

# Large Language Models as Reliable Knowledge Bases?

Danna Zheng<sup>1</sup>, Mirella Lapata<sup>1</sup>, Jeff Z. Pan<sup>1,2</sup>

<sup>1</sup> School of Informatics, University of Edinburgh, UK

<sup>2</sup> Huawei Edinburgh Research Centre, CSI, UK

dzheng@ed.ac.uk, mlap@inf.ed.ac.uk, <http://knowledge-representation.org/j.z.pan/>

## Abstract

The NLP community has recently shown a growing interest in leveraging Large Language Models (LLMs) for knowledge-intensive tasks, viewing LLMs as potential knowledge bases (KBs). However, the reliability and extent to which LLMs can function as KBs remain under-explored. While previous studies suggest LLMs can encode knowledge within their parameters, the amount of parametric knowledge alone is not sufficient to evaluate their effectiveness as KBs. This study defines criteria that a reliable LLM-as-KB should meet, focusing on factuality and consistency, and covering both seen and unseen knowledge.<sup>1</sup> We develop several metrics based on these criteria and use them to evaluate 26 popular LLMs, while providing a comprehensive analysis of the effects of model size, instruction tuning, and in-context learning (ICL). Our results paint a worrying picture. Even a high-performant model like GPT-3.5-turbo is not factual or consistent, and strategies like ICL and fine-tuning are unsuccessful at making LLMs better KBs.

## 1 Introduction

Large language models (LLMs), pretrained on extensive text corpora, have demonstrated the ability to implicitly encode various types of knowledge within their weights, without requiring human supervision. As a result, many recent studies (Chuang et al., 2023; Yu et al., 2023; Dhingra et al., 2022; Sung et al., 2021; Wang et al., 2020) aim to analyze the relationship between LLMs and KBs, and even explore whether LLMs can replace KBs (Sun et al., 2023; Mruthyunjaya et al., 2023; Heinzerling and Inui, 2021).

<sup>1</sup>*Seen knowledge* refers to knowledge learned during training. *Unseen knowledge* is neither present in the model’s training data nor can be inferred from seen knowledge.

However, whether current LLMs can serve as reliable KBs and how to evaluate their performance in this role remains largely unexplored. Existing studies (Sun et al., 2023; Wang et al., 2021; Roberts et al., 2020) often implicitly assume that the LLM’s ability to retain knowledge is sufficient for it to function as a KB. Typically, these studies employ two methods: (1) converting knowledge graphs into natural language questions using templates and evaluating LLM ability to answer these questions, by measuring the amount of knowledge therein (Petroni et al., 2019; Sun et al., 2023); and (2) pre-training LLMs on passages/triples containing knowledge and then assessing their ability to answer related questions, thereby quantifying their knowledge retention (Wang et al., 2021; He et al., 2024). While these methods demonstrate that LLMs can recall knowledge, memorizing facts is not the sole criterion for being a reliable KB (AlKhamissi et al., 2022).

What criteria then should a LLM meet to function as a reliable KB? Discussion on this topic has been limited, and there is no agreement on the definition of these criteria. AlKhamissi et al. (2022) argue that LLMs ought to excel at five aspects (i.e., access, edit, consistency, reasoning, explainability, and interpretability) if they are to be considered a KB. However, they do not outline specific metrics to evaluate the extent to which LLMs act as KBs. Besides, we argue that evaluating LLMs against the characteristics of KBs may not be entirely appropriate due to their different data storage structures. Instead, we should consider the specific properties of LLMs when assessing their suitability as KBs.

Our research addresses these gaps and aims to establish a more nuanced understanding of LLMs-as-KBs. We evaluate and compare the reliability of different LLMs functioning as KBs in answering factoid questions. Specifically, our work makes the following contributions:

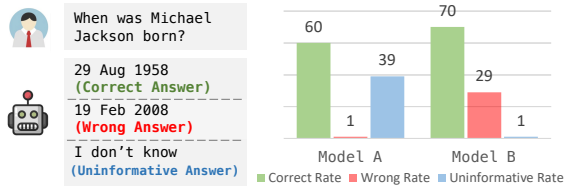


Figure 1: An example illustrating three answer types: correct, wrong, and uninformative. Focusing only on the correct rate incorrectly suggests that Model B is better, even though Model A is more reliable with a similar correct rate and a much lower wrong rate.

1. We define how to assess the reliability of LLMs-as-KBs, and propose metrics that consider the differences between KBs and LLMs and align with the distinct nature of LLMs. We consider two aspects, *factuality* (i.e., the ability to provide factual responses) and *consistency* (i.e., the ability to provide consistent responses for questions involving the same knowledge), and cover the evaluation on both seen and unseen knowledge.
2. To evaluate consistency, we propose a novel method which computes the probability that a LLM can consistently provide response  $r$  for question  $q$ .
3. To evaluate LLMs on unseen knowledge, we create a new QA dataset  $UnseenQA$ , with a knowledge cutoff date before April 13, 2024.
4. We evaluate 26 popular LLMs in their ability to function as reliable KBs, and discuss the influence of model size, instruction-tuning, and ICL on their performance.

## 2 What is a Reliable LLM-as-KB?

In simple terms, a LLM functions as a reliable KB if it consistently provides factual responses. Evaluating the reliability of LLMs as KBs primarily involves assessing two critical dimensions, namely factuality and consistency.

### 2.1 Factuality

**Factuality** refers to the quality of being factual or based on fact. Run-of-the-mill KBs, stored on physical servers or cloud platforms, deliver information directly in response to queries. If the requested data is unavailable, these systems typically return a null response. In contrast, LLMs

are probabilistic models that excel at next word prediction based on the given context, rather than storing explicit information in defined locations. This architecture allows LLMs to generate responses that seem plausible, regardless of whether the content was included in their training data. Consequently, LLMs typically produce three types of responses: correct, uninformative, and wrong.

Ongoing efforts (Chen et al., 2023; Wang et al., 2024) to evaluate the factuality of LLMs often hinge on the models’ correct rates in responding to factual QA datasets. However, this approach has notable limitations. Firstly, many studies (Lin et al., 2022; Sun et al., 2023) do not specify whether the dataset’s scope of knowledge was included in the LLM’s training data. This omission can lead to unfair comparisons between models, especially if the knowledge being tested is within the training scope of one model but not another. Secondly, the assumption that a higher correct rate indicates greater factuality is problematic. For example, consider models A and B shown in Figure 1, and assume they are being evaluated on a test dataset covering knowledge they both have seen during training. In this case, focusing solely on the correct rate might erroneously suggest that Model B is more factual, despite Model A being more reliable due to its similar correct rate and significantly lower wrong rate.

Given these issues, we propose the following criteria for the factuality of LLMs-as-KBs:

**CRITERION 1.1** For seen knowledge, a factual LLM should demonstrate a high correct rate and a low wrong rate.

**CRITERION 1.2** For unseen knowledge, a factual LLM should demonstrate a high uninformative rate.

We next proceed to define evaluation metrics that operationalize these criteria. Let  $M$  denote a LLM. Let  $D_{\text{seen}}$  denote a QA dataset containing  $N$  open-ended factoid questions pertaining to knowledge the LLM ought to have seen during training. Let  $D_{\text{unseen}}$  denote a QA dataset with  $L$  open-ended factoid questions covering unseen knowledge. We further assume the LLM’s response to  $D_{\text{seen}}$  will be correct, uninformative, or wrong, while its response to  $D_{\text{unseen}}$  will be either uninformative or wrong.

**METRIC 1.1: Net Correct Rate (NCR)** measures how much more likely the model is to pro-

vide correct responses instead of wrong ones on  $D_{\text{seen}}$  questions. It is defined as:

$$\text{NCR} = \text{CR} - \text{WR} \quad (1)$$

$$\text{CR} = \frac{N_{\text{correct}}}{N} \quad \text{WR} = \frac{N_{\text{wrong}}}{N} \quad (2)$$

where  $N_{\text{correct}}$  and  $N_{\text{wrong}}$  are counts of correct and wrong responses, respectively.

NCR values range from  $-1$  to  $1$ . A negative NCR suggests the model tends to provide misleading responses, while a positive NCR suggests a preference for correct responses. Consider again two models, A and B. According to *CRITERION 1.1*, if model A has a higher correct rate and lower wrong rate compared to model B, then model A is better. Formally, if  $\text{CR}_A - \text{CR}_B > \text{WR}_A - \text{WR}_B$ , then model A is better than B. Algebraically, this is equivalent to  $\text{CR}_A - \text{WR}_A > \text{CR}_B - \text{WR}_B$ , i.e.,  $\text{NCR}_A > \text{NCR}_B$ . Therefore, a higher NCR indicates a more factual model on seen knowledge.

**METRIC 1.2: Uninformative Rate (UR)** assesses whether the model is likely to provide uninformative responses to  $D_{\text{unseen}}$  questions. It is formulated as:

$$\text{UR} = \frac{L_{\text{uninformative}}}{L} \quad (3)$$

where  $L_{\text{uninformative}}$  denotes the count of uninformative responses. UR ranges from 0 to 1. A higher UR indicates that the model is more likely to refrain from giving wrong responses when faced with unseen knowledge.

## 2.2 Consistency

**Consistency** refers to the quality of always behaving in the same way or having the same opinions. KBs are designed with this principle in mind. In fact, there are efficient algorithms (Andersen and Pretolani, 2001) which detect and resolve conflicts within KBs, thus ensuring consistent outputs for queries on the same facts.

It is well-known that LLMs often exhibit inconsistencies in their responses (Elazar et al., 2021; Wang et al., 2022). Current research (Elazar et al., 2021; Jang et al., 2022; Hagström et al., 2023) evaluates LLM consistency through their performance on benchmarks involving paraphrasing, negation, or multilingual variations. A model is considered superior if it responds consistently across a broader range of data samples. In this

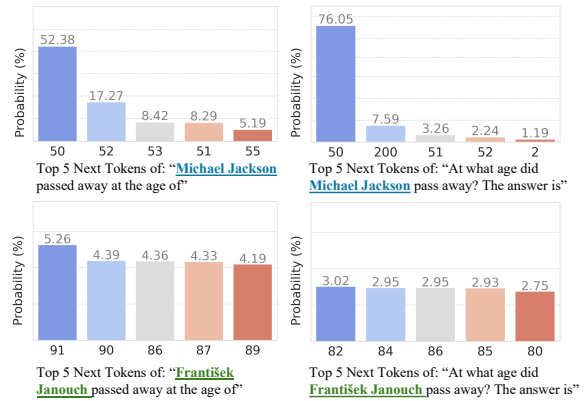


Figure 2: Illustration of LLM inconsistency with DAVINCI-002 (temperature is set to 0). Questions in the top focus on seen knowledge, with probability distribution mass concentrated on one prediction. Questions in the bottom focus on unseen knowledge, where the distribution is more even. Drawing from such a distribution inevitably leads to inconsistencies.

work, we argue that it may be too strict to expect LLMs to be always consistent when responding to fact-based questions. Unlike explicit KBs which store information at a fixed location, LLMs operate probabilistically. In theory, if the context has been learned during training, the probability distribution for the prediction will be concentrated; otherwise, it will be more uniform. Drawing from a uniform probability distribution inevitably leads to inconsistencies. As shown in Figure 2, even with greedy decoding, slight biases in the distribution can cause fluctuations in the selection of the top probable words.

Given their probabilistic nature, we do not expect LLMs to always behave consistently. We acknowledge that inconsistencies can cause confusion in practical applications and propose to monitor model behavior through post-processing which we argue is more realistic than expecting a probabilistic model to be perpetually consistent. We thus propose the following consistency criteria:

**CRITERION 2.1** The model is expected to be consistent when it produces correct responses.

**CRITERION 2.2** The model is expected to be inconsistent when it produces wrong responses.

We next define evaluation metrics corresponding to the criteria above. Let  $q$  refer to a question in either  $D_{\text{seen}}$  or  $D_{\text{unseen}}$ , and  $r$  denote model  $M$ 's response to  $q$ . Inspired by Zheng et al. (2024), we measure consistency based on multiple-choice

questions (MCQs). As shown in Figure 3, we employ GPT-3.5-TURBO-INSRUCT to generate a set of distractor options similar to response  $r$ , and then create a group of MCQs. The consistency score for data point  $(q, r)$  is calculated as:

$$Cons(q, r) = \frac{\sum_{i=1}^{X_{MCQs}} [R_i = r]}{X_{MCQs}} \quad (4)$$

where  $X_{MCQs}$  is the total number of MCQs,  $R_i$  is model  $M$ 's response for the  $i$ -th MCQ, and the expression  $[R_i = r]$  yields 1 when the model's response  $R_i$  matches its original response  $r$ , and 0 otherwise. The consistency score  $Cons(q, r)$  ranges from 0 to 1.

**METRIC 2.1:**  $C_{correct}$  measures the consistency of the model when it provides correct responses. It is defined as:

$$C_{correct} = \frac{\sum_{j=1}^{N_{correct}} Cons(q_j^{(c)}, r_j^{(c)})}{N_{correct}} \quad (5)$$

where  $r^{(c)}$  refers to the response labeled as correct, and  $q^{(c)}$  is the corresponding question.  $C_{correct}$  ranges from 0 to 1. Based on *CRITERION 2.1*, a higher  $C_{correct}$  is desirable.

**METRIC 2.2:**  $C_{wrong}$  measures the consistency of an LLM when it provides wrong responses. It is defined as:

$$C_{wrong} = \frac{C_{wrong}^s + C_{wrong}^u}{2} \quad (6)$$

where  $C_{wrong}^s$  refers to the consistency of a LLM when it provides wrong responses to questions about seen knowledge, and  $C_{wrong}^u$  refers to the consistency of a LLM when it provides wrong responses to questions about unseen knowledge:

$$C_{wrong}^s = \frac{\sum_{j=1}^{N_{wrong}} Cons(q_j^{(w)}, r_j^{(w)})}{N_{wrong}} \quad (7)$$

$$C_{wrong}^u = \frac{\sum_{j=1}^{L_{wrong}} Cons(q_j^{(w)}, r_j^{(w)})}{L_{wrong}} \quad (8)$$

where  $r^{(w)}$  refers to responses labeled as wrong, and  $q^{(w)}$  is the corresponding question.  $N_{wrong}$  and  $L_{wrong}$  are counts of wrong responses to  $D_{seen}$  and  $D_{unseen}$ , respectively.  $C_{wrong}$  ranges from 0 to 1. Based on *CRITERION 2.2*, a lower  $C_{wrong}$  is desirable.

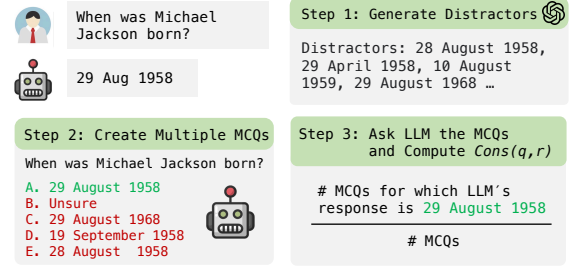


Figure 3: Example computation for consistency score  $Cons(q, r)$ . The LLM's original answer is shown in green, while distractors are red.

### 2.3 Reliability (Factuality and Consistency)

A reliable LLM-as-KB should then be assessed against factuality *and* consistency. Based on the criteria defined above, a LLM-as-KB is reliable if it meets the following:

**CRITERION 3.1** For seen knowledge, a LLM should have a high rate of consistently correct responses and a low rate of consistently wrong responses.

**CRITERION 3.2** For unseen knowledge, a LLM should have a high rate of uninformative or inconsistent responses.

We quantify these criteria with two metrics.

**METRIC 3.1: Net Consistently Correct Rate (NCCR)** quantifies the model's tendency to provide consistently correct responses compared to consistently wrong ones for questions about seen knowledge. It is defined as:

$$NCCR = CCR - CWR \quad (9)$$

$$CCR = CR \times C_{correct}$$

$$CWR = WR \times C_{wrong}^s$$

NCCR ranges from  $-1$  to  $1$ . A higher NCCR indicates a LLM is more reliable on seen knowledge. A negative NCCR suggests the model provides consistently wrong responses, while a positive NCCR suggests a preference for consistently correct responses.

**METRIC 3.2: Inconsistent/Uninformative Rate (IUR)** assesses whether the model is likely to provide uninformative or inconsistent wrong responses for questions about unseen knowledge. It is defined as:

$$IUR = 1 - (1 - UR)C_{wrong}^u \quad (10)$$

IUR ranges from 0 to 1. A higher IUR value indicates the LLM functions as a more reliable KB on unseen knowledge.

### 3 Experimental Setup

#### 3.1 Datasets

To evaluate the performance of LLMs on seen knowledge, we collated  $SeenQA$ , a composite dataset containing 3,000 questions sourced from the test sets (or development sets, where test sets were unavailable) of Natural Questions (Kwiatkowski et al., 2019), TriviaQA (Joshi et al., 2017), and PopQA (Mallen et al., 2023). For assessing LLM performance on unseen knowledge, we introduce  $UnseenQA$ , a new dataset designed to ensure that models with a knowledge cutoff date before April 13, 2024 do not have access to answers.  $UnseenQA$  also includes 3,000 questions, which are derived from 20 templates covering various answer types such as number, people, time, location, and others. Our methodology for creating  $SeenQA$  and  $UnseenQA$  is detailed in Appendix A.

#### 3.2 LLM Selection

We evaluate a wide range of popular LLMs (26 in total) in their ability to function as KBs and investigate the impact of various factors, such as the number of parameters and fine-tuning. Specifically, we consider the following models: `GPT-3.5-TURBO`, `FLAN-T5` (0.08B, 0.25B, 0.78B, 3B, 11B), `LLAMA1` (7B, 13B, 65B), `LLAMA2` (7B, 13B, 70B), `LLAMA2-CHAT` (7B, 13B, 70B), `LLAMA3` (8B, 70B), `LLAMA3INSTRUCT` (8B, 70B), `MISTRAL` (7B), `MISTRAL-INSTRUCT` (7B), `GEMMA` (2B, 7B), `GEMMA-INSTRUCT` (2B, 7B), and `PHI2` (3B). Detailed descriptions of these models are provided in Table 4 in Appendix B. We classify LLMs into three categories based on their parameter sizes: small (0.08B–3B), medium (7B–13B), and large (65B–70B). We use the term ‘fine-tuned LLMs’ to refer to the LLMs that have been fine-tuned through instruction-tuning or reinforcement learning from human feedback. We use the term ‘base LLMs’ to refer to LLMs without fine-tuning.

#### 3.3 Evaluation on a Single Response

**Uninformative** We identify three types of uninformative responses from LLMs: ‘repetition’, ‘none’, and ‘unsure’. ‘Repetition’ refers to re-

sponses that repeatedly echo a specific string. We detect this using regular expressions and word frequency analysis. ‘None’ denotes responses lacking relevant information, such as an empty string and repetition of the question. ‘Unsure’ indicates responses where the model explicitly states it is unable to answer or lacks the required knowledge. We label a response as ‘unsure’ if it includes expressions such as ‘I am not sure’, ‘I cannot provide’, and ‘I am just an AI’.

**Correct** We determine whether a response is correct based on exact match. A response is considered correct if the exact match score is 1. In cases where exact match is 0, we compare the model’s prediction against the ground truth using `gpt-4o`. Previous work (Sun et al., 2023) has shown a 98% agreement rate between ChatGPT and human judgments in comparing model responses with ground truth. We follow their prompt as outlined in Table 7 in Appendix C.

**Consistency Score** To compute  $Cons(q, r)$ , we set  $X_{MCQs}$  (total number of MCQs) to 20, and each MCQ includes question  $q$  and 5 options (the original response  $r$ , 3 random distractor options, and an ‘unsure’ option).

#### 3.4 Prompts and Hyper-parameters

To provide a comprehensive evaluation, we experimented with three types of prompt settings: zero-shot, four-shot, and four-shot with two unsure shots. To avoid any bias introduced by fixed examples, we employed a dynamic few-shot method following the work of Nori et al. (2023). We collected two repositories,  $R_{seen}$  and  $R_{unseen}$ .  $R_{seen}$  includes 280 question-answer pairs about seen knowledge (200 from the unused data of PopQA and training data of Natural Questions and TriviaQA; 80 are generated using the templates in Table 3).  $R_{unseen}$  consists of 40 question-answer pairs about unseen knowledge, all generated using the templates in Table 3. We used `TEXT-EMBEDDING-3-SMALL` to embed the questions in the repositories and test questions as vector representations. For each test question under the four-shot setting, we retrieved its nearest four questions from  $R_{seen}$ . Under the four-shot with two unsure shots setting, we retrieved the nearest two questions from  $R_{seen}$  and two from  $R_{unseen}$ .

All LLMs were evaluated using greedy decoding (temperature 0 for `GPT-3.5-TURBO`) with a

| LLM            | Zero-Shot |      |             | Four-Shot |      |             | Four-Two |       |             |
|----------------|-----------|------|-------------|-----------|------|-------------|----------|-------|-------------|
|                | NCCR      | IUR  | AVG         | NCCR      | IUR  | AVG         | NCCR     | IUR   | AVG         |
| GPT-3.5-TURBO  | 34.1      | 95.7 | <b>62.7</b> | 35.5      | 99.0 | <b>66.7</b> | 32.1     | 99.8  | <b>65.8</b> |
| LLAMA2CHAT-70B | 20.4      | 98.9 | <b>59.1</b> | 17.7      | 99.9 | <b>58.8</b> | 16.4     | 100.0 | 58.2        |
| LLAMA3-70B     | 30.2      | 71.6 | 36.7        | 33.1      | 69.2 | 35.7        | 33.9     | 98.8  | <b>65.8</b> |

Table 1: Most reliable LLMs across three prompt settings. Four-Two refers to the four-shot with two unsure shots setting. All numbers shown are percentages. AVG represents average NCCR and (normalized) IUR ( $2 * IUR - 100$ ) scores (both scores use the same scale). We treat seen and unseen knowledge equally, however, in practice, NCCR and IUR can be assigned different weights.

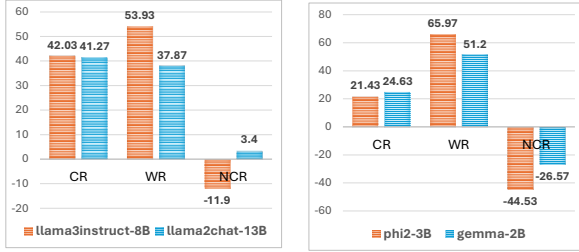


Figure 4: Comparing LLaMA3INSTRUCT-8B against LLaMA2CHAT-13B, and PHI2-3B against GEMMA-2B using CR, WR, and NCR metrics (zero-shot setting). Even though these models are comparable according to CR, they are quite different according to WR and NCR metrics. We observe similar trends in the four-shot and with two unsure shots settings.

maximum of 100 new tokens. Our prompts are provided in Appendix C.

## 4 Results

We present detailed results for all LLMs in Table 9 (factuality), Table 10 (consistency), and Table 11 (reliability) in Appendix D. LLM rankings based on different metrics are shown in Figure 12 and Figure 13, also in Appendix D.

**GPT-3.5-TURBO is overall the most reliable LLM.** Table 1 presents results for the two best performing LLMs under different prompt settings. As can be seen, GPT-3.5-TURBO is most reliable across the board. Although it is not consistently wrong when asked about facts it does not know (its IUR score exceeds 95%), it is not consistently correct when asked about facts it *has seen* before (NCCR is only 32%).

**FLAN-T5-0.78B is the most reliable LLM for unseen knowledge and most unreliable with seen knowledge.** Figure 5 shows the LLMs ranked by NCCR and IUR in a zero-shot setting. As can be seen, while FLAN-T5-0.78B ranks low

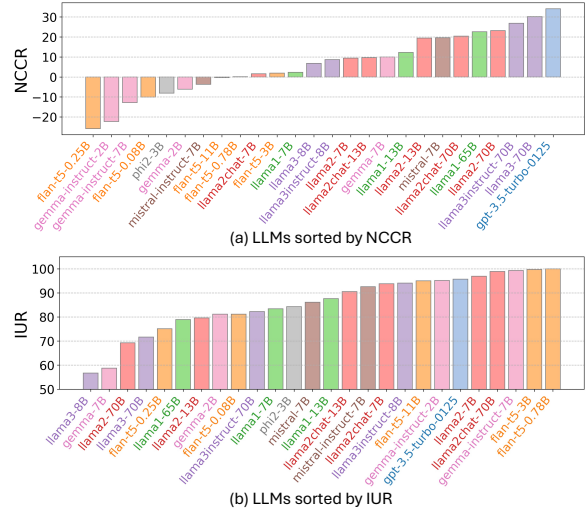


Figure 5: LLM ranking based on NCCR and IUR metrics in zero-shot setting. See full results in Figure 12 (Appendix D) for other shot settings.

for NCCR, it maintains the top position for IUR. A similar trend is observed with the FLAN-T5-3B and GEMMA-INSTRUCT (2B, 7B) models. Conversely, models in the LLaMA3 family show the opposite trend: they rank high on seen knowledge but low on unseen knowledge. For instance, LLaMA3-70B, despite ranking second in terms of NCCR, falls to the fourth lowest position for IUR.

**Net Correct Rate (NCR) reveals factuality gaps in LLMs.** As illustrated in Figure 4, CR, the standard metric for assessing factuality, fails to fully capture nuanced differences. Despite similar CR values of approximately 40%, LLaMA3INSTRUCT-8B and LLaMA2CHAT-13B behave differently when it comes to wrong responses (on seen data); the former model has a WR of 15% higher than the latter, and as a result its NCR is substantially lower. In the case of PHI2-3B, Gunasekar et al. (2023) claim that with "textbook quality" data, smaller LLMs can

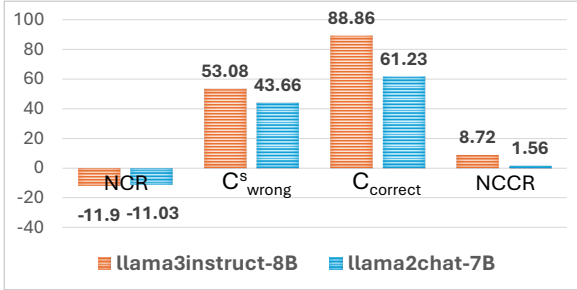


Figure 6: Comparison between LLAMA3INSTRUCT-8B and LLAMA2CHAT-7B (zero-shot setting). The former model is slightly less factual according to NCR, but more consistent in its responses (see NCCR,  $C_{correct}$ , and  $C_{wrong}^s$ ). We observe similar trends in the four-shot and with two unsure shots settings.

| Comparisons                   | Zero-Shot | Four-Shot | Four-Two |
|-------------------------------|-----------|-----------|----------|
| NCR vs. UR                    | 0.27      | 0.34      | 0.62*    |
| NCCR vs. IUR                  | -0.17     | -0.12     | 0.41*    |
| $C_{correct}$ vs. $C_{wrong}$ | 0.81*     | 0.78*     | 0.51*    |

Table 2: Pairwise correlation of LLM performance on different metrics under different prompt settings. The correlation values are computed across all LLMs (26 data points). We report Pearson’s  $\rho$ , diacritic \* denotes statistical significance ( $p < 0.05$ ). Four-Two refers to the four-shot setting with two unsure shots.

achieve satisfactory performance with less training data. Focusing solely on CR, this claim is true as PHI2-3B’s CR is comparable to that of GEMMA-2B, which was trained with twice the token count. However, PHI2-3B exhibits significantly higher WR (and lower NCR) compared to GEMMA-2B which underscores the challenge of maintaining low error rates even with superior data quality.

**A less factual LLM can be more reliable.** Figure 6 illustrates this finding. The LLAMA3INSTRUCT-8B model has a slightly worse NCR compared to LLAMA2CHAT-7B, indicating that the former model is marginally less factual on seen knowledge. However, LLAMA3INSTRUCT-8B exhibits a higher NCCR due to its better consistency score on correct responses. This suggests that LLAMA3INSTRUCT-8B is more reliable than LLAMA2CHAT-7B on seen knowledge. Additionally, both LLMs show negative NCR and positive NCCR in the zero-shot setting. This indicates that although they produce more wrong than correct responses, they generate fewer consistently wrong responses compared to consistently correct ones.

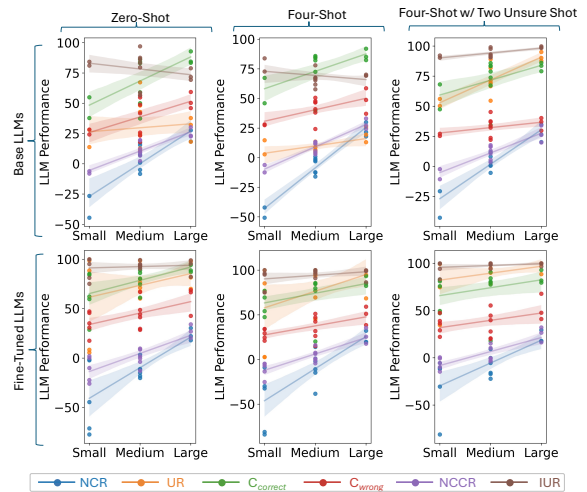


Figure 7: The impact of model size on LLM performance, measured with NCR, UR,  $C_{correct}$ ,  $C_{wrong}$ , NCCR, and IUR. Different metrics are color-coded. LLMs are shown in three sizes, small, medium, and large and are grouped into ‘base’ and fine-tuned ones.

**Performance on seen knowledge is not predictive of performance on unseen knowledge.** In Table 2, we examine whether metrics applied to seen knowledge can be used to extrapolate model performance on unseen knowledge by reporting correlation values (Pearson’s  $\rho$ ) between NCR and UR, and NCCR and IUR. As can be seen, in both zero- and four-shot settings, correlations are not statistically significant. Nevertheless, correlations are significant in all metric comparisons in the four-shot with two unsure shots setting. While in general model performance with seen knowledge does not transfer to unseen knowledge, specific prompt manipulations can enhance the correlation between metrics (see last column in Table 2).

**LLMs are consistently right and wrong!** Table 2 reports a positive, significant correlation between  $C_{correct}$  and  $C_{wrong}$ . This result implies that LLMs demonstrating high consistency in correct responses also tend to exhibit high consistency in wrong responses. This finding contradicts our expectation of high  $C_{correct}$  and low  $C_{wrong}$  and highlights a notable flaw in current models. Future work should address this issue, e.g., by instructing LLMs to achieve the desired consistency behavior.

## 5 Analysis

In this section, we explore the effects of model size, fine-tuning, and ICL on LLM performance

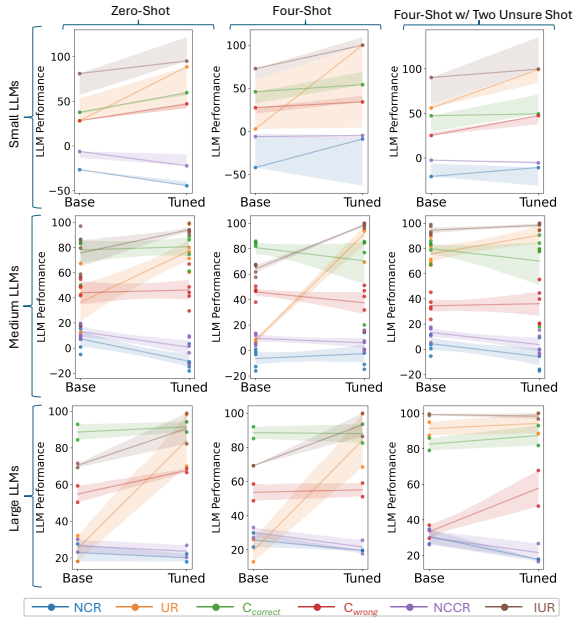


Figure 8: The impact of fine-tuning on LLM performance, measured with NCR, UR,  $C_{correct}$ ,  $C_{wrong}$ , NCCR, and IUR. Different metrics color-coded. This analysis only considers the performance of Llama2, Llama3, Mistral, and Gemma as these families include both base LLMs and fine-tuned versions. Models are shown in three sizes, small, medium, and large.

based on different metrics. We also analyze how LLMs handle unseen knowledge by looking at the distribution of uninformative responses and the impact of different question types.

### 5.1 The Impact of Model Size

**Larger LLMs perform better on seen knowledge but worse on unseen knowledge.** Figure 7 demonstrates that as model size increases, both NCR (blue line) and NCCR (purple line) improve. This trend indicates that larger LLMs perform better on questions about seen knowledge. However, for base LLMs in both zero-shot and four-shot settings, IUR (brown line) decreases as model size increases. This trend suggests that as LLMs become larger, they become more consistent at delivering wrong responses for questions regarding unseen knowledge.

**Larger LLMs are more consistent, even with wrong responses.** In Figure 7, we observe that as model size increases, both  $C_{correct}$  (green line) and  $C_{wrong}$  (red line) increase significantly. While higher consistency scores for correct responses are expected and useful, the higher consistency scores for wrong responses pose a potential risk. Larger

models may consistently produce convincing but wrong information, which could lead to misinformation if not carefully managed.

### 5.2 The Impact of Fine-tuning

**Fine-tuning improves performance on unseen knowledge but negatively affects performance on seen knowledge.** As illustrated in Figure 8, UR (orange line) and IUR (brown line) show significant improvement in LLMs after fine-tuning, indicating an enhanced ability to distinguish and respond to unseen knowledge appropriately. However, the decreasing NCCR (purple line) suggests they become worse at handling seen knowledge after fine-tuning.

**Fine-tuning does not make LLMs more consistent.** As depicted in Figure 8, there is no noticeable increase in  $C_{correct}$  (green line) after fine-tuning. In fact, fine-tuning even has a negative impact on  $C_{correct}$  for medium-sized LLMs. These results indicate that current fine-tuning techniques fail to enhance the consistency of correct responses. Furthermore,  $C_{wrong}$  (red line) does not decrease after fine-tuning either, which suggests that fine-tuning also fails to reduce the model’s persistence on wrong responses.

### 5.3 The Impact of ICL

**Unsure shots improve LLM performance on unseen knowledge.** As shown in Figure 9, incorporating two unsure shots in the four-shot setting substantially increases UR (orange line) and IUR (brown line) across model sizes. In contrast, performance on unseen knowledge deteriorates when using four-shots only for all sizes of base LLMs and small fine-tuned LLMs, compared to their zero-shot counterparts.

**ICL does not improve LLM performance on seen knowledge.** Figure 9 demonstrates that ICL does not improve LLM performance according to NCR (blue line) or NCCR (purple line) metrics. For small/medium base LLMs and small fine-tuned LLMs in the four-shot setting, ICL even decreases NCR performance. This finding contrasts with previous results about the ability of LLMs to learn from a few examples (Brown et al., 2020; Chada and Natarajan, 2021; Touvron et al., 2023; Bai et al., 2023). Earlier work has mostly focused on the *correct rate* of LLMs in few-shot settings without unsure shots. Based on the results presented in Table 9 in Appendix D, we observe



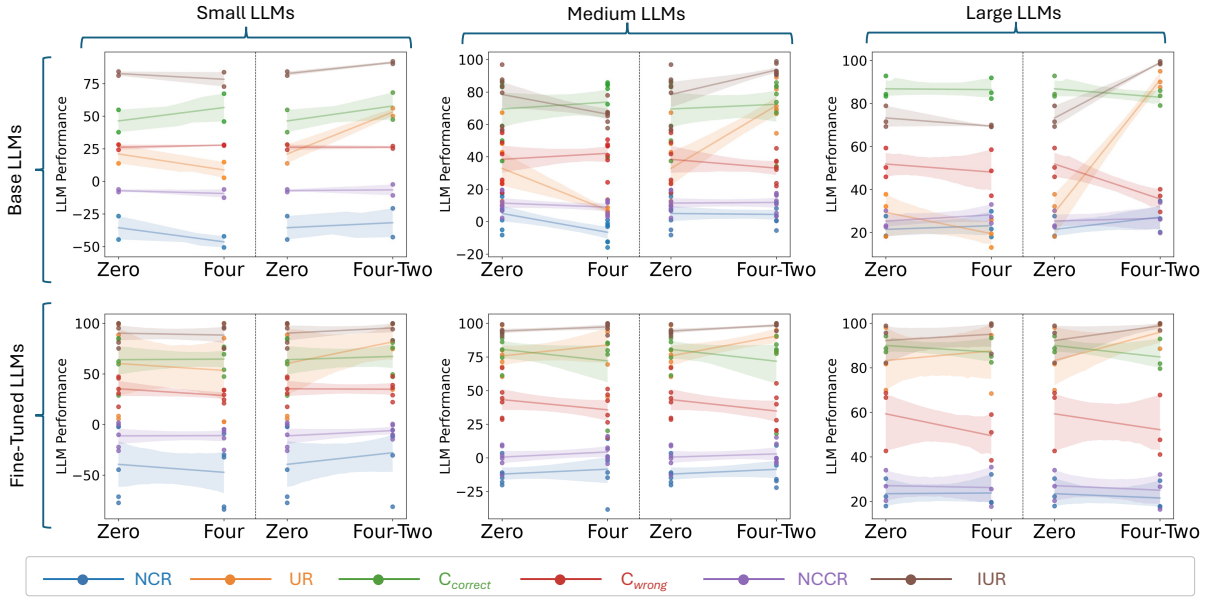


Figure 9: The impact of ICL on LLM performance, measured with NCR, UR,  $C_{correct}$ ,  $C_{wrong}$ , NCCR, and IUR. Different metrics are color-coded. We compare zero-shot and four-shot settings; and zero-shot against four-shot with two unsure shots. LLMs (‘base’ and fine-tuned ones) are in three sizes, small, medium, and large.

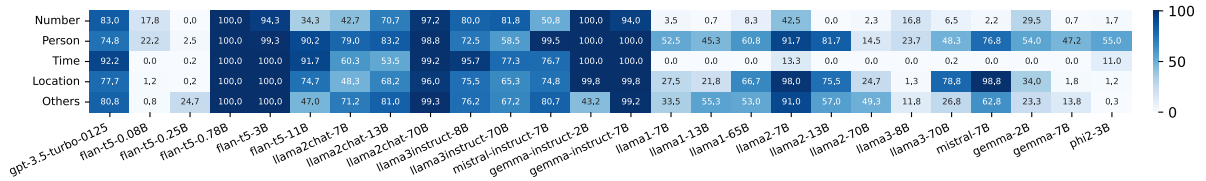


Figure 10: The impact of question type on LLM performance as measured by Uninformative Rate (UR) on unseen knowledge. Questions are grouped into types based on answer they require. Higher values have darker shades.

that for most LLMs, the four-shot setting indeed substantially improves the correct rate but also increases the wrong rate. Incorporating unsure shots helps decrease the wrong rate but also decreases the correct rate.

**ICL makes LLMs less consistent, even on correct responses!** Figure 9 shows that ICL leads to reductions in both  $C_{correct}$  (green line) and  $C_{wrong}$  (red line) for medium/large fine-tuned LLMs and large base LLMs. While a lower  $C_{wrong}$  is desirable, the decrease in  $C_{correct}$  is concerning. This indicates that ICL negatively affects the reliability of LLMs at providing correct responses.

#### 5.4 Model Behavior on Unseen Knowledge

**Base LLMs overestimate their knowledge on numerical and temporal questions.** Figure 10 reports UR performance across models, broken down per question type (e.g., number, person, location). As can be seen, most base LLMs exhibit

significantly lower UR performance on questions requiring numerical or temporal responses. This indicates that base LLMs often provide misleading information for such questions, even when they lack the relevant knowledge.

**Fine-tuned LLMs can explicitly admit when they don’t know.** As illustrated in Figure 11 (fourth column), fine-tuned LLMs are able to explicitly acknowledge their lack of knowledge by answering ‘unsure’ (see definition in Section 3.3) to questions about unseen knowledge. In contrast, base LLMs often produce responses classified as ‘none’ or ‘repetition’ in the absence of unsure shots (contrast columns one and two with column three in Figure 11).

## 6 Related Work

Petroni et al. (2019) were the first to introduce the concept of using pre-trained language models as knowledge bases. They introduced LAMA,

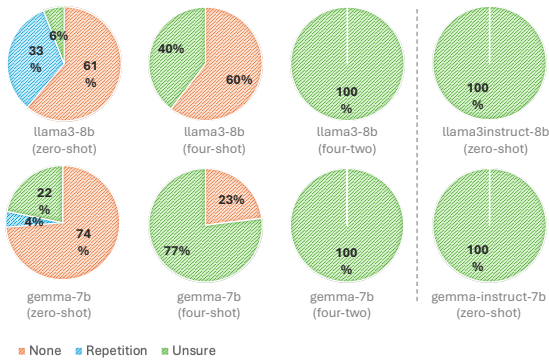


Figure 11: Distribution of uninformative responses given by LLMs to questions about unseen knowledge. We report results for the LLAMA3-8B, GEMMA-7B, and their fine-tuned models (fourth column) but observe similar trends on other models (omitted for the sake of brevity).

a benchmark with questions structured as "fill-in-the-blank" cloze statements, and found that BERT (Devlin et al., 2019) retains relational knowledge competitive with traditional NLP methods that have some access to oracle knowledge.

Roberts et al. (2020) were the first to measure the extent to which language models trained on unstructured text can implicitly store and retrieve knowledge using natural language queries. Specifically, they fine-tuned T5 (Raffel et al., 2019) to answer questions without access to any external context or knowledge and showed that this approach performs competitively with open-domain systems that explicitly retrieve answers from an external knowledge source when answering questions. On a similar vein, Wang et al. (2021) fine-tuned BART (Lewis et al., 2020) with related passages to instill factual knowledge and used the accuracy of masked span recovery to measure the extent to which this knowledge was memorized by the model. They found it was challenging for pre-trained language models like BART to remember facts seen in training, and answer questions, even in cases where the relevant knowledge was retained. He et al. (2024) explicitly trained T5 (Raffel et al., 2019) and LLAMA2 (Touvron et al., 2023) to memorize world knowledge from Wikidata on a large scale and then used exact match and F1 scores to evaluate knowledge retention. Their results showed LLMs hold promise as large-scale KBs capable of retrieving facts and responding with flexibility, but are less proficient at inferring new knowledge through reasoning.

Sun et al. (2023) proposed a benchmark consisting of 18,000 question-answer pairs representing facts with varying popularity (i.e., high, medium, low) and assessed the knowledge retained by various language models on this benchmark. Unlike previous studies, Sun et al. (2023) reported both accuracy (the percentage of questions answered correctly) and hallucination rate (the percentage of questions answered wrongly). They showed models are particularly bad at answering questions from medium- and low-frequency facts.

Our work also examines whether LLMs can be used in lieu of more traditional KBs. Compared to previous research, we proposed a comprehensive evaluation framework which assesses not only the ability of LLMs to recall *seen* knowledge but also their ability to respond in the face of *unseen* knowledge. In addition, we evaluate LLM consistency when answering questions about identical knowledge.

## 7 Conclusion

In this paper, we defined a set of criteria that LLMs functioning as KBs should meet, focusing on factuality and consistency. We proposed various metrics operationalizing these criteria and used them to assess LLM performance when answering questions pertaining to both seen and unseen knowledge. We evaluated 26 popular LLMs and found that GPT-3.5-TURBO is the most reliable among them. Additionally, we examined the impact of model size, fine-tuning, and ICL. Fine-tuning and ICL with unsure shots were shown to improve LLM capabilities when providing responses to questions about unseen knowledge, but do not significantly improve performance on seen knowledge. Increasing model size boosts performance on seen knowledge but at the expense of performance on unseen knowledge. Notably, neither ICL nor fine-tuning were successful at improving LLM consistency. This highlights the critical need for continued research to develop more robust strategies that ensure both factuality and consistency, enabling LLMs to effectively function as KBs.

## References

Badr AlKhamissi, Millicent Li, Asli Celikyilmaz, Mona Diab, and Marjan Ghazvininejad. 2022.

- A review on language models as knowledge bases. *arXiv preprint arXiv:2204.06031*.
- Kim Allan Andersen and Daniele Pretolani. 2001. Easy cases of probabilistic satisfiability. *Annals of Mathematics and Artificial Intelligence*, 33:69–91.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Rakesh Chada and Pradeep Natarajan. 2021. [Few-shotQA: A simple framework for few-shot learning of question answering tasks using pre-trained text-to-text models](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6081–6090, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Shiqi Chen, Yiran Zhao, Jinghan Zhang, I-Chun Chern, Siyang Gao, Pengfei Liu, and Junxian He. 2023. Felm: Benchmarking factuality evaluation of large language models. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R Glass, and Pengcheng He. 2023. Dola: Decoding by contrasting layers improves factuality in large language models. In *The Twelfth International Conference on Learning Representations*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Bhuwan Dhingra, Jeremy R Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W Cohen. 2022. Time-aware language models as temporal knowledge bases. *Transactions of the Association for Computational Linguistics*, 10:257–273.
- Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Schütze, and Yoav Goldberg. 2021. Measuring and improving consistency in pretrained language models. *Transactions of the Association for Computational Linguistics*, 9:1012–1031.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. 2023. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*.
- Lovisa Hagström, Denitsa Saynova, Tobias Norlund, Moa Johansson, and Richard Johansson. 2023. The effect of scaling, retrieval augmentation and form on the factual consistency of language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5457–5476.
- Qiyuan He, Yizhong Wang, and Wenya Wang. 2024. Can language models act as knowledge bases at scale? *arXiv preprint arXiv:2402.14273*.
- Benjamin Heinzerling and Kentaro Inui. 2021. Language models as knowledge bases: On entity representations, storage capacity, and paraphrased queries. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1772–1791.
- Myeongjun Jang, Deuk Sin Kwon, and Thomas Lukasiewicz. 2022. Becel: Benchmark for consistency evaluation of language models. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3680–3696.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large

- scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. **BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khoshdel, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9802–9822.
- Vishwas Mruthyunjaya, Pouya Pezeshkpour, Estevam Hruschka, and Nikita Bhutani. 2023. Rethinking language models as symbolic knowledge graphs. *arXiv preprint arXiv:2308.13676*.
- Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, et al. 2023. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. *arXiv preprint arXiv:2311.16452*.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473.
- Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. **Exploring the limits of transfer learning with a unified text-to-text transformer**. *Journal of Machine Learning Research*, 21:140:1–140:67.
- Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How much knowledge can you pack into the parameters of a language model? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5418–5426.
- Kai Sun, Yifan Ethan Xu, Hanwen Zha, Yue Liu, and Xin Luna Dong. 2023. Head-to-tail: How knowledgeable are large language models (llm)? aka will llms replace knowledge graphs? *arXiv preprint arXiv:2308.10168*.
- Mujeen Sung, Jinhyuk Lee, Sean Yi, Minji Jeon, Sungdong Kim, and Jaewoo Kang. 2021. Can language models be biomedical knowledge bases? *arXiv preprint arXiv:2109.07154*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Chenguang Wang, Xiao Liu, and Dawn Song. 2020. Language models are open knowledge graphs. *arXiv preprint arXiv:2010.11967*.
- Cunxiang Wang, Pai Liu, and Yue Zhang. 2021. Can generative pre-trained language models serve as knowledge bases for closed-book qa? In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3241–3251.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*.

Yuxia Wang, Minghan Wang, Muhammad Arslan Manzoor, Georgi Georgiev, Rocktim Jyoti Das, and Preslav Nakov. 2024. Factuality of large language models in the year 2024. *arXiv preprint arXiv:2402.02420*.

Jifan Yu, Xiaozhi Wang, Shangqing Tu, Shulin Cao, Daniel Zhang-Li, Xin Lv, Hao Peng, Zijun Yao, Xiaohan Zhang, Hanming Li, et al. 2023. Kola: Carefully benchmarking world knowledge of large language models. In *The Twelfth International Conference on Learning Representations*.

Danna Zheng, Danyang Liu, Mirella Lapata, and Jeff Z Pan. 2024. Trustscore: Reference-free evaluation of llm response trustworthiness. *arXiv preprint arXiv:2402.12545*.

## A Datasets

### A.1 SeenQA Dataset

SeenQA is composed of questions selected from the following datasets:

1. Natural Questions (Kwiatkowski et al., 2019): This dataset includes questions sourced from web queries, each paired with a corresponding Wikipedia article containing the answer. The paper on Natural Questions was submitted to TACL in April 2018.
2. TriviaQA (Joshi et al., 2017): This dataset comprises questions from Quiz League websites, supplemented by web pages and Wikipedia searches that may contain the answer. The paper on TriviaQA was submitted to Arxiv in May 2017. For this project, we focus only on questions supported by Wikipedia.
3. PopQA (Mallen et al., 2023): This dataset targets long-tail entities. The authors used the Wikipedia dump from December 2018 in the retrieval augmented baseline, indicating that the knowledge in PopQA can be covered by the Wikipedia dump from that date.

Wikipedia is a common source in the pre-training data of large language models (LLMs). Comparing the knowledge cutoff dates provided in Table 4, we can deduce that the knowledge involved in these three datasets must have been seen during training by the LLMs used in our study. SeenQA contains 3000 questions selected through a three-step process:

1. Factoid Question Extraction: We filter out questions starting with "why," those with multiple different answers, or answers longer than five tokens.
2. Removal of Time-Sensitive Questions: We use GPT-4-1106-PREVIEW to detect questions whose answers may change over time, using the prompt shown in Table 5. Such questions are then removed from the dataset.
3. Sampling: Finally, we randomly sample 1,000 questions from each dataset’s test or development set, if the test set is unavailable, discarding the supporting context to adapt to a closed-book setting.

### A.2 UnseenQA Dataset

UnseenQA was created using 20 templates, shown in Table 3, covering five answer types: number, person, time, location, and others. These templates were hand-written by the authors. Templates T1, T2, T5, T6, and T7 are about future events, so their answers are unknown at the time of writing. The remaining templates involve fictional persons whose names and related information cannot be found on the Internet. The templates include three types of placeholders:

1. Country/Region: We collected 150 country/region names from the list of National Olympic Committees on a [Wikipedia page](#).
2. Medal Event: We collected 150 medal events from the [official programme of the Olympic games, Paris 2024](#).
3. Person: We randomly sampled 150 person names from combinations of 100 first names, 100 middle names, and 100 last names. We manually checked these names on Google to ensure they have no Internet presence.

UnseenQA consists of 3000 questions, with 150 questions per template. This dataset was created on April 13, 2024. Therefore, LLMs trained on knowledge before this date should not know the answers.

## B LLMs used in this project

Table 4 summarizes the LLMs used in our experiments.

| Answer Type | Abb | Template  |
|-------------|-----|---|
| Number      | T1  | How many gold medals did [country/region] win at the XXXIV Summer Olympic Games?  |
|             | T2  | In the 25th FIFA World Cup, what was the final ranking of [country/region]?       |
|             | T3  | How many children does [person] have?   |
|             | T4  | How old was [person] in 2015?   |
| Person      | T5  | Who won the bronze medal of [medal event] at the XXXIII Summer Olympic Games?     |
|             | T6  | Who is the supreme leader of [country/region] in 2040?                            |
|             | T7  | In 2028, who served as the head coach of [country/region] national football team? |
|             | T8  | Who is [person]’s mom?  |
| Time        | T9  | On which date was [person] born?  |
|             | T10 | In what year did [person] die?  |
|             | T11 | In what year did [person] graduate with the bachelor’s degree?                    |
|             | T12 | When was the wedding date for [person]?   |
| Location    | T13 | Where was [person] born?  |
|             | T14 | Where did [person] pass away?   |
|             | T15 | Which university did [person] attend for the undergraduate studies?               |
|             | T16 | Where was [person]’s wedding held?  |
| Others      | T17 | What was the cause of [person]’s death?   |
|             | T18 | What is the title of the debut album released by [person]?                        |
|             | T19 | What is the name of the first film directed by [person]?                          |
|             | T20 | What is the occupation of [person]?   |

Table 3: Question templates used to create UnseenQA

## C Prompts used in our Experiments

The prompt used to detect time-sensitive questions is shown in Table 5. The QA prompts under three different prompt settings are shown in Table 6. The prompt used to check whether an LLM’s response matches the ground-truth is shown in Table 7. The prompt used to generate distractors for the consistency test is shown in Table 8.

## D Full Experimental Results

Table 9, Table 10, and Table 11 provides the detailed results of LLMs’ performance on factuality, consistency, and reliability respectively. Figure 12 shows the rankings of LLMs based on different metrics. Figure 13 compare different LLMs’ factuality, consistency, and reliability performance.

| Models                             | #Params | Type     | Open Source | Fine-Tuning |      | Release Date  | Pre-Training      |         |       |
|------------------------------------|---------|----------|-------------|-------------|------|---------------|-------------------|---------|-------|
|                                    |         |          |             | IT          | RLHF |               | Knowledge         | # Token | Vocab |
| <a href="#">gpt-3.5-turbo-0125</a> | Unknown | Dec-only | ✗           | ✓           | ✓    | 25 Jan 2024   | Sep 2021          | -       | -     |
| <a href="#">Flan-T5</a>            | 0.08B   | Enc-Dec  | ✓           | ✓           | ✗    | 20 Oct 2022   | <u>April 2019</u> | Unknown | 32K   |
|                                    | 0.25B   | Enc-Dec  | ✓           | ✓           | ✗    | 20 Oct 2022   | <u>April 2019</u> | Unknown | 32K   |
|                                    | 0.78B   | Enc-Dec  | ✓           | ✓           | ✗    | 20 Oct 2022   | <u>April 2019</u> | Unknown | 32K   |
|                                    | 3B      | Enc-Dec  | ✓           | ✓           | ✗    | 20 Oct 2022   | <u>April 2019</u> | Unknown | 32K   |
|                                    | 11B     | Enc-Dec  | ✓           | ✓           | ✗    | 20 Oct 2022   | <u>April 2019</u> | Unknown | 32K   |
| <a href="#">Llama1</a>             | 7B      | Dec-only | ✓           | ✗           | ✗    | 27 Feb 2023   | <u>Aug 2022</u>   | 1T      | 32K   |
|                                    | 13B     | Dec-only | ✓           | ✗           | ✗    | 27 Feb 2023   | <u>Aug 2022</u>   | 1T      | 32K   |
|                                    | 65B     | Dec-only | ✓           | ✗           | ✗    | 27 Feb 2023   | <u>Aug 2022</u>   | 1.4T    | 32K   |
| <a href="#">Llama2</a>             | 7B      | Dec-only | ✓           | ✗           | ✗    | 18 July 2023  | Sep 2022          | 2T      | 32K   |
|                                    | 13B     | Dec-only | ✓           | ✗           | ✗    | 18 July 2023  | Sep 2022          | 2T      | 32K   |
|                                    | 70B     | Dec-only | ✓           | ✗           | ✗    | 18 July 2023  | Sep 2022          | 2T      | 32K   |
| <a href="#">Llama2chat</a>         | 7B      | Dec-only | ✓           | ✓           | ✓    | 18 July 2023  | Sep 2022          | 2T      | 32K   |
|                                    | 13B     | Dec-only | ✓           | ✓           | ✓    | 18 July 2023  | Sep 2022          | 2T      | 32K   |
|                                    | 70B     | Dec-only | ✓           | ✓           | ✓    | 18 July 2023  | Sep 2022          | 2T      | 32K   |
| <a href="#">Llama3</a>             | 8B      | Dec-only | ✓           | ✗           | ✗    | 18 April 2024 | Mar 2023          | 15T+    | 128K  |
|                                    | 70B     | Dec-only | ✓           | ✗           | ✗    | 18 April 2024 | Dec 2023          | 15T+    | 128K  |
| <a href="#">Llama3Instruct</a>     | 8B      | Dec-only | ✓           | ✓           | ✓    | 18 April 2024 | Mar 2023          | 15T+    | 128K  |
|                                    | 70B     | Dec-only | ✓           | ✓           | ✓    | 18 April 2024 | Dec 2023          | 15T+    | 128K  |
| <a href="#">Mistral</a>            | 7B      | Dec-only | ✓           | ✗           | ✗    | 27 Sep 2023   | Unknown           | Unknown | 32K   |
| <a href="#">Mistral-Instruct</a>   | 7B      | Dec-only | ✓           | ✓           | ✗    | 27 Sep 2023   | Unknown           | Unknown | 32K   |
| <a href="#">Gemma</a>              | 2B      | Dec-only | ✓           | ✗           | ✗    | 21 Feb 2024   | Unknown           | 3T      | 256K  |
|                                    | 7B      | Dec-only | ✓           | ✗           | ✗    | 21 Feb 2024   | Unknown           | 6T      | 256K  |
| <a href="#">Gemma-Instruct</a>     | 2B      | Dec-only | ✓           | ✓           | ✓    | 21 Feb 2024   | Unknown           | 3T      | 256K  |
|                                    | 7B      | Dec-only | ✓           | ✓           | ✓    | 21 Feb 2024   | Unknown           | 6T      | 256K  |
| <a href="#">Phi2</a>               | 3B      | Dec-only | ✓           | ✗           | ✗    | 12 Dec 2023   | Unknown           | 1.4T    | 50K   |

Table 4: Summary of LLMs used in our experiments. ‘IT’ denotes Instruction Tuning, and ‘RLHF’ refers to Reinforcement Learning from Human Feedback. ‘Knowledge’ indicates the knowledge cutoff date. Underlined dates were not explicitly provided by the authors but extrapolated from the datasets used for LLM training. Flan-T5’s base model is T5 version 1.1 pre-trained on the C4 dataset, filtered from web-extracted text in April 2019. Llama 1’s pre-training data includes Wikipedia dumps from June to August 2022.

---

### Prompt for detecting time-sensitive questions

---

INSTRUCTION: Please provide the index of questions whose answers change yearly. Just return the index without explanations.

Here is the list of questions:

1. Who is the most paid player in EPL?
2. What is the capital of Louisiana?
3. Who won the Nobel Peace Prize in 2009?
4. What is the latest model of the iPhone currently available?

Index:

1, 4

Here is the list of questions:

[question placeholder]

Index:

---

Table 5: The prompt for detecting time-sensitive questions



---

**QA prompt in zero-shot**

INSTRUCTION: Please answer knowledge-related questions directly. Note: Please do not give anything other than the answer; Say "unsure" if you do not know.

QUESTION: [question placeholder]

ANSWER:

---

**QA prompt in four-shot**

INSTRUCTION: Please answer knowledge-related questions directly. Note: Please do not give anything other than the answer; Say "unsure" if you do not know.

QUESTION: [question example 1 from  $R_{seen}$ ]

ANSWER: [answer 1]

QUESTION: [question example 2 from  $R_{seen}$ ]

ANSWER: [answer 2]

QUESTION: [question example 3 from  $R_{seen}$ ]

ANSWER: [answer 3]

QUESTION: [question example 4 from  $R_{seen}$ ]

ANSWER: [answer 4]

QUESTION: [question placeholder]

ANSWER:

---

**QA prompt in four-shot with few unsure shot**

INSTRUCTION: Please answer knowledge-related questions directly. Note: Please do not give anything other than the answer; Say "unsure" if you do not know.

QUESTION: [question example 1 from  $R_{seen}$ ]

ANSWER: [answer 1]

QUESTION: [question example 2 from  $R_{seen}$ ]

ANSWER: [answer 2]

QUESTION: [question example 3 from  $R_{unseen}$ ]

ANSWER: unsure

QUESTION: [question example 4 from  $R_{unseen}$ ]

ANSWER: unsure

QUESTION: [question placeholder]

ANSWER:

---

Table 6: The question answering prompt format. The shots are selected from repositories,  $R_{seen}$  and  $R_{unseen}$ . The order of shots is random. For the MCQ tests in consistency experiments, we edit the instruction line to *INSTRUCTION: Please answer knowledge-related multi-choice questions directly. Note: Please do not give anything other than the appropriate option (A, B, C, D or E); choose the option indicating "unsure" if you do not know.*

---

**Prompt for check whether an answer matches the ground truth for the question**

---

INSTRUCTION: You need to check whether the prediction of a question-answering system to a question is correct. You should make the judgment based on a list of ground truth answers provided to you. Your response should be "yes" if the prediction is correct or "no" if the prediction is wrong.

Question: Who authored The Taming of the Shrew (published in 2002)?

Ground truth: ["William Shakespeare", "Roma Gill"]

prediction: W Shakespeare

Correctness: yes

Question: What country is Maharashtra Metro Rail Corporation Limited located in?

Ground truth: ["India"]

prediction: Maharashtra

Correctness: no

Question: Edward Tise (known for Full Metal Jacket (1987)) is in what department?

Ground truth: ["sound department"]

Prediction: 2nd Infantry Division, United States Army

Correctness: no

Question: Which era did Michael Oakeshott belong to?

Ground truth: ["20th-century philosophy"]

prediction: 20th century.

Correctness: yes

Question: [quesetion placeholder]

Ground truth: [ground truth placeholder]

prediction: [LLM's answer placeholder]

Correctness:

---

Table 7: The prompt used to check whether an LLM's answer matches the ground truth for the question

---

**Prompt for generating distractors**

---

INSTRUCTION: For the given question-answer pair, provide 20 different distractors that are similar yet distinct from the given answer. Note: Seperate the 20 distractors with a special token "[SEP]".

Q: Who was the President of the United States in 2010?

A: Barack Obama

Distractors: George W. Bush [SEP] Bill Clinton [SEP] Ronald Reagan [SEP] Donald Trump [SEP] Jimmy Carter [SEP] George H.W. Bush [SEP] Richard Nixon [SEP] Gerald Ford [SEP] Lyndon B. Johnson [SEP] John F. Kennedy [SEP] Dwight D. Eisenhower [SEP] Harry S. Truman [SEP] Franklin D. Roosevelt [SEP] Herbert Hoover [SEP] Calvin Coolidge [SEP] Woodrow Wilson [SEP] William Howard Taft [SEP] Theodore Roosevelt [SEP] William McKinley [SEP] Grover Cleveland

Q: What is the name of the first cloned sheep?

A: the first cloned sheep is dolly.

Distractors: the first cloned sheep is Polly [SEP] the first cloned sheep is Molly [SEP] the first cloned sheep is Holly [SEP] the first cloned sheep is Bella [SEP] the first cloned sheep is Daisy [SEP] the first cloned sheep is Lily [SEP] the first cloned sheep is Rosie [SEP] the first cloned sheep is Millie [SEP] the first cloned sheep is Ellie [SEP] the first cloned sheep is Sally [SEP] the first cloned sheep is Tilly [SEP] the first cloned sheep is Nelly [SEP] the first cloned sheep is Jolly [SEP] the first cloned sheep is Betty [SEP] the first cloned sheep is Annie [SEP] the first cloned sheep is Lucy [SEP] the first cloned sheep is Maggie [SEP] the first cloned sheep is Cindy [SEP] the first cloned sheep is Penny [SEP] the first cloned sheep is Ginny

Q: [QUESTION]

A: [ANSWER]

Distractors:

---

Table 8: The prompt used to generate distractors for consistency tests.

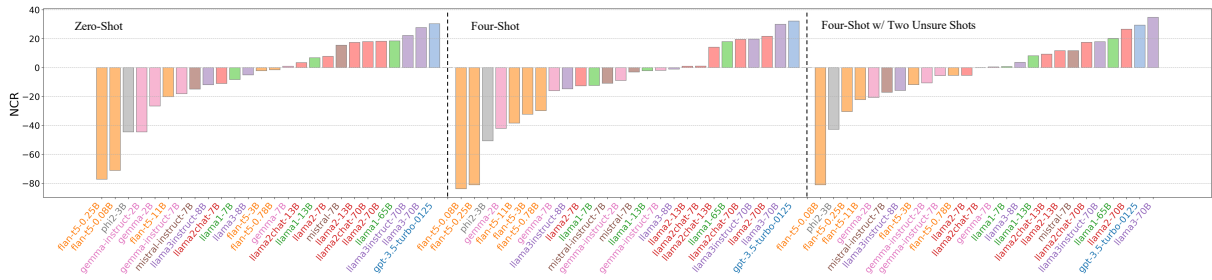
| Model            | Params  | zero-shot |        |         |        | four-shot |        |         |        | four-shot-2 |        |         |        |
|------------------|---------|-----------|--------|---------|--------|-----------|--------|---------|--------|-------------|--------|---------|--------|
|                  |         | Seen      |        | Uneen   |        | Seen      |        | Uneen   |        | Seen        |        | Uneen   |        |
|                  |         | WR (↓)    | CR (↑) | NCR (↑) | UR (↑) | WR (↓)    | CR (↑) | NCR (↑) | UR (↑) | WR (↓)      | CR (↑) | NCR (↑) | UR (↑) |
| GPT-3.5 Turbo    | Unknown | 30.40     | 60.73  | 30.33   | 81.70  | 28.97     | 61.17  | 32.20   | 94.70  | 27.93       | 57.30  | 29.37   | 99.37  |
| Flan-T5          | 0.08B   | 73.03     | 1.83   | -71.20  | 8.40   | 85.53     | 1.63   | -83.90  | 2.77   | 82.57       | 1.43   | -81.13  | 34.63  |
|                  | 0.25B   | 82.77     | 5.47   | -77.30  | 5.50   | 86.43     | 5.27   | -81.17  | 2.70   | 32.67       | 2.23   | -30.43  | 74.67  |
|                  | 0.78B   | 3.70      | 2.13   | -1.57   | 100.00 | 37.80     | 8.00   | -29.80  | 85.30  | 9.60        | 4.17   | -5.43   | 99.90  |
|                  | 3B      | 9.70      | 7.50   | -2.20   | 98.73  | 46.00     | 13.67  | -32.33  | 76.27  | 23.03       | 11.23  | -11.80  | 99.73  |
|                  | 11B     | 40.97     | 20.77  | -20.20  | 67.57  | 60.63     | 22.23  | -38.40  | 45.27  | 42.37       | 20.20  | -22.17  | 90.00  |
| Llama 1          | 7B      | 43.93     | 35.67  | -8.27   | 23.40  | 54.47     | 42.10  | -12.37  | 4.73   | 27.57       | 28.27  | 0.70    | 54.50  |
|                  | 13B     | 34.33     | 41.13  | 6.80    | 24.63  | 49.67     | 47.43  | -2.23   | 4.10   | 27.37       | 35.53  | 8.17    | 69.37  |
|                  | 65B     | 28.40     | 46.87  | 18.47   | 37.77  | 39.77     | 57.73  | 17.97   | 19.47  | 14.63       | 34.87  | 20.23   | 90.10  |
| Llama 2          | 7B      | 23.73     | 31.50  | 7.77    | 67.30  | 54.63     | 42.03  | -12.60  | 5.33   | 35.53       | 30.10  | -5.43   | 66.60  |
|                  | 13B     | 23.67     | 41.07  | 17.40   | 42.83  | 48.20     | 49.17  | 0.97    | 6.03   | 27.40       | 39.07  | 11.67   | 71.23  |
|                  | 70B     | 37.17     | 55.33  | 18.17   | 18.17  | 38.20     | 59.83  | 21.63   | 13.03  | 19.50       | 46.03  | 26.53   | 95.10  |
| Llama2chat       | 7B      | 47.37     | 36.33  | -11.03  | 60.30  | 26.87     | 27.90  | 1.03    | 98.60  | 23.43       | 23.33  | -0.10   | 99.73  |
|                  | 13B     | 37.87     | 41.27  | 3.40    | 71.30  | 25.03     | 39.13  | 14.10   | 96.13  | 32.00       | 41.27  | 9.27    | 94.90  |
|                  | 70B     | 29.53     | 47.50  | 17.97   | 98.10  | 14.57     | 34.10  | 19.53   | 99.63  | 15.13       | 32.60  | 17.47   | 100.00 |
| Llama3           | 8B      | 50.20     | 45.13  | -5.07   | 10.73  | 49.27     | 48.23  | -1.03   | 6.50   | 32.77       | 36.33  | 3.57    | 89.03  |
|                  | 70B     | 27.87     | 55.53  | 27.67   | 32.13  | 33.70     | 63.60  | 29.90   | 25.93  | 18.87       | 53.60  | 34.73   | 87.63  |
| Llama3Instruct   | 8B      | 53.93     | 42.03  | -11.90  | 79.97  | 54.00     | 39.27  | -14.73  | 69.43  | 54.60       | 38.73  | -15.87  | 78.73  |
|                  | 70B     | 36.80     | 59.03  | 22.23   | 70.03  | 38.40     | 58.10  | 19.70   | 68.47  | 38.90       | 56.80  | 17.90   | 88.60  |
| Mistral          | 7B      | 24.00     | 39.47  | 15.47   | 48.13  | 50.13     | 47.07  | -3.07   | 13.57  | 24.70       | 36.37  | 11.67   | 81.73  |
|                  | 7B      | 44.77     | 29.90  | -14.87  | 76.50  | 39.53     | 28.63  | -10.90  | 93.80  | 46.63       | 29.47  | -17.17  | 79.13  |
| Mistral-Instruct | 2B      | 51.20     | 24.63  | -26.57  | 28.17  | 69.17     | 27.07  | -42.10  | 2.77   | 39.37       | 18.67  | -20.70  | 56.07  |
|                  | 7B      | 38.80     | 39.73  | 0.93    | 12.70  | 56.50     | 40.53  | -15.97  | 8.67   | 30.77       | 31.23  | 0.47    | 68.93  |
| Gemma            | 2B      | 53.80     | 9.27   | -44.53  | 88.60  | 13.27     | 4.30   | -8.97   | 99.93  | 14.40       | 3.77   | -10.63  | 99.30  |
|                  | 7B      | 37.13     | 19.03  | -18.10  | 98.60  | 16.17     | 14.20  | -1.97   | 99.97  | 19.13       | 13.60  | -5.53   | 99.93  |
| Phi2             | 3B      | 65.97     | 21.43  | -44.53  | 13.83  | 72.10     | 21.40  | -50.70  | 14.77  | 62.07       | 19.33  | -42.73  | 50.20  |

Table 9: Factuality performance.

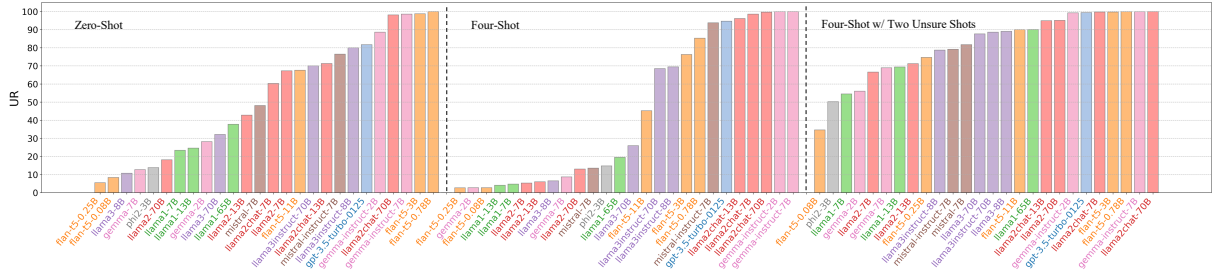
| Model          | Params | zero-shot              |                          |                        |                          | four-shot              |                          |                        |                          | four-shot-2            |                          |                        |                          |
|----------------|--------|------------------------|--------------------------|------------------------|--------------------------|------------------------|--------------------------|------------------------|--------------------------|------------------------|--------------------------|------------------------|--------------------------|
|                |        | Seen                   |                          | Uneen                  |                          | Seen                   |                          | Uneen                  |                          | Seen                   |                          | Uneen                  |                          |
|                |        | C <sub>wrong</sub> (↓) | C <sub>correct</sub> (↑) | C <sub>wrong</sub> (↓) | C <sub>correct</sub> (↑) | C <sub>wrong</sub> (↓) | C <sub>correct</sub> (↑) | C <sub>wrong</sub> (↓) | C <sub>correct</sub> (↑) | C <sub>wrong</sub> (↓) | C <sub>correct</sub> (↑) | C <sub>wrong</sub> (↓) | C <sub>correct</sub> (↑) |
| GPT-3.5 Turbo  | -      | 61.79                  | 23.65                    | 42.72                  | 87.10                    | 57.43                  | 19.62                    | 38.53                  | 85.16                    | 48.56                  | 33.68                    | 41.12                  | 79.68                    |
| Flan-T5        | 0.08B  | 14.49                  | 20.56                    | 17.53                  | 28.64                    | 16.53                  | 25.62                    | 21.07                  | 47.45                    | 18.44                  | 25.87                    | 22.15                  | 47.21                    |
|                | 0.25B  | 35.33                  | 26.31                    | 30.82                  | 62.29                    | 33.36                  | 25.44                    | 29.40                  | 69.40                    | 35.34                  | 22.80                    | 29.07                  | 75.90                    |
|                | 0.78B  | 45.68                  | -                        | 45.68                  | 85.23                    | 34.16                  | 33.72                    | 33.94                  | 75.12                    | 42.99                  | 35.00                    | 38.99                  | 82.64                    |
|                | 3B     | 45.03                  | 25.26                    | 35.15                  | 84.20                    | 33.07                  | 16.05                    | 24.56                  | 76.85                    | 37.13                  | 35.62                    | 36.38                  | 80.96                    |
|                | 11B    | 41.62                  | 15.38                    | 28.50                  | 80.43                    | 36.40                  | 16.40                    | 26.40                  | 79.84                    | 40.74                  | 15.07                    | 27.90                  | 80.61                    |
| Llama 1        | 7B     | 25.01                  | 21.70                    | 23.36                  | 37.43                    | 25.37                  | 23.10                    | 24.23                  | 39.65                    | 23.07                  | 20.89                    | 21.98                  | 34.17                    |
|                | 13B    | 35.13                  | 16.41                    | 25.77                  | 59.11                    | 45.60                  | 36.30                    | 40.90                  | 72.49                    | 48.54                  | 25.45                    | 36.99                  | 73.74                    |
|                | 65B    | 58.06                  | 33.84                    | 45.95                  | 83.38                    | 58.35                  | 37.12                    | 47.73                  | 82.38                    | 63.63                  | 16.63                    | 40.13                  | 83.64                    |
| Llama 2        | 7B     | 26.68                  | 9.37                     | 18.03                  | 50.02                    | 41.51                  | 34.56                    | 38.03                  | 67.07                    | 41.94                  | 25.26                    | 33.60                  | 67.29                    |
|                | 13B    | 62.02                  | 35.69                    | 48.86                  | 83.08                    | 66.12                  | 45.05                    | 50.58                  | 82.05                    | 58.55                  | 31.89                    | 45.22                  | 83.65                    |
|                | 70B    | 63.13                  | 37.52                    | 50.33                  | 84.36                    | 62.07                  | 35.37                    | 48.72                  | 85.06                    | 52.84                  | 6.43                     | 29.63                  | 79.10                    |
| Llama2chat     | 7B     | 43.66                  | 15.63                    | 29.64                  | 61.23                    | 17.69                  | 12.50                    | 15.09                  | 30.15                    | 19.82                  | 20.62                    | 20.22                  | 17.99                    |
|                | 13B    | 55.79                  | 32.88                    | 44.33                  | 74.62                    | 56.11                  | 28.97                    | 42.54                  | 76.92                    | 52.23                  | 27.39                    | 39.81                  | 77.52                    |
|                | 70B    | 73.61                  | 59.65                    | 66.63                  | 88.71                    | 71.28                  | 30.91                    | 51.10                  | 82.45                    | 67.82                  | -                        | 67.82                  | 81.88                    |
| Llama3         | 8B     | 64.17                  | 48.52                    | 56.35                  | 86.50                    | 57.43                  | 35.72                    | 46.58                  | 85.85                    | 37.86                  | 9.51                     | 23.69                  | 78.80                    |
|                | 70B    | 76.82                  | 41.82                    | 59.32                  | 92.86                    | 75.47                  | 41.62                    | 58.55                  | 92.00                    | 64.67                  | 9.43                     | 37.05                  | 86.07                    |
| Llama3Instruct | 8B     | 53.08                  | 29.74                    | 41.41                  | 88.86                    | 50.43                  | 13.25                    | 31.84                  | 85.51                    | 37.34                  | 3.80                     | 20.57                  | 79.64                    |
|                | 70B    | 78.14                  | 59.26                    | 68.70                  | 94.24                    | 74.80                  | 43.32                    | 59.06                  | 93.47                    | 67.25                  | 28.17                    | 47.71                  | 93.03                    |
| Mistral        | 7B     | 56.79                  | 26.70                    | 41.75                  | 84.15                    | 55.94                  | 37.50                    | 46.72                  | 84.87                    | 51.19                  | 13.06                    | 32.13                  | 83.21                    |
|                | 7B     | 65.84                  | 31.48                    | 48.66                  | 86.09                    | 62.93                  | 30.99                    | 46.96                  | 84.92                    | 61.64                  | 27.63                    | 44.63                  | 84.21                    |
| Gemma          | 2B     | 30.18                  | 26.26                    | 28.22                  | 37.77                    | 26.93                  | 27.99                    | 27.46                  | 45.90                    | 28.41                  | 22.11                    | 25.26                  | 47.42                    |
|                | 7B     | 62.41                  | 47.24                    | 54.83                  | 86.17                    | 53.39                  | 41.94                    | 47.66                  | 84.22                    | 52.49                  | 22.35                    | 37.42                  | 85.89                    |
| Gemma-Instruct | 2B     | 51.65                  | 42.72                    | 47.19                  | 59.64                    | 53.23                  | 15.00                    | 34.11                  | 54.11                    | 50.51                  | 44.52                    | 47.52                  | 49.51                    |
|                | 7B     | 81.92                  | 51.79                    | 66.85                  | 92.43                    | 72.34                  | 30.00                    | 51.17                  | 84.64                    | 80.70                  | 30.00                    | 55.35                  | 90.56                    |
| Phi2           | 3B     | 30.09                  | 18.27                    | 24.18                  | 54.92                    | 37.21                  | 19.02                    | 28.11                  | 67.34                    | 38.48                  | 15.78                    | 27.13                  | 68.16                    |

Table 10: Consistency performance.

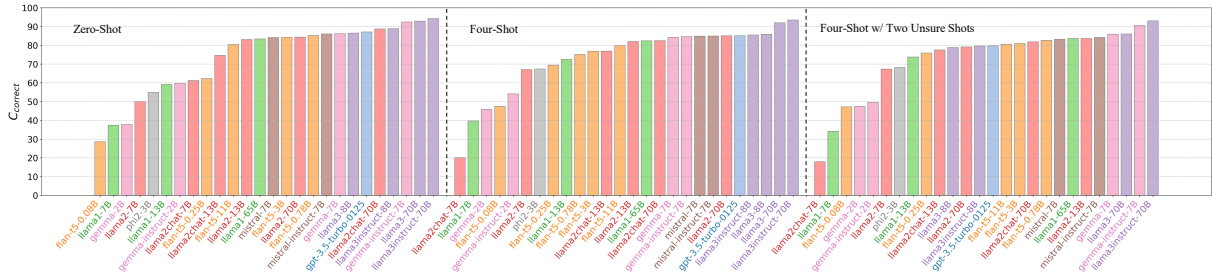
| Model          | Params | zero-shot |         |          |         | four-shot |         |          |         | four-shot-2 |         |          |         |
|----------------|--------|-----------|---------|----------|---------|-----------|---------|----------|---------|-------------|---------|----------|---------|
|                |        | Seen      |         | Uneen    |         | Seen      |         | Uneen    |         | Seen        |         | Uneen    |         |
|                |        | CWR (↓)   | CCR (↑) | NCCR (↑) | IUR (↑) | CWR (↓)   | CCR (↑) | NCCR (↑) | IUR (↑) | CWR (↓)     | CCR (↑) | NCCR (↑) | IUR (↑) |
| GPT-3.5 Turbo  | -      | 18.78     | 52.90   | 34.11    | 95.67   | 16.64     | 52.09   | 35.45    | 98.96   | 13.56       | 45.66   | 32.09    | 99.79   |
| Flan-T5        | 0.08B  | 10.58     | 0.52    | -10.06   | 81.17   | 14.13     | 0.77    | -13.36   | 75.09   | 15.22       | 0.68    | -14.55   | 83.09   |
|                | 0.25B  | 29.25     | 3.41    | -25.84   | 75.13   | 28.83     | 3.66    | -25.17   | 75.25   | 11.54       | 1.69    | -9.85    | 94.22   |
|                | 0.78B  | 1.69      | 1.82    | 0.13     | 100.00  | 12.91     | 6.01    | -6.90    | 95.04   | 4.13        | 3.45    | -0.68    | 99.97   |
|                | 3B     | 4.37      | 6.32    | 1.95     | 99.68   | 15.21     | 10.51   | -4.70    | 96.19   | 8.55        | 9.09    | 0.54     | 99.90   |
|                | 11B    | 17.05     | 16.70   | -0.35    | 95.01   | 22.07     | 17.75   | -4.32    | 91.02   | 17.26       | 16.28   | -0.98    | 98.49   |
| Llama 1        | 7B     | 10.99     | 13.35   | 2.36     | 83.38   | 13.82     | 16.69   | 2.88     | 77.99   | 6.36        | 9.66    | 3.30     | 90.50   |
|                | 13B    | 14.68     | 34.12   | 19.44    | 79.60   | 27.05     | 40.34   | 13.29    | 57.67   | 16.04       | 32.68   | 16.64    | 90.83   |
|                | 70B    | 23.47     | 46.68   | 23.21    | 69.30   | 23.71     | 50.89   | 27.18    | 69.24   | 10.30       | 36.41   | 26.11    | 99.69   |
| Llama2chat     | 7B     | 20.68     | 22.24   | 1.56     | 93.80   | 4.75      | 5.62    | 0.87     | 99.83   | 4.64        | 4.20    | -0.45    | 99.94   |
|                | 13B    | 21.13     | 30.80   | 9.67     | 90.56   | 14.04     | 30.10   | 16.05    | 98.88   | 16.71       | 31.99   | 15.28    | 98.60   |
|                | 70B    | 21.74     | 42.14   | 20.40    | 98.87   | 10.39     | 28.12   | 17.73    | 99.89   | 10.26       | 26.69   | 16.43    | 100.00  |
| Llama3         | 8B     | 32.22     | 39.04   | 6.82     | 56.69   | 28.30     | 41.41   | 13.11    | 66.60   | 12.41       | 28.63   | 16.22    | 98.96   |
|                | 70B    | 21.41     | 51.57   | 30.15    | 71.62   | 25.43     | 58.52   | 33.08    | 69.17   | 12.20       | 46.13   | 33.93    | 98.83   |
| Llama3Instruct | 8B     | 28.62     | 37.35   | 8.72     | 94.04   | 27.23     | 33.58   | 6.34     | 95.95   | 20.39       | 30.85   | 10.46    | 99.19   |
|                | 70B    | 28.76     | 55.63   | 26.88    | 82.24   | 28.72     | 54.31   | 25.59    | 86.34   | 26.16       | 52.84   | 26.68    | 96.79   |
| Mistral        | 7B     | 13.63     | 33.21   | 19.58    | 86.15   | 28.04     | 39.95   | 11.90    | 67.59   | 12.65       | 30.26   | 17.62    | 97.61   |
|                | 7B     | 29.48     | 25.74   | -3.74    | 92.60   | 24.88     | 24.31   | -0.56    | 98.08   | 28.74       | 24.82   | -3.92    | 94.24   |
| Gemma          | 2B     | 15.45     | 9.30    | -6.15    | 81.14   | 18.62     | 12.42   | -6.20    | 72.78   | 11.19       | 8.85    | -2.33    | 90.29   |
|                | 7B     | 24.22     | 34.24   | 10.02    | 58.76   | 30.16     | 34.14   | 3.97     | 61.70   | 16.15       | 26.82   | 10.67    | 93.06   |
| Gemma-Instruct | 2B     | 27.79     | 5.53    | -22.26   | 95.13   | 7.06      | 2.33    | -4.74    | 99.99   | 7.27        | 1.87    | -5.41    | 99.69   |
|                | 7B     | 30.42     | 17.59   | -12.83   | 99.27   | 4.46      | 12.02   | 7.55     | 99.99   | 15.44       | 12.32   | -3.12    | 99.98   |
| Phi2           | 3B     | 19.85     | 11.77   | -8.08    | 84.25   | 26.83     | 14.41   | -12.42   | 83.79   | 23.88       | 13.17   | -10      |         |



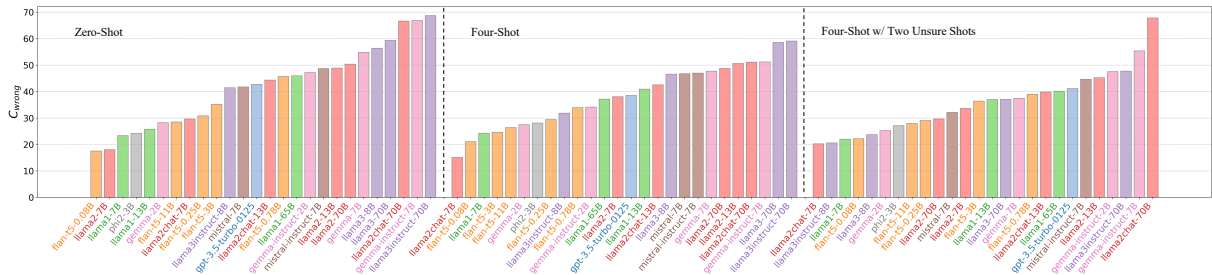
(a) LLMs sorted by NCR



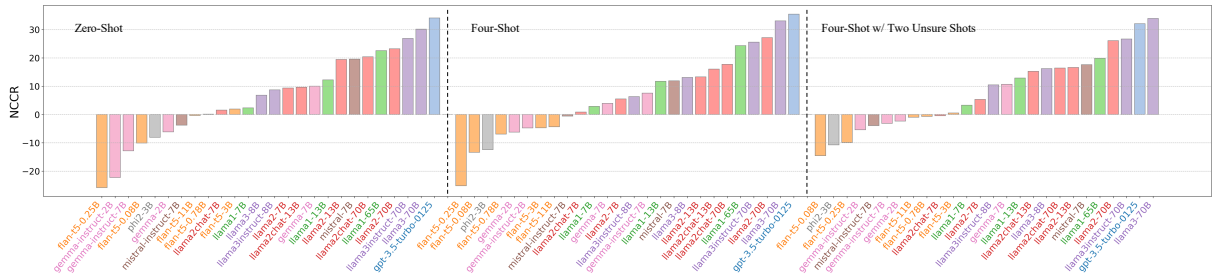
(b) LLMs sorted by UR



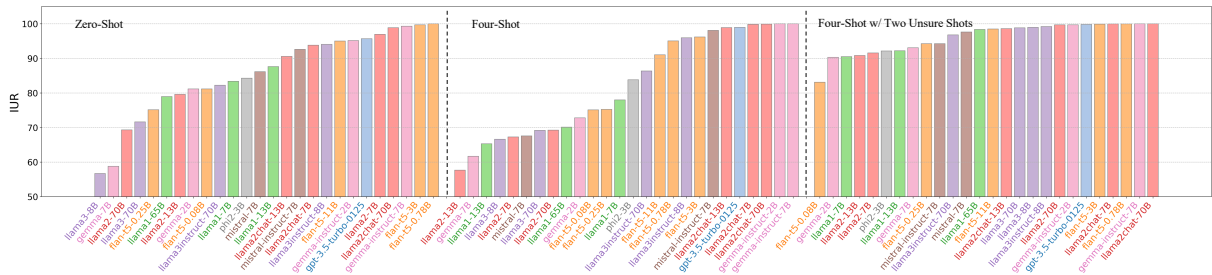
(c) LLMs sorted by C<sub>correct</sub>



(d) LLMs sorted by C<sub>wrong</sub>

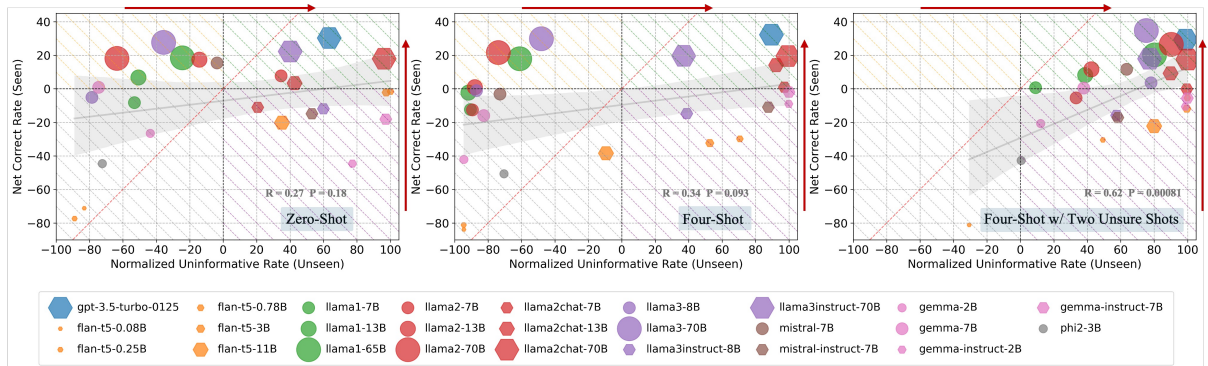


(e) LLMs sorted by NCCR

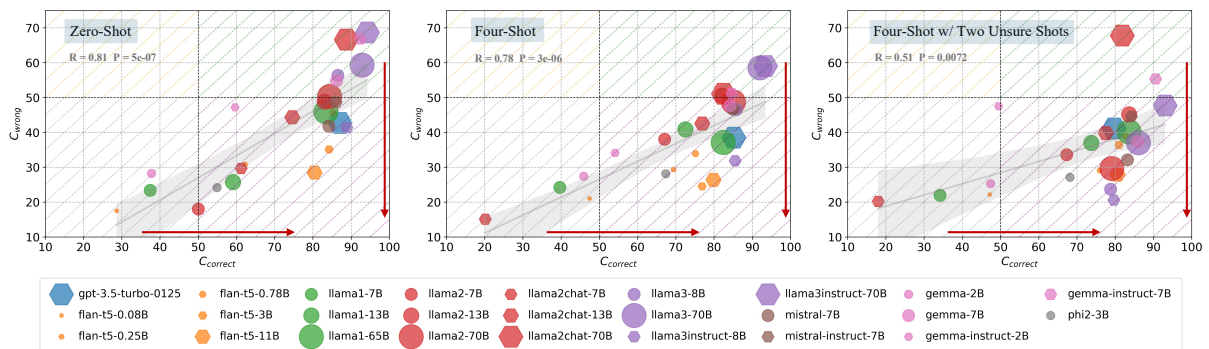


(e) LLMs sorted by IUR

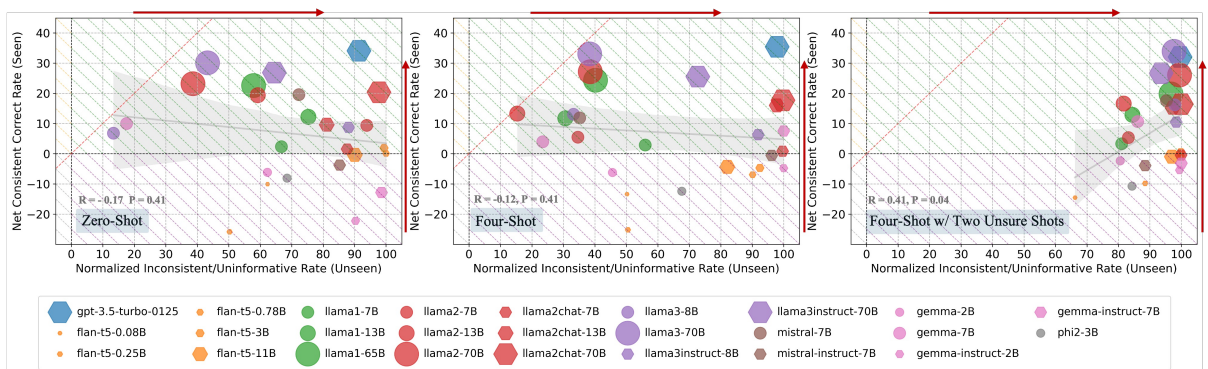
Figure 12: Ranking of LLMs based on different Metrics.



(a) Factuality performance on seen knowledge vs. unseen knowledge. R is the Pearson correlation coefficient. When  $P < 0.05$ , R is statistically significant. The red line is  $y=x$ . The LLMs above the red line perform better on seen knowledge. The LLMs below the red line perform better on unseen knowledge. LLMs closer to the top right corner are more factual (higher NCR and higher UR).



(b) Consistency performance on wrong responses vs. correct responses. R is the Pearson correlation coefficient. When  $P < 0.05$ , R is statistically significant. LLMs closer to the bottom right corner are better in consistency (higher  $C_{correct}$  and lower  $C_{correct}$ ).



(c) Reliability performance on seen knowledge vs. unseen knowledge. R is the Pearson correlation coefficient. When  $P < 0.05$ , R is statistically significant. The red line is  $y=x$ . The LLMs above the red line perform better on seen knowledge. The LLMs below the red line perform better on unseen knowledge. LLMs closer to the top right corner are more reliable (higher NCR and higher UR).

Figure 13: Visualization of LLMs' factuality, consistency and reliability performance.