu Wang. 2025.		1008
	<a href="#">Know your limits: A survey of abstention in large</a>		1009
	<a href="#">language models</a> . <i>Preprint</i> , arXiv:2407.18418.		1010
	Siye Wu, Jian Xie, Jiangjie Chen, Tinghui Zhu, Kai		1011
	Zhang, and Yanghua Xiao. 2024. <a href="#">How easily do</a>		1012
	<a href="#">irrelevant inputs skew the responses of large language</a>		1013
	<a href="#">models?</a> In <i>First Conference on Language Modeling</i> .		1014
	Yuan Xia, Jingbo Zhou, Zhenhui Shi, Jun Chen, and		1015
	Haifeng Huang. 2024. <a href="#">Improving retrieval aug-</a>		1016
	<a href="#">mented language model with self-reasoning</a> . <i>arXiv</i>		1017
	<i>preprint arXiv: 2407.19813</i> .		1018
	Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and		1019
	Yu Su. 2023. <a href="#">Adaptive chameleon or stubborn sloth:</a>		1020
	<a href="#">Revealing the behavior of large language models in</a>		1021
	<a href="#">knowledge conflicts</a> . <i>International Conference on</i>		1022
	<i>Learning Representations</i> .		1023
	Siheng Xiong, Ali Payani, Ramana Kompella, and Fara-		1024
	marz Fekri. 2024. <a href="#">Large language models can learn</a>		1025
	<a href="#">temporal reasoning</a> . In <i>Proceedings of the 62nd An-</i>		1026
	<i>annual Meeting of the Association for Computational</i>		1027
	<i>Linguistics (Volume 1: Long Papers)</i> , pages 10452–		1028
	10470, Bangkok, Thailand. Association for Compu-		1029
	tational Linguistics.		1030
	Hongshen Xu, Zichen Zhu, Situo Zhang, Da Ma, Shuai		1031
	Fan, Lu Chen, and Kai Yu. 2024a. <a href="#">Rejection im-</a>		1032
	<a href="#">proves reliability: Training LLMs to refuse unknown</a>		1033
	<a href="#">questions using RL from knowledge feedback</a> . In		1034
	<i>First Conference on Language Modeling</i> .		1035
	Rongwu Xu, Zehan Qi, Zhijiang Guo, Cunxiang Wang,		1036
	Hongru Wang, Yue Zhang, and Wei Xu. 2024b.		1037
	<a href="#">Knowledge conflicts for LLMs: A survey</a> . In <i>Pro-</i>		1038
	<i>ceedings of the 2024 Conference on Empirical Meth-</i>		1039
	<i>ods in Natural Language Processing</i> , pages 8541–		1040
	8565, Miami, Florida, USA. Association for Compu-		1041
	tational Linguistics.		1042
	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui,		1043
	Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu,		1044
	Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jian-		1045
	hong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang,		1046



1047	Jingren Zhou, Junyang Lin, Kai Dang, and 22 others. 2024. Qwen2.5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	<i>ciation for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 4225–4237, Mexico City, Mexico. Association for Computational Linguistics.	1104
1048			1105
1049			1106
1050	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. <a href="#">Hotpotqa: A dataset for diverse, explainable multi-hop question answering</a> . <i>Preprint</i> , arXiv:1809.09600.		1107
1051			
1052			
1053			
1054			
1055	Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. 2023. <a href="#">Making retrieval-augmented language models robust to irrelevant context</a> . <i>Preprint</i> , arXiv:2310.01558.		
1056			
1057			
1058			
1059	Tian Yu, Shaolei Zhang, and Yang Feng. 2024a. <a href="#">Truth-aware context selection: Mitigating hallucinations of large language models being misled by untruthful contexts</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 10862–10884, Bangkok, Thailand. Association for Computational Linguistics.		
1060			
1061			
1062			
1063			
1064			
1065			
1066	Wenhao Yu, Hongming Zhang, Xiaoman Pan, Peixin Cao, Kaixin Ma, Jian Li, Hongwei Wang, and Dong Yu. 2024b. <a href="#">Chain-of-note: Enhancing robustness in retrieval-augmented language models</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 14672–14685, Miami, Florida, USA. Association for Computational Linguistics.		
1067			
1068			
1069			
1070			
1071			
1072			
1073			
1074	Xiaowei Yuan, Zhao Yang, Yequan Wang, Shengping Liu, Jun Zhao, and Kang Liu. 2024. <a href="#">Discerning and resolving knowledge conflicts through adaptive decoding with contextual information-entropy constraint</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 3903–3922, Bangkok, Thailand. Association for Computational Linguistics.		
1075			
1076			
1077			
1078			
1079			
1080			
1081			
1082	Hanning Zhang, Shizhe Diao, Yong Lin, Yi Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng Ji, and Tong Zhang. 2024a. <a href="#">R-tuning: Instructing large language models to say ‘I don’t know’</a> . 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 7113–7139, Mexico City, Mexico. Association for Computational Linguistics.		
1083			
1084			
1085			
1086			
1087			
1088			
1089			
1090			
1091	Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou, Lifeng Jin, Linfeng Song, Haitao Mi, and Helen Meng. 2024b. <a href="#">Self-alignment for factuality: Mitigating hallucinations in LLMs via self-evaluation</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1946–1965, Bangkok, Thailand. Association for Computational Linguistics.		
1092			
1093			
1094			
1095			
1096			
1097			
1098			
1099	Zheng Zhao, Emilio Monti, Jens Lehmann, and Haytham Assem. 2024. <a href="#">Enhancing contextual understanding in large language models through contrastive decoding</a> . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Asso-</i>		
1100			
1101			
1102			
1103			

## Appendix

### A UniKnow: Implementation Details

#### A.1 Datasets

The total number of samples for each dataset is in Table 3. Each sample includes a question, original answer, conflicting answer, and four types of context: original, conflicting, random, and incorrectly retrieved contexts. We provide a detailed description of the datasets used in our study below.

#### NaturalQuestions (Kwiatkowski et al., 2019)

Questions consist of real queries issued to the Google search engine. From a Wikipedia page from the top 5 search results, annotators select a long answer containing enough information to completely infer the answer to the question, and a short answer that comprises the actual answer. The long answer becomes the context matched with the question, while the short answer is used as the answer.

**TriviaQA (Joshi et al., 2017)** Question-answer pairs are authored by trivia enthusiasts and independently gathered evidence documents that provide high quality supervision for answering the questions. The web version of TriviaQA is used, where the contexts are retrieved from the results of a Bing search query.

**HotpotQA (Yang et al., 2018)** Questions are diverse and not constrained to any pre-existing knowledge base. Multi-hop reasoning is required to solve the questions. Paragraphs that provide supporting facts required for reasoning, are given along with the question. In the original setting, additional distractor paragraphs are augmented in order to increase the difficulty of inference. However, these distractor paragraphs are not used in this setting.

**SQuAD (Rajpurkar et al., 2016)** Paragraphs from Wikipedia are presented to crowdworkers, and they are asked to write questions that entail extractive answers. The answer to each question is a segment of text from the corresponding reading passage. To remove the uncertainty that excessively long paragraphs bring, QA pairs that do not align with the first 800 tokens are discarded in this setting.

**BioASQ (Tsatsaronis et al., 2015)** BioASQ is a challenge that assesses the ability of systems to semantically index large numbers of biomedical scientific articles and return concise answers to given natural language questions. Each question

Dataset	Train	Test
NQ	83,787	3,994
TriviaQA	61,177	7,712
HotpotQA	-	4,760
SQuAD	-	7,918
Bioasq	-	697
TextbookQA	-	1,056
RelationExtraction	-	1,974
<b>Total</b>	144,964	28,111

Table 3: Number of samples for each dataset.

Answer the following questions:

<few-shots>

Question: <question>

Answer:

Table 4: Template used in closed-book generation.

is linked to multiple related science articles. The full abstract of each linked article is used as an individual context. Abstracts that do not exactly contain the answer are discarded.

**TextbookQA (Kembhavi et al., 2017)** TextbookQA aims at answering multimodal questions when given a context in formats of text, diagrams and images. This dataset is collected from lessons from middle school Life Science, Earth Science, and Physical Science textbooks. Questions that are accompanied with a diagram and "True or False" questions are not used in this setting.

**RelationExtraction (Levy et al., 2017)** Given labeled slot-filling examples, relations between entities are transformed into QA pairs using templates. Multiple templates for each type of relation are utilized. The zero-shot benchmark split of this dataset, which showed that generalization to unseen relations is possible at lower accuracy levels, is used.

#### A.2 Predefined Abstention Words

The predefined abstain words (Amayuelas et al., 2024a) used in evaluations are: [ 'unanswerable', 'unknown', 'no known', 'not known', 'do not know', 'uncertain', 'unclear', 'no scientific evidence', 'no definitive answer', 'no right answer', 'no concrete answer', 'no public information', 'debate', 'impossible to know', 'impossible to answer', 'difficult to predict', 'not sure', 'irrelevant', 'not relevant']

---

Answer the following questions:  
<few-shots>  
Context: <context>  
Question: <question>  
Answer:

---

Table 5: Template for the naïve open-book generation.

---

Answer an entity of the same type as the given keyword. Please note that the keyword is from the given context, and consider the part of speech of the keyword inside the context. You should not give a synonym or alias of the given keyword. The entity and given keyword must have different meanings. Only answer the entity itself without any extra phrases.  
<few-shots>  
Keyword: <original-answer>  
Context: <context>  
Answer:

---

Table 6: Template used when instructing the model to generate a conflicting answer, given the original answer and context.

### A.3 Details on UniKnow Construction

To ensure context informativeness and maintain experimental controllability, we have processed the original contexts from the MRQA benchmark by limiting their length and ensuring that the ground-truth answer span is always included. For each occurrence span of the ground-truth answer in the raw context, we take a 100-word portion surrounding that span and consider it a candidate context. We then compute the NLI (BART-LARGE, Lewis et al., 2020) score between the question-answer pair and each candidate context, and select the context with the highest NLI score as the original context.

To generate conflicting answers, Template 6 is employed. For retrieved-uninformative contexts, a Wikipedia dump from December 2018 is used as a database. Each context is chunked into 100 words. As a retriever model, CONTRIEVER-MSMARCO (Izacard et al., 2022) is utilized. The number of samples per scenario and model is provided in Table 10. Template 4 is used to perform closed-book generation for estimating the presence of parametric knowledge.

## B Knowledge Utilization Methods

### B.1 Details on Methods

For naïve open-book generation, Template 5 is used. The instruction template used in the prompting approach is in Template 7.

---

Answer the following questions. The context may or may not be helpful. If the context is unhelpful and you are not knowledgeable about the question, it is appropriate to say, "<UNKNOWN>".  
<few-shots>  
Context: <context>  
Question: <question>  
Answer:

---

Table 7: Instruction for LMs to abstain if unknown.

---

Answer the following questions. If you are not knowledgeable about the question, it is appropriate to say, "<UNKNOWN>".  
<few-shots>  
Question: <question>  
Answer:

---

Table 8: Instruction used in COIECD<sub>Prompt</sub> for LMs to abstain if unknown under closed-book generation.

**COIECD** For COIECD, which requires two hyperparameters, we adopt the values reported in the original paper ( $\alpha = 1.0$  and  $\lambda = 0.25$ ), as Yuan et al. (2024) shows that these values generalize well across models and datasets.

**COIECD<sub>Prompt</sub>** In COIECD<sub>Prompt</sub>, we use Template 7 for input with context and Template 8 for input without context.

**KAFT** Unlike Li et al., 2023, which treats the parametric answers as gold-standard for irrelevant contexts, we use the original answer to ensure fair evaluation in the U scenario.

### B.2 Hyperparameters for Training

We use the same setting for every training-based approach, including RetRobust, KAFT, and LM<sub>UniKnow</sub>. We train the model for three epochs using the AdamW (Loshchilov and Hutter, 2017) optimizer with a learning rate of 0.0001 and a batch size of 16. For efficient fine-tuning, we employ QLoRA (Dettmers et al., 2023) with rank  $r=4$  and  $\alpha=16$ . Training is conducted on two NVIDIA RTX A6000.

## C Additional Results

In this section, we provide exact values of figures and additional results for models not included in Section 6.1 and Section 7.

**Main Results** The EM scores corresponding to Figure 4 are provided in Table 11. Also, Figure 9 visualizes the EM scores of LLAMA2 7B and 13B

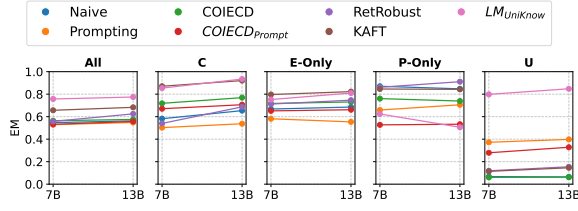


Figure 9: EM scores of LLAMA2 models across different sizes, averaged over all datasets within UniKnow.

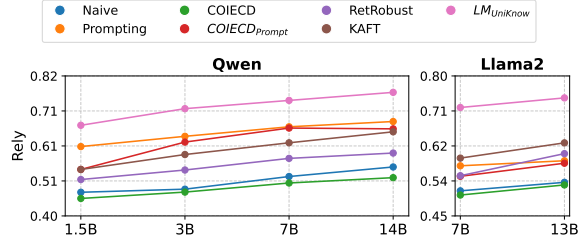


Figure 10: Rely scores of QWEN and LLAMA2 across model sizes.

across different knowledge scenarios. Figure 10 illustrates the impact of model scale with Rely metric for LLAMA2 and QWEN2.5.

The exact values for the Acc and Rely scores presented in Figure 7 are listed in Table 12 per dataset. While Figure 7 presents overall trends averaged across all datasets, Figure 11 and Figure 12 break down the results by in-domain and out-of-domain datasets, respectively. They further highlight that the overall trend across methods holds consistently and generalizes well to out-of-domain settings.

**Error Analysis** We present the error type distribution for each knowledge scenario across different models. Results for LLAMA2-7B, MISTRAL-7B, and QWEN2.5-7B are shown in Table 13, Table 14, and Table 15, respectively.

**Ablation Study** Figure 16 shows the effect of varying the proportion of abstention data on the performance across datasets. These results align with the averaged trend discussed in Section 7.1, confirming that the observed pattern holds consistently across datasets. Table 9 shows the impact of context type diversity on additional datasets beyond those reported in Table 2.

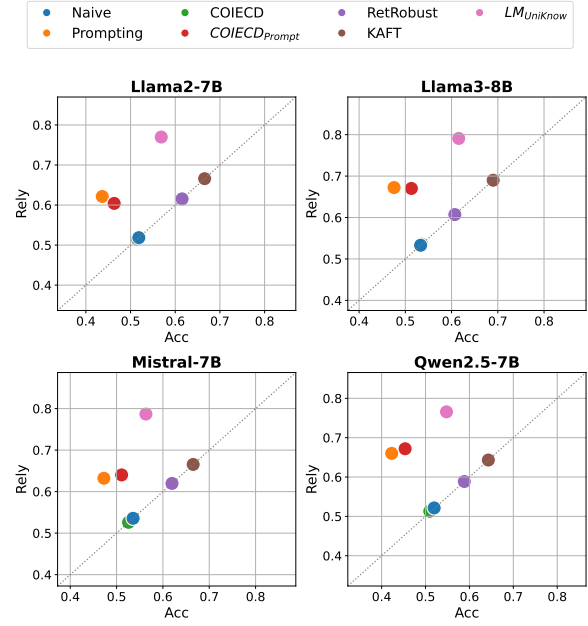


Figure 11: Acc and Rely scores averaged over in-domain datasets. Each point represents a method averaged over all datasets. The dotted line indicates equal values of Acc and Rely.

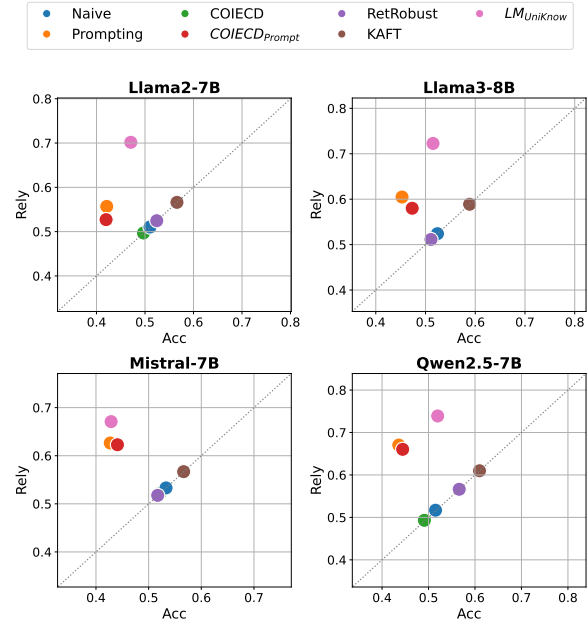


Figure 12: Acc and Rely scores averaged over out-of-domain datasets. Each point represents a method averaged over all datasets. The dotted line indicates equal values of Acc and Rely.



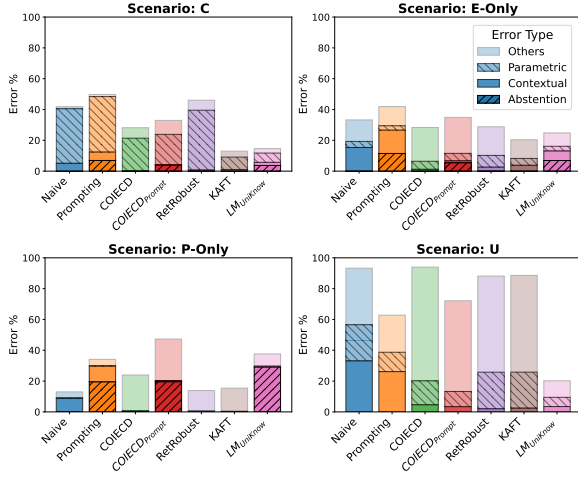


Figure 13: Stacked error type distributions across methods for each knowledge scenario. Transparency reflects error type. Evaluated using LLAMA2-7B.

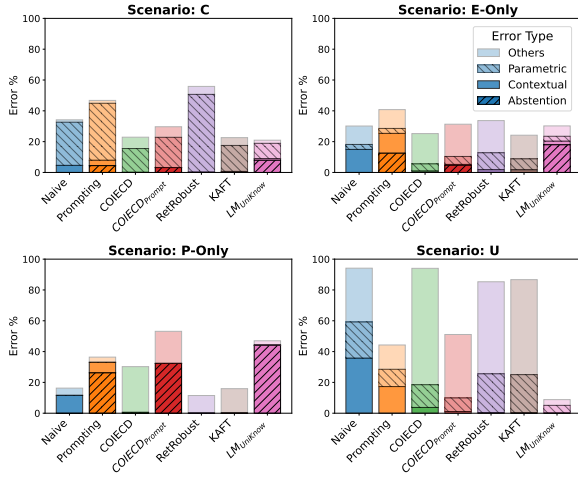


Figure 14: Stacked error type distributions across methods for each knowledge scenario. Transparency reflects error type. Evaluated using MISTRAL-7B.

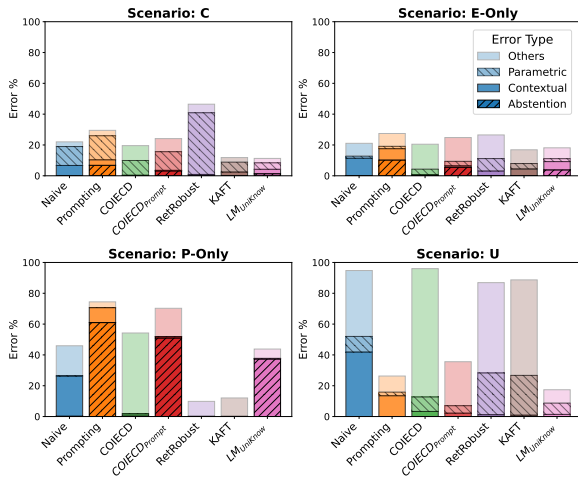


Figure 15: Stacked error type distributions across methods for each knowledge scenario. Transparency reflects error type. Evaluated using QWEN2.5-7B.

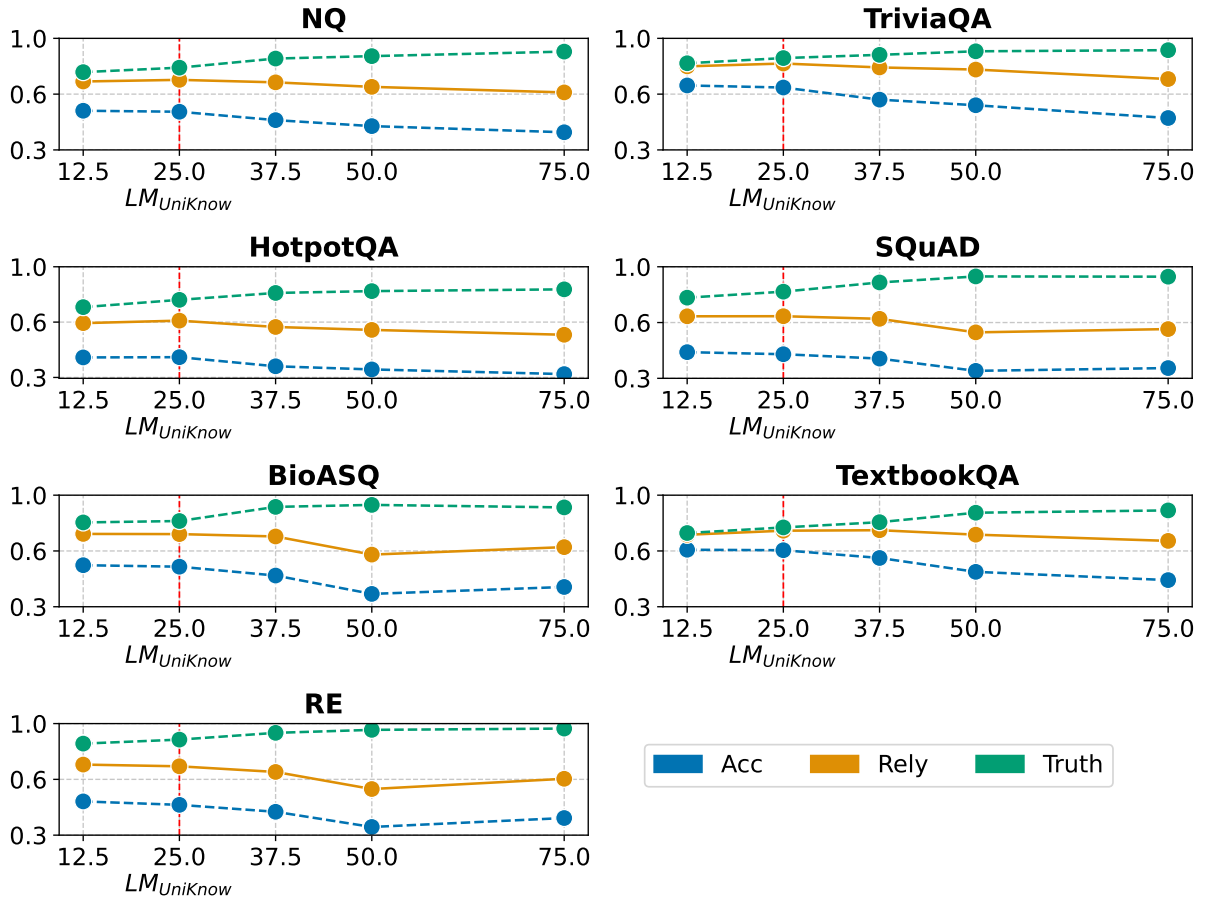


Figure 16: Effect of varying the proportion of abstention data on model performance for LLAMA3-8B for each dataset. The red dashed line indicates the proportion used in  $LM_{UniKnow}$ .

Dataset Metric	HotpotQA			BioASQ			SQuAD			TextbookQA			RE		
	Acc	Truth	Rely	Acc	Truth	Rely	Acc	Truth	Rely	Acc	Truth	Rely	Acc	Truth	Rely
LM <sub>UniKnow</sub>	<b>0.4282</b>	<b>0.7908</b>	<b>0.6593</b>	0.5513	<b>0.8379</b>	<b>0.7557</b>	0.4513	<b>0.8438</b>	<b>0.6897</b>	<b>0.6546</b>	<b>0.7976</b>	<b>0.7771</b>	0.4901	<b>0.8994</b>	<b>0.7319</b>
−C	0.4506	0.4984	0.4961	<b>0.5821</b>	0.6449	0.6410	0.4970	0.5715	0.5659	0.6044	0.6416	0.6402	0.5591	0.7247	0.6973
−IR	0.4402	0.5030	0.4990	0.5760	0.6080	0.6069	<b>0.5129</b>	0.5565	0.5546	0.5978	0.6089	0.6088	0.5678	0.6331	0.6288
−C, IR	0.4579	0.4960	0.4946	0.5918	0.5940	0.5940	0.5063	0.5224	0.5221	0.6108	0.6158	0.6157	<b>0.5784</b>	0.6312	0.6284

Table 9: Ablation study on context types in the training data for LLAMA3-8B, measuring the impact of excluding conflicting contexts (−C), incorrectly retrieved contexts (−IR), or both (−C, IR). **Bold** indicates the best.

Model	Scenario (↓)	NQ	TriviaQA	HotpotQA	SQuAD	BioASQ	TextbookQA	RE
LLAMA2-7B	C	221	2,442	160	303	74	175	145
	P-Only	442	4,884	320	606	148	350	290
	E-Only	5,090	3,676	6,878	11,088	626	694	2,422
	U	5,090	3,676	6,878	11,088	626	694	2,422
LLAMA2-13B	C	361	3,050	299	431	74	191	207
	P-Only	722	6,100	598	862	148	382	414
	E-Only	4,556	2,812	6,514	10,480	604	632	2,306
	U	4,556	2,812	6,514	10,480	604	632	2,306
LLAMA3-8B	C	273	3,231	317	462	101	193	233
	P-Only	546	6,462	634	924	202	386	466
	E-Only	4,766	3,076	6,444	10,360	448	580	2,150
	U	4,766	3,076	6,444	10,360	448	580	2,150
MISTRAL-7B	C	326	3,282	302	473	116	220	197
	P-Only	652	6,564	604	946	232	440	394
	E-Only	4,756	3,196	6,530	10,656	494	628	2,462
	U	4,756	3,196	6,530	10,656	494	628	2,462
QWEN-1.5B	C	119	1,011	80	157	59	158	78
	P-Only	238	2,022	160	314	118	316	156
	E-Only	6,202	9,246	7,774	12,292	856	802	2,964
	U	6,202	9,246	7,774	12,292	856	802	2,964
QWEN-3B	C	188	1,472	167	270	92	184	118
	P-Only	376	2,944	334	540	184	368	236
	E-Only	5,624	7,266	7,254	11,584	580	626	2,722
	U	5,624	7,266	7,254	11,584	580	626	2,722
QWEN-7B	C	315	2,485	231	401	167	282	187
	P-Only	630	4,970	462	802	334	564	374
	E-Only	5,068	5,458	6,924	10,694	422	502	2,460
	U	5,068	5,458	6,924	10,694	422	502	2,460
QWEN-14B	C	334	3,284	363	633	202	303	233
	P-Only	668	6,568	726	1,266	404	606	466
	E-Only	4,692	3,808	6,328	9,630	316	502	2,254
	U	4,692	3,808	6,328	9,630	316	502	2,254

Table 10: Number of samples in each scenario.

Scenario	Method ( $\downarrow$ )	LLAMA2-7B	LLAMA2-13B	LLAMA3-8B	MISTRAL-7B	QWEN-1.5B	QWEN-3B	QWEN-7B	QWEN-14B
All	Naïve	.5467	.5628	.5430	.5632	.5384	.5284	.5406	.5419
	Prompting	.5288	.5486	.5727	.5795	.5916	.5880	.6059	.6019
	COIECD	.5642	.5753	.5600	.5691	.5276	.5168	.5243	.5114
	COIECD <sub>Prompt</sub>	.5321	.5572	.5881	.5870	.5516	.5819	.6130	.5976
	RetRobust	.5580	.6249	.5061	.5342	.5723	.5629	.5757	.5617
	KAFT	<u>.6568</u>	<u>.6829</u>	<u>.6565</u>	<u>.6265</u>	<u>.6344</u>	<u>.6445</u>	<u>.6764</u>	<u>.6831</u>
	LM <sub>UniKnow</sub>	<b>.7571</b>	<b>.7745</b>	<b>.7672</b>	<b>.7326</b>	<b>.7120</b>	<b>.7551</b>	<b>.7735</b>	<b>.7918</b>
C	Naïve	.5817	.6538	.5911	.6585	.7280	.7538	.7799	.7400
	Prompting	.5026	.5373	.5064	.5324	.7234	.7314	.7051	.6610
	COIECD	.7185	.7691	.7254	.7711	.7754	.7775	.8043	.7487
	COIECD <sub>Prompt</sub>	.6707	.7061	.6979	.7033	.7591	.7515	.7587	.7112
	RetRobust	.5398	.6873	.3633	.4415	.6116	.5613	.5355	.5057
	KAFT	<b>.8706</b>	<u>.9207</u>	<u>.8556</u>	<u>.7741</u>	<u>.8181</u>	<u>.8091</u>	<u>.8818</u>	<u>.9058</u>
	LM <sub>UniKnow</sub>	<u>.8539</u>	<b>.9343</b>	<b>.8780</b>	<b>.7906</b>	<b>.8290</b>	<b>.8327</b>	<b>.8873</b>	<b>.9287</b>
P-Only	Naïve	<b>.8703</b>	<u>.8474</u>	.8518	.8371	.6031	.5379	.5407	.5768
	Prompting	.6591	.7056	.7295	.6360	.4098	.2348	.2557	.3000
	COIECD	.7606	.7388	.7379	.6977	.5391	.4656	.4578	.4842
	COIECD <sub>Prompt</sub>	.5272	.5329	.6051	.4685	.4152	.2888	.2975	.3447
	RetRobust	<u>.8612</u>	<b>.9109</b>	<b>.9081</b>	<b>.8851</b>	<b>.8245</b>	<b>.8522</b>	<b>.9015</b>	<b>.8831</b>
	KAFT	.8459	.8429	<u>.8792</u>	<u>.8407</u>	<u>.8060</u>	<u>.8493</u>	<u>.8795</u>	<u>.8633</u>
	LM <sub>UniKnow</sub>	.6239	.5059	.6713	.5299	.5549	.7233	.5620	.5668
E-Only	Naïve	.6677	.6855	.6623	.6987	.7704	.7795	.7893	.7694
	Prompting	.5813	.5536	.5937	.5923	.7480	.7203	.7255	.7184
	COIECD	.7171	.7309	.7077	<u>.7478</u>	.7594	.7841	.7952	.7580
	COIECD <sub>Prompt</sub>	.6514	.6617	.6854	.6869	.7416	.7314	.7518	.7182
	RetRobust	.7122	.7467	.6132	.6634	.7654	.7346	.7349	.7183
	KAFT	<b>.7967</b>	<b>.8225</b>	<u>.7663</u>	<b>.7583</b>	<b>.8248</b>	<b>.8221</b>	<b>.8314</b>	<u>.8341</u>
	LM <sub>UniKnow</sub>	<u>.7523</u>	<u>.8099</u>	<b>.7737</b>	.6979	<u>.7935</u>	<u>.8171</u>	<u>.8186</u>	<b>.8408</b>
U	Naïve	.0674	.0644	.0668	.0587	.0519	.0426	.0523	.0816
	Prompting	<u>.3724</u>	<u>.3980</u>	<u>.4611</u>	<u>.5572</u>	<u>.4852</u>	<b>.6654</b>	<u>.7371</u>	<u>.7283</u>
	COIECD	.0606	.0623	.0690	.0597	.0366	.0400	.0399	.0548
	COIECD <sub>Prompt</sub>	.2790	.3282	.3641	.4891	.2904	.5560	.6442	.6164
	RetRobust	.1189	.1546	.1396	.1469	.0877	.1034	.1310	.1396
	KAFT	.1139	.1456	.1248	.1331	.0885	.0977	.1131	.1293
	LM <sub>UniKnow</sub>	<b>.7984</b>	<b>.8479</b>	<b>.7460</b>	<b>.9118</b>	<b>.6705</b>	<u>.6472</u>	<b>.8262</b>	<b>.8308</b>

Table 11: The exact value of EM score for each scenario, across models. **Bold** indicates the best, and the underline indicates the second best.



Model	Method (↓)	NQ		TriviaQA		HotpotQA		SQuAD		BioASQ		TextbookQA		RE	
		Acc	Rely	Acc	Rely	Acc	Rely	Acc	Rely	Acc	Rely	Acc	Rely	Acc	Rely
LLAMA2-7B	Naïve	.4177	.4177	.6194	.6194	.4342	.4342	<u>.4856</u>	.4859	<u>.5402</u>	.5402	.5604	.5604	.5313	.5313
	Prompting	.3309	.5665	.5425	.6762	.3591	<u>.4675</u>	.3748	.5134	.3849	.5776	.4799	.6318	.5067	<u>.5944</u>
	COIECD	.4328	.4328	.5982	.5983	<u>.4355</u>	.4356	.4818	.4822	.5147	.5147	.5284	.5284	.5234	.5236
	COIECD <sub>Prompt</sub>	.3845	.5620	.5421	.6463	.3643	.4316	.3906	<u>.5215</u>	.3630	.5487	.4633	.5795	.5172	.5540
	RetRobust	<u>.5587</u>	.5588	<u>.6719</u>	.6720	.4277	.4277	.4715	.4722	.5319	.5330	<u>.6241</u>	.6241	<u>.5660</u>	.5661
	KAFT	<b>.5990</b>	<u>.5991</u>	<b>.7327</b>	<u>.7327</u>	<b>.4485</b>	.4486	<b>.5169</b>	.5176	<b>.5900</b>	<u>.5904</u>	<b>.6870</b>	<u>.6870</u>	<b>.5869</b>	.5869
	LM <sub>UniKnow</sub>	.5167	<b>.7236</b>	.6207	<b>.8160</b>	.3817	<b>.6219</b>	.4401	<b>.6844</b>	.4982	<b>.7232</b>	.5642	<b>.7628</b>	.4695	<b>.7160</b>
LLAMA2-13B	Naïve	.4474	.4475	.6556	.6556	.4503	.4503	.5062	.5064	.5674	.5674	.5691	.5691	.5412	.5415
	Prompting	.3678	.5357	.5649	.7067	.3993	.4528	.4148	.6083	.2991	.5487	.5208	.6164	.4933	<u>.6471</u>
	COIECD	.4594	.4594	.6361	.6362	.4509	.4510	.4959	.4961	.5739	.5739	.5533	.5535	.5222	.5223
	COIECD <sub>Prompt</sub>	.4322	.5736	.6023	.6539	.3995	.4654	.4378	<u>.6172</u>	.3311	.5752	.4979	.5793	.4829	.6078
	RetRobust	<u>.6259</u>	.6262	<u>.7461</u>	.7462	<u>.4752</u>	.4753	<u>.5099</u>	.5107	<u>.5925</u>	.5925	<u>.6889</u>	.6892	<u>.5997</u>	.6001
	KAFT	<b>.6445</b>	<u>.6446</u>	<b>.7890</b>	<u>.7890</u>	<b>.4934</b>	<u>.4936</u>	<b>.5448</b>	.5457	<b>.6234</b>	<u>.6248</u>	<b>.7294</b>	<u>.7294</u>	<b>.6016</b>	.6017
	LM <sub>UniKnow</sub>	.5416	<b>.7510</b>	.6348	<b>.8351</b>	.4143	<b>.6589</b>	.4708	<b>.7145</b>	.4598	<b>.7051</b>	.5859	<b>.7963</b>	.5138	<b>.7548</b>
LLAMA3-8B	Naïve	.4443	.4444	.6218	.6218	.4529	.4529	<u>.4943</u>	.4944	.5656	.5656	.5627	.5627	.5447	.5447
	Prompting	.4200	<u>.6347</u>	.5312	.7100	.3590	<u>.4936</u>	.4209	.6063	.4914	.6454	.4934	.6173	.4994	<u>.6613</u>
	COIECD	.4724	.4726	.5984	.5984	<u>.4534</u>	.4537	.4893	.4896	<u>.5857</u>	.5864	.5301	.5301	.5230	.5230
	COIECD <sub>Prompt</sub>	.4407	.6316	.5855	.7087	.3860	.4493	.4565	<u>.6181</u>	.5294	.6143	.4882	.6047	.5061	.6138
	RetRobust	<u>.5536</u>	.5536	.6606	.6606	.4089	.4502	.4507	.5430	.5441	.6063	.6063	<u>.5464</u>	<u>.5464</u>	.5465
	KAFT	<b>.6144</b>	.6144	<b>.7657</b>	<u>.7657</u>	<b>.4764</b>	.4764	<b>.5124</b>	.5129	<b>.6485</b>	<u>.6489</u>	<b>.7145</b>	<u>.7145</u>	<b>.5916</b>	.5917
	LM <sub>UniKnow</sub>	.5396	<b>.7396</b>	<u>.6915</u>	<b>.8421</b>	.4282	<b>.6593</b>	.4513	<b>.6897</b>	.5513	<b>.7557</b>	<u>.6546</u>	<b>.7771</b>	.4901	<b>.7319</b>
MISTRAL-7B	Naïve	.4444	.4444	.6270	.6270	<u>.4586</u>	.4586	<b>.5109</b>	.5111	<u>.5911</u>	.5911	.5658	.5658	<u>.5386</u>	.5386
	Prompting	.3304	.5615	.6149	.7028	.3459	.5471	.3806	<u>.6145</u>	.4634	<u>.6677</u>	.4761	.6244	.4695	<b>.6786</b>
	COIECD	.4601	.4603	.5917	.5919	.4575	.4575	<u>.5011</u>	.5015	.5653	.5653	.5457	.5457	.5351	.5352
	COIECD <sub>Prompt</sub>	.4039	.6053	.6179	.6750	.3525	<u>.5535</u>	.4327	<b>.6389</b>	.4516	.6269	.4967	.6267	.4705	<u>.6681</u>
	RetRobust	<u>.5873</u>	.5875	<u>.6519</u>	.6520	.4231	.4232	.4386	.4393	.5728	.5735	<u>.6139</u>	.6139	.5376	.5380
	KAFT	<b>.6055</b>	<u>.6057</u>	<b>.7251</b>	<u>.7252</u>	<b>.4631</b>	.4633	.4942	.4948	<b>.5983</b>	.5994	<b>.7131</b>	<u>.7131</u>	<b>.5637</b>	.5638
	LM <sub>UniKnow</sub>	.5129	<b>.7466</b>	.6138	<b>.8270</b>	.3830	<b>.6329</b>	.3434	.5845	.4602	<b>.7071</b>	.5533	<b>.7792</b>	.4040	.6507
QWEN-1.5B	Naïve	.4300	.4306	.5005	.5014	.4056	.4057	<u>.4615</u>	.4622	.4727	.4738	.5192	.5194	.5023	.5037
	Prompting	.4067	<u>.5683</u>	.4780	<u>.6333</u>	.3723	.5370	.4462	<u>.5830</u>	.4225	<u>.6465</u>	.4427	.6174	.4547	<u>.6714</u>
	COIECD	.4152	.4161	.5009	.5035	.3560	.3566	.4539	.4560	.4476	.4494	.4870	.4877	.4938	.4978
	COIECD <sub>Prompt</sub>	.3769	.4919	.4845	.5681	.3383	.4280	.4297	.5218	.4362	.5668	.4657	.5653	.4743	.6341
	RetRobust	<u>.4706</u>	.4708	<u>.5502</u>	.5503	.4087	.4087	.4603	.4612	<u>.5312</u>	.5319	<u>.6113</u>	.6115	<u>.5299</u>	.5299
	KAFT	<b>.4904</b>	.4905	<b>.5795</b>	.5796	<b>.4315</b>	.4317	<b>.4904</b>	.4913	<b>.5861</b>	.5861	<b>.6536</b>	<b>.6539</b>	<b>.5426</b>	.5427
	LM <sub>UniKnow</sub>	.4506	<b>.6414</b>	.5342	<b>.7351</b>	.3787	<b>.6171</b>	.4420	<b>.6547</b>	.4878	<b>.7078</b>	.5713	<u>.6323</u>	.4783	<b>.7154</b>
QWEN-3B	Naïve	.4388	.4393	.5137	.5145	.4192	.4194	.4680	.4688	.4993	.4996	.5116	.5116	.5061	.5086
	Prompting	.3878	<u>.6106</u>	.4299	.6608	.3497	<u>.5903</u>	.4279	<u>.6460</u>	.4275	.6577	.4025	.6241	.4512	<u>.6825</u>
	COIECD	.4263	.4278	.5130	.5174	.4066	.4075	.4617	.4636	.4803	.4831	.4858	.4889	.5048	.5125
	COIECD <sub>Prompt</sub>	.3782	.5787	.4681	<u>.6744</u>	.3566	.5739	.4309	.6247	.4336	.6278	.4325	.5966	.4639	.6748
	RetRobust	.5134	.5135	.5804	.5805	.4251	.4251	.4688	.4698	.5864	.5868	.6477	.6477	<u>.5410</u>	.5410
	KAFT	<b>.5397</b>	.5399	<b>.6356</b>	.6357	<b>.4620</b>	.4622	<b>.5106</b>	.5114	<b>.6768</b>	<u>.6786</u>	<b>.7029</b>	<u>.7029</u>	<b>.5589</b>	.5589
	LM <sub>UniKnow</sub>	<u>.5136</u>	<b>.6925</b>	<u>.6172</u>	<b>.7710</b>	<u>.4407</u>	<b>.6601</b>	<u>.4817</u>	<b>.6847</b>	<u>.5897</u>	<b>.7675</b>	<u>.6612</u>	<b>.7423</b>	.5072	<b>.7349</b>
QWEN-7B	Naïve	.4523	.4529	.5870	.5900	<u>.4413</u>	.4450	.4905	.4929	.5735	.5742	.5573	.5578	.5132	.5147
	Prompting	.3788	.6230	.4672	.6973	.3759	.6191	.4493	<u>.6798</u>	.4487	<u>.6860</u>	.4422	.6659	.4639	<u>.7010</u>
	COIECD	.4410	.4413	.5785	.5847	.4183	.4190	.4772	.4807	.5215	.5222	.5329	.5331	.5057	.5106
	COIECD <sub>Prompt</sub>	.4096	<u>.6386</u>	.4975	<u>.7052</u>	.3612	.5871	.4469	.6687	.4864	.6856	.4517	.6620	.4758	.6984
	RetRobust	<u>.5613</u>	.5616	<u>.6155</u>	.6156	.4325	.4326	<u>.5139</u>	.5147	<u>.6270</u>	.6270	<u>.6610</u>	.6612	<u>.5959</u>	.5959
	KAFT	<b>.5928</b>	.5929	<b>.6935</b>	.6935	<b>.4772</b>	.4773	<b>.5334</b>	.5341	<b>.6844</b>	.6844	<b>.7502</b>	<u>.7502</u>	<b>.6033</b>	.6033
	LM <sub>UniKnow</sub>	.5074	<b>.7280</b>	.5881	<b>.8032</b>	.4083	<b>.6561</b>	.4750	<b>.7161</b>	.5624	<b>.7768</b>	.6409	<b>.7887</b>	.5122	<b>.7571</b>
QWEN-14B	Naïve	.4743	.4746	.6522	.6523	.4614	.4616	.5051	.5063	<u>.6385</u>	.6385	.5661	.5663	.5268	.5285
	Prompting	.4141	<u>.6470</u>	.5146	.7240	.4036	<u>.6293</u>	.4510	<u>.6888</u>	.4867	<u>.7220</u>	.4709	.6681	.4630	<u>.7035</u>
	COIECD	.4421	.4427	.6220	.6228	.4294	.4302	.4740	.4766	.5768	.5789	.5379	.5489	.5016	.5037
	COIECD <sub>Prompt</sub>	.4045	.6170	.5550	.7167	.3852	.5826	.4467	.6645	.4900	.6992	.4922	.6559	.4704	.6931
	RetRobust	<u>.6158</u>	.6162	<u>.6753</u>	.6754	.4741	.4741	<u>.5194</u>	.5205	.5699	.5703	<u>.6795</u>	.6797	<u>.5845</u>	.5845
	KAFT	<b>.6415</b>	.6416	<b>.7682</b>	<u>.7682</u>	<b>.5252</b>	.5252	<b>.5637</b>	.5646	<b>.6962</b>	.6962	<b>.7640</b>	<u>.7640</u>	<b>.6066</b>	.6069
	LM <sub>UniKnow</sub>	.5453	<b>.7555</b>	.6543	<b>.8448</b>	.4529	<b>.6912</b>	.4926	<b>.7299</b>	.5452	<b>.7729</b>	.6456	<b>.8335</b>	.5342	<b>.7676</b>

Table 12: Exact values of Acc and Rely for each method and model across datasets. **Bold** indicates the best, and the underline indicates the second best.