

# Knowing What LLMs DO NOT Know: A Simple Yet Effective Self-Detection Method

Anonymous ACL submission

## Abstract

Large Language Models (LLMs) have shown great potential in Natural Language Processing (NLP) tasks. However, recent literature reveals that LLMs hallucinate intermittently, which impedes their reliability for further utilization. In this paper, we propose a novel self-detection method to detect which questions an LLM does not know. Our proposal is empirical and applicable for continually upgrading LLMs compared with state-of-the-art methods. Specifically, we examine the divergence of the LLM’s behaviors on different verbalizations for a question and examine the atypicality of the verbalized input. We combine the two components to identify whether the model generates a non-factual response to the question. The above components can be accomplished by utilizing the LLM itself without referring to any other external resources. We conduct comprehensive experiments and demonstrate the effectiveness of our method for recently released LLMs involving Llama 2, Vicuna, ChatGPT, and GPT-4 across factoid question-answering, arithmetic reasoning, and commonsense reasoning tasks.

## 1 Introduction

With the significant improvements in large language models (LLMs) such as PaLM (Chowdhery et al., 2022), ChatGPT (Ouyang et al., 2022), GPT-4 (OpenAI, 2023), LLAMA 2 (Touvron et al., 2023), and Vicuna (Chiang et al., 2023), LLMs have been applied in various natural language tasks. Unfortunately, LLMs still produce unexpected falsehoods (Bang et al., 2023; Li et al., 2023), i.e., they are unaware of what they do not know and generate responses indiscriminately. For example, ChatGPT generates falsehoods for a knowledge quiz and math problem, as shown in Table 1. These intermittent errors can severely hinder the LLMs’ reliability in practice, which makes detecting what they do not know an important research problem (Hendrycks et al., 2021; Lin et al.,

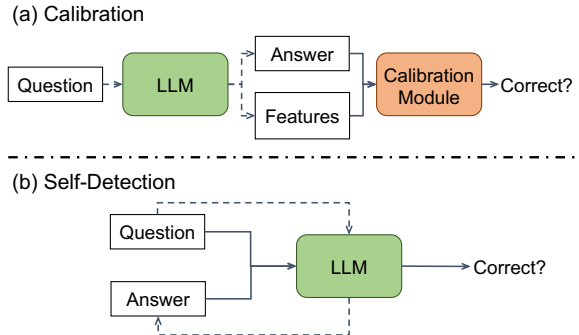


Figure 1: Two paradigms for detecting hallucinations. The dashed lines denote the LLM generation process. The solid lines denote non-factuality detection.

2022; Kadavath et al., 2022).

There are two main paradigms to detect non-factuality: the calibration-based methods and the self-detection methods. The first class of methods calibrates the model confidence to better detect falsehoods of the generations (See Figure 1(a)). Among them, Mielke et al. (2022) train auxiliary calibrators, Lin et al. (2022) and Jiang et al. (2021) improve the calibration through fine-tuning the language model. We propose a self-detection method that does not require further fine-tuning.

The self-detection methods directly leverage the LLMs themselves to detect whether they hallucinate (See Figure 1(b)). For example, Kadavath et al. (2022) prompt the LLMs to predict the confidence score on whether their responses are true, and Si et al. (2023) directly utilizes the token probabilities of the generations as the confidence score; Wang et al. (2023) and Manakul et al. (2023) detect the falsehoods by sampling answers with a high temperature and examining self-consistency between them. However, the performance of these works is limited as LLMs tend to be overconfident about their own outputs and these work would be less effective after the models are trained more aligned (OpenAI, 2023; Ouyang et al., 2022; Zhao et al., 2023).

<b>Knowledge Quiz:</b> The lyricist of <i>Kadam Kadam Badhaye Ja Khushi ke Geet Gaye Ja</i>	
Q1: Who has written the lyrics to patriotic song kadam kadam badhaye ja khushi ke geet gaye ja?	Q2: Who is the lyricist of the patriotic song Kadam Kadam Badhaye Ja Khushi ke Geet Gaye Ja?
R1: The lyrics of the patriotic song Kadam Kadam Badhaye Ja Khushi Ke Geet Gaye Ja were written by <b>Shri Ravi Shankar Sharma</b> .	R2: The lyricist of the patriotic song Kadam Kadam Badhaye Ja Khushi ke Geet Gaye Ja is <b>Shri Pyarelal Santoshi</b> .
Correct Answer: <b>Vanshidhar Shukla</b> .	
<b>Math Problem</b>	
Q1: Tom's restaurant gets 6 reservations a night. They normally order 2 meals that cost \$5 each and a \$5 bottle of wine. How much do they make a week if they are open 2 days a week?	Q2: Kanan's restaurant gets 6 reservations a night. They normally order 2 meals that cost \$5 each and a \$5 bottle of wine. How much do they make a week if they are open 2 days a week?
R1: They make $6 \times 2 = 12$ meals a night. So they make $12 \times 10 = 120$ a night. <b>That means they make <math>120 \times 2 = 240</math> a week.</b>	R2: They make $2 \times 6 = 12$ meals a night. So they make $12 \times 2 = 24$ on meals. They also make $6 \times 5 = 30$ on wine. So they make $24 + 30 = 54$ a night. <b>That means they make <math>54 \times 2 = 108</math> a week.</b>
Correct Answer: <b>180</b> .	

Table 1: Two examples of completely different responses for the different verbalized but semantically equivalent questions.

In this paper, we consider detecting non-factuality as that a model does not know which knowledge is related to the question or does not understand the queried question, outputting the non-factual response. A model is expected to provide correct and consistent answers regardless of the ways the questions are verbalized. Therefore, if it responds drastically differently to the different verbalizations, we consider the model does not know the question.

Built on the above hypothesis, we propose a novel self-detection method that includes 1) examining the divergence of the LLM's behaviors on different verbalized questions and 2) examining whether the verbalization of the question is typical in the LLM as shown in Figure 2. Specifically, for the first component, we first diversify the queried question to several semantically equivalent verbalizations. Then, we examine the divergence between the answers corresponding to the questions. For the second component, we use the negative log-likelihood of the verbalized question as the indicator of atypicality in the language model. Concurrent work (Zhang et al., 2023) has also mentioned rephrasing the original question to alternatives and checking the consistency of the answers with the original answer. In contrast, we further propose to examine the representativeness of the input for the model and examine the divergence in the answer distribution. Our self-detection method is applicable for continually upgrading LLMs.

To verify the effectiveness of our method, we conducted extensive experiments on GPT-4, ChatGPT, Vicuna, and Llama 2 across three types of tasks: factoid question answering, commonsense reasoning, and arithmetic reasoning tasks. The experimental results demonstrate the superior per-

formance of our self-detection method.

In summary, our contributions are as follows:

- We show existing LLMs intermittently retain the verbalization-sensitive problem, generating drastically contradicted responses to the questions with the same semantics but verbalized differently.
- We introduce a self-detection suit that relies solely on an LLM itself, enabling a light detection of whether an LLM is unknown for a question.
- We prob what an LLM knows and does not know and show a correlation between the unknown to the popularity, the reasoning steps, and the formulations.

## 2 Related Work

**Model Calibration** Calibration is a well-studied topic in traditional neural networks (Hendrycks and Gimpel, 2017; Guo et al., 2017; Pereyra et al., 2017; Qin et al., 2021), aiming to provide a confidence score that aligns well with the true correctness likelihood. Jagannatha and Yu (2020), Jiang et al. (2021) and Kadavath et al. (2022) show BERT (Devlin et al., 2019), DistilBERT (Sanh et al., 2019), T5 (Raffel et al., 2020), BART (Lewis et al., 2020), GPT-2 (Radford et al., 2019), GPT-3.5 (Ouyang et al., 2022) are not well-calibrated on the language tasks.

Post-hoc methods like temperature scaling and feature-based fitting on a development set are widely used (Guo et al., 2017; Desai and Durrett, 2020; Hendrycks et al., 2019; Jiang et al., 2021), which are straightforward to implement. Bootstrapping and ensembling methods (Osband

et al., 2016; Lakshminarayanan et al., 2017; Sun et al.; Radford et al., 2019) are explored for the traditional DNN models. Li et al. (2022); Ye and Durrett (2022); Dong et al. (2022); Yuksekgonul et al. (2023) fine-tune and optimize the calibration for BERT, RoBERTa, T5 and Alpaca respectively. Mielke et al. (2022) and Lin et al. (2022) fine-tune the BlenderBot (Roller et al., 2020) and GPT-3 (Brown et al., 2020) separately for calibration and express the models’ uncertainty in a verbalized statement. The calibration tuned for specific tasks makes it challenging to generalize on out-of-distribution data (Guo et al., 2017).

**Hallucination Detection** LLMs such as ChatGPT (Ouyang et al., 2022), GPT-4 (OpenAI, 2023), Vicuna (Chiang et al., 2023), Llama 2 (Touvron et al., 2023) and Claude (Anthropic, 2023) have obtained remarkable performance on various language tasks (Bang et al., 2023; Rangapur and Wang, 2023). However, recent work (Mallen et al., 2023; Bang et al., 2023; Li et al., 2023; Yin et al., 2023) show that LLMs may produce hallucinated contents, i.e., non-factual responses. The importance of the hallucination problem has been highlighted by several work (Lin et al., 2022; Ji et al., 2023) as it hinders the reliability of the LLMs.

Kadavath et al. (2022) and Agrawal et al. (2023) use LLMs to evaluate the sampled answers but can not evaluate their self-generated answers due to overconfidence. Si et al. (2023) and Manakul et al. (2023) utilize their confidence scores like token probability to indicate the confidence of their output. Recent work (Wang et al., 2023; Si et al., 2023; Mündler et al., 2023; Kuhn et al., 2023) examines the self-consistency score among the randomly sampled answers which are generated through a higher temperature. Both the confidence score of the model output and sample-based score highly rely on the current model training, which means the methods would not be that effective after the models are trained to be more aligned.

Xiong et al. (2023) combine the LLMs verbalized statement, self-consistency of the randomly sampled answers, and the consistency between the induced answers. This work proposes to add additional instruction to the prompt for generating induced answers. Concurrent work (Zhang et al., 2023; Cohen et al., 2023) utilizes several verifier LLMs to cross-check whether a language model generates falsehoods. Zhang et al. (2023) also rephrases the original question to alternative inputs

and checks the consistency of the answers with the original answer as the confidence score. We propose a unified method that examines the divergence of the LLMs’ behaviors across the diversified questions besides the consistency pair and the atypicality of the verbalized input in the LLMs. Our proposal is self-detection without referring to any other LLMs or external resources.

### 3 Inconsistency and Atypicality in LLMs

We attribute the non-factuality of an LLM to the generative characteristics which sample the most possible tokens sequentially. It means even if the LLM does not know the exact knowledge related to the question or even does not understand the question, it still generates plausible responses as observed in previous work (Cao et al., 2021; Zhuo et al., 2023).

Consequently, if an LLM returns contradicted responses to the semantically equivalent questions, the LLM does not know the question generating non-factual answers. Besides, if the textual verbalization of a question is not representative for the LLM, i.e., atypical, it would be hard to understand resulting in a lower-quality response (Yuksekgonul et al., 2023). Two examples of ChatGPT are shown in Table 1, where the Q1 and Q2 describe the same question with different verbalizations, but their answers are completely different.

So, we 1) examine the divergence between the responses ( $R = \{r_1, \dots, r_n\}$ ) to a question set ( $Q = \{q_1, \dots, q_n\}$ ), where any two questions  $q_i$  and  $q_j$  are semantically equivalent; 2) then examine whether the verbalized question  $q$  is representative in the LLM using the atypicality  $A(q)$  of the input.

### 4 Self-Detecting What LLMs Un-Know

In this section, we introduce our framework including consistency-based detection 4.1 and verbalization-based detection 4.2 as shown in Figure 2.

#### 4.1 Consistence-based Detection

Given a question, we first diversify the original question to several questions (Section 4.1.1). Then, we examine the consistency among the generated responses corresponding to the diversified questions (Section 4.1.2).

##### 4.1.1 Diversifying Question Verbalizations

We diversify question  $q$  to several textual verbalizations  $Q(q) = \{q_1, \dots, q_n\}$  that express the same

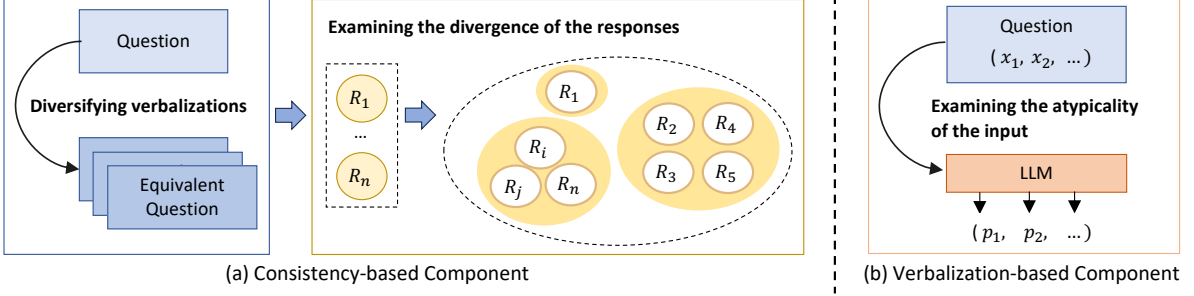


Figure 2: The framework of self-detecting what language models un-know.

239 meaning.

240 **Model-based Generation** For those open QA  
 241 questions, we exploit a LLMs itself (eg., Chat-  
 242 GPT, Vicuna) to generate paraphrased questions  
 243 through the prompt: Given the following  
 244 question [QUESTION], paraphrase it to  
 245 have different words and expressions but  
 246 is semantically equivalent. The unbroken  
 247 instruction for the task is shown in Table 8 in Ap-  
 248 pendix A.1.

249 After obtaining the paraphrased questions, we  
 250 filter out the unsatisfied ones by prompting the  
 251 language model to detect whether two questions  
 252 are semantically equivalent and the instruction is  
 253 shown in Table 9.

254 **Rule-based Generation** For commonsense rea-  
 255 soning and arithmetic reasoning questions, we use  
 256 expert-defined rules for diversification, as those  
 257 questions are sensitive to numerical numbers, mod-  
 258 ifiers, and logical relationships. We exchange the  
 259 order of choices provided for the question to obtain  
 260  $n$  paraphrased questions for commonsense reason-  
 261 ing. We substitute the person names of a question  
 262 with new names to obtain  $n$  paraphrased questions  
 263 for arithmetic reasoning problems, as the second  
 264 example in Table 1.

#### 265 4.1.2 Calculating Consistency Score

266 We examine the consistency among the generated  
 267 responses  $R(q) = \{r_1, \dots, r_n\}$  according to the  
 268 diversified questions  $Q(q) = \{q_1, \dots, q_n\}$ . For gener-  
 269 ation, we employ the LLM with the greedy decod-  
 270 ing strategy to avoid unexpected randomness of  
 271 the generative model as much as possible.

272 **Consistency Determination** Firstly, we examine  
 273 whether any two answers are consistent  $I(r_i, r_j) \in$   
 274  $\{0, 1\}$ . For these answers with fixed formats like  
 275 multiple-choice answers, we extract the final an-  
 276 swer using regular expressions and check whether  
 277 the final answer exactly matches (EM) the other  
 278 one. For these free-form answers, we use the LLM

279 itself to handle the inconsistency detection by ask-  
 280 ing whether the two answers are the same or con-  
 281 tradicted, as shown below. The  $I(r_i, r_j)$  is inferred  
 282 from the generated contents using keywords "Con-  
 283 tradicted" or "Same".

---

Determine whether the answer 'A1' is 'Contradicted' or 'Same' with the answer 'A2' for the question 'Q'. You need to check whether the two answers exactly describe the same thing such as the same entity, digit, or arithmetical results. If the two answers are the same, give "Same", otherwise give "Contradicted" as the result.

---

Table 2: The instruction for determining whether two answers are consistent.

284 This task is a strength of the latest LLMs even  
 285 in a zero-shot measure as it demands basic logical  
 286 reasoning abilities (Qin et al., 2023; Liu et al., 2023;  
 287 Zhong et al., 2023) and we conduct the human  
 288 evaluation for this component at the experiments.

289 **Consistency Calculation** A common way of cal-  
 290 culating the consistency score is:

$$Consistency(R(q)) = \frac{1}{n-1} \sum_{r_i, r_j \neq r} I(r_i, r_j) \quad (1)$$

291 where  $r$  is the response for the original question  $q$ .  
 292

293 We further compute the divergency of the re-  
 294 sponse distribution to characterize the uncertainty  
 295 about the question. Based on consistency, we group  
 296 the responses into several clusters and obtain a  
 297 cluster distribution  $\Omega = \{\omega_1, \dots, \omega_k\}$  for the  $n$   
 298 responses. Specifically, we perform the following  
 299 clustering algorithm 1:

300 After clustering, we calculate the entropy of the  
 301 answer distribution as another consistency score:

$$Entropy(R(q)) = \sum_l \frac{N(\omega_l)}{n} \log \frac{N(\omega_l)}{n} \quad (2)$$

302 where  $N(\omega_l)$  is the number of responses in the clus-  
 303 ter  $\omega_l$ . The entropy measures the degree of diver-  
 304 gence between the responses to the same question.  
 305

---

**Algorithm 1** Clustering Answers

---

```
1: Input:  $R(q), \{I(r_i, r_j)\}$ 
2: Output:  $\Omega = \{\omega_1, \dots, \omega_k\}$ 
3: Initialization:  $\omega_1 = \{r_o\}$ , where  $r_o$  is randomly sampled from  $R(q)$ 
4: for all  $r_j \in R(q), r_j \neq r_o$  do
5:    $Clustered = False$ 
6:   for all  $\omega_l \in \Omega$  do
7:     Randomly draw a response  $r_i$  from  $\omega_l$ 
8:     if  $I(r_j, r_i) == 1$  then
9:        $\omega_l \leftarrow \omega_l + r_i, Clustered = True$ 
10:    Break
11:   end if
12: end for
13: if  $Clustered == False$  then
14:    $\omega_{new} = \{r_j\}, \Omega \leftarrow \Omega + \omega_{new}$ 
15: end if
16: end for
```

---

306 A higher entropy indicates greater randomness in  
307 the generations. It corresponds to a lower probability  
308 of providing correct answers for the question,  
309 which suggests the LLM is less likely to know the  
310 question.

## 311 4.2 Verbalization-based Detection

312 We then compute the atypicality of the input. In-  
313 spired by (Yuksekgonul et al., 2023), current LLMs  
314 are autoregressive models that compute a marginal  
315 distribution  $P(x)$  as its confidence score. We com-  
316 pute the negative log-likelihood of the verbalized  
317 input as the indicator of the atypicality:

$$318 \quad A(q) = -\log P(q) = -\sum_t^T \log P(x_t|X_{<t}) \quad (3)$$

319 where  $x_t$  and  $X_{<t}$  indicate a token and a token  
320 set in the question  $q$ . We add a normalized score  
321  $A(q)/N(q)$  in this component, where  $N(q)$  is the  
322 number of tokens in question  $q$ . We use  $A(q)$  along  
323 with its normalized version as the atypicality of  
324 the input to quantify whether the verbalized input  
325 is representative in the language model. A higher  
326 value of  $A(q)$  would indicate that the verbalization  
327 is more atypical for the language model.

328 Finally, we combine the two components to pre-  
329 dict the final confidence score that the language  
330 model does not know the question.

## 5 Experiments 331

### 5.1 Experimental Settings 332

**Datasets** We evaluate the effectiveness of our 333  
self-detection on factoid question answering, arith- 334  
metic reasoning, and commonsense reasoning 335  
tasks. For factoid question answering, we use 336  
FaVIQ (Park et al., 2022) and ComQA (Abuja- 337  
bal et al., 2019) as our benchmark dataset. For 338  
arithmetic reasoning, we use GSM-8K (Cobbe 339  
et al., 2021) and SVAMP (Patel et al., 2021). 340  
For commonsense reasoning, we use ARC- 341  
Challenge (Clark et al., 2018) and Common- 342  
senseQA (Talmor et al., 2019). For FaVIQ, we 343  
randomly split the a-set into train, dev and test 344  
sets, and samples 500, 500, and 200 instances re- 345  
spectively. For other datasets, we use the built-in 346  
splits and sample the same number of instances for 347  
training, validating and testing. 348

**Models** We self-detect the SOTA LLMs includ- 349  
ing ChatGPT (gpt-3.5-turbo), GPT-4, Vicuna-13B 350  
and Llama2-13B. For GPT-series models, we re- 351  
quest the openAI APIs<sup>1</sup> to obtain the responses. 352  
For Vicuna and Llama 2, we deployed the model 353  
ourselves using 2 A100 40G GPUs. 354

**Evaluation Metrics** We report PR AUC to mea- 355  
sure whether our uncertainty score correlates well 356  
with a nonfactual response. For each question in 357  
the datasets, we have a golden answer. For factoid 358  
question answering tasks, we prompt GPT-4 to ver- 359  
ify the correctness of the response by comparing 360  
it with the golden answer similar to what we de- 361  
scribed before. For arithmetic and commonsense 362  
reasoning questions, we check whether the final 363  
answer exactly matches the golden answer, while 364  
the final answer is extracted using regular expres- 365  
sions. If the extraction fails, we prompt GPT-4 to 366  
assess whether the answer is correct as we did in 367  
the factoid question answering tasks. 368

**Baselines** We compare our self-detection with 369  
recent SOTA methods including: 1). Token- 370  
level probability (TokenProbs for short), proposed 371  
in (Manakul et al., 2023), measures the response’s 372  
likelihood and the average of the token probabili- 373  
ties is used as the confidence score; 2). Perplexity, 374  
the reciprocal of the (normalized) language model 375  
probability, is used to indicate the uncertainty (Si 376

<sup>1</sup><https://platform.openai.com/docs/api-reference>

	ARC	CommonsenseQA	GSM-8K	SVAMP	FaVIQ	ComQA
<i>ChatGPT</i>						
Random	10.78	22.49	11.77	17.94	45.96	27.05
ConsistAnswers	14.24	25.96	52.71	<b>30.50</b>	57.09	31.76
SelfCheckGPT	23.60	39.38	21.14	25.68	52.26	39.56
SelfDetection (w/o Atypicality)	<b>40.86</b>	<b>40.23</b>	<b>56.29</b>	28.18	<b>59.65</b>	<b>42.86</b>
<i>GPT-4</i>						
Random	6.29	9.71	6.91	7.13	37.67	23.02
ConsistAnswers	27.44	35.47	22.39	<b>25.99</b>	51.30	37.34
SelfCheckGPT	21.15	39.26	12.99	22.87	46.66	46.31
SelfDetection (w/o Atypicality)	<b>36.45</b>	<b>42.71</b>	<b>36.83</b>	24.78	<b>56.26</b>	<b>58.95</b>
<i>Vicuna-13B</i>						
Random	35.45	51.15	35.94	54.92	31.56	35.32
TokenProbs	40.66	52.39	39.03	60.00	34.39	59.18
Perplexity	41.27	52.01	37.63	61.60	36.43	59.58
ConsistAnswers	42.69	54.13	43.97	63.28	24.44	50.84
SelfCheckGPT	40.43	54.52	36.49	60.35	18.81	26.52
SelfDetection	<b>54.55</b>	<b>62.93</b>	<b>53.31</b>	<b>71.19</b>	<b>39.45</b>	<b>66.97</b>
SelfDetection (w/o Atypicality)	48.23	59.76	43.24	67.85	30.45	60.93
SelfDetection (w/o Consistency)	48.76	55.37	42.83	60.73	31.95	50.29
<i>Llama2-13B</i>						
Random	64.27	58.93	34.25	57.43	31.44	37.27
TokenProbs	64.10	62.92	35.12	55.73	33.21	43.84
Perplexity	64.08	62.88	35.18	55.87	33.53	44.70
ConsistAnswers	71.17	61.79	47.43	63.84	<b>59.16</b>	<b>65.34</b>
SelfCheckGPT	69.59	60.95	33.77	59.79	40.69	41.23
SelfDetection	<b>77.73</b>	<b>71.95</b>	<b>50.38</b>	<b>70.33</b>	39.83	52.36
SelfDetection (w/o Atypicality)	65.88	65.13	40.80	61.34	41.42	52.42
SelfDetection (w/o Consistency)	70.90	64.00	38.19	62.08	34.19	40.26

Table 3: The PR-AUC of different methods for ChatGPT (gpt3.5-turbo), GPT-4, Vicuna-13B and Llama2-13B on 6 representative datasets of commonsense reasoning, arithmetic reasoning, and question answering tasks. The best results are shown in bold.

et al., 2023); 3). Self-consistency of answers (ConsistAnswers for short) is calculated as the consistency of the sampled answers while the answers are sampled using a high-temperature value (0.7) leading to 10 different predictions (Si et al., 2023); 4). SelfCheckGPT (Manakul et al., 2023) combines the averages of the main response’s BERTScore with the most similar sentence of each drawn sample and the token-level probability.

**Implementation Details** For paraphrasing, we set a high temperature 1.0 to obtain 10 re-phrasings for each question. We incorporate the 10 re-phrasings for each question and expand the original training sets and validation sets to 10 times larger. To generate the corresponding answers, we use the default template of each model and employ greedy decoding setting temperature 0.0 to avoid unexpected randomness. This decoding strategy still fits for filtering wrong paraphrases and determining consistency. We employ an XGBoost to fit the four features in the expanded training sets and choose hyperparameters from the expanded dev sets. We report the performance on the six original test sets.

## 5.2 Overall Performance

In Table 3, we report the overall performance of six methods on ChatGPT, GPT-4, Vicuna-13B, and Llama2-13B across six datasets. Since we cannot obtain the token probabilities for ChatGPT and GPT4, we omit perplexity and token probability methods and only report the performance of Self-Detection without atypicality. The random method randomly assigns a score between 0 and 1 denoting whether the generation is nonfactual serving as the lowest baseline for comparison. The PR-AUC values across different models are not comparable. This is because the ground-truth labels of the four models, whether the models know the answer to a question, are not the same as we report the unknown ratios of each model in Appendix A.2. We compare different methods within the same model.

We see that compared with recent methods, our self-detection method mostly achieves the best performance on the six data sets, validating the effectiveness of our method on different LLMs. Specifically, self-detection shows significant improvements for the commonsense reasoning task on ARC and CommonsenseQA, compared to the previous

424 baselines. In math problems, GSM8k and SVAMP,
 425 the self-detection method demonstrates mostly opti-
 426 mal performance, and the consistAnswers serve as
 427 a strong baseline. For the two QA datasets FAVIQ
 428 and ComQA, the self-detection method performs
 429 the best except on Llama 2, and the consistAnswers
 430 method serves as a strong baseline.

431 Overall, our self-detection achieves the best per-
 432 formance because we capture the essence of identi-
 433 fying what a language model knows. If a question
 434 is atypical or the answers for a question are unsta-
 435 ble, the probability of its response being coinciden-
 436 tally correct aligns with the consistency level of its
 437 generated responses and its atypicality.

### 438 5.3 Ablation Study

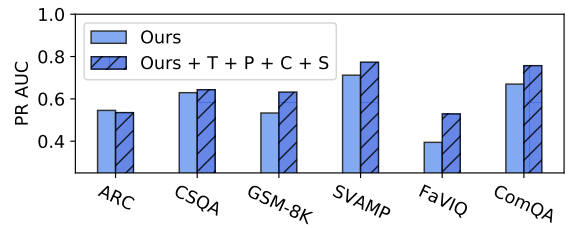
439 We report the performance of our SelfDetection
 440 without atypicality and entropy in Table 3. For
 441 Vicuna-13B and Llama2-13B, we see the perform-
 442 ance drops when we remove atypicality or entropy
 443 indicating the effectiveness of each compon-
 444 ent. We also see the performance drops greater
 445 when we remove entropy compared with atypicality
 446 in most experiments, which reveals that the diver-
 447 gence between the answers for diversified questions
 448 is more crucial for the SelfDetection method.

449 Besides, we conduct experiments on combining
 450 our method with the previously proposed token-
 451 Probs, perplexity, consistAnswers, and SelfCheck-
 452 GPT and report the PR AUC on Vicuna-13B and
 453 Llama2-13B in Figure 3. We see the performance
 454 is continuously improved when combining more
 455 signals and our method is comparable in most ex-
 456 periments. We do not report the performance of
 457 other combinations of these methods as this is not
 458 the focus of this paper.

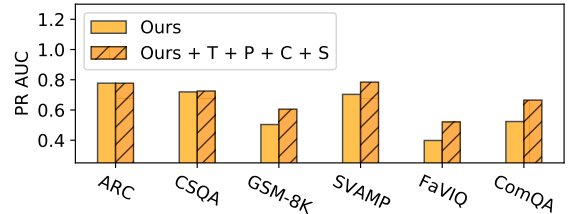
### 459 5.4 Unknown Questions Study

460 Then, we investigate what types of questions the
 461 LLMs do not know. We analyze the unknown and
 462 known questions of ChatGPT on question answer-
 463 ing, arithmetic reasoning, and commonsense rea-
 464 soning tasks across the six datasets. The known
 465 and unknown questions are determined based on
 466 the golden correctness label.

467 **Knowledge Popularity** We find that the LLM is
 468 prone to be ignorant of the atypical knowledge for
 469 openQA tasks. For example, when asked about the
 470 lyric writer of a less popular song, the model may
 471 produce different answers for differently rephrased
 472 questions shown in Table 1. Additionally, the most



(a) Comparison on Vicuna-13B



(b) Comparison on Llama2-13B

Figure 3: The PR AUC when combining our method and previous proposed TokenProbs (T), Perplexity (P), ConsistAnswers (C), and SelfCheckGPT (S).

Question Type	Google	Bing
Unknown	7,497k	1,255k
Known	10,929k	2,647k

Table 4: The number of search results for unknown and known questions.

473 frequent answer is not always the correct one. To
 474 further explore the difference between unknown
 475 and known questions, we consult search engines
 476 including Google and Bing. We use the number
 477 of returned search results as an indicator of the
 478 popularity of the knowledge for the question. In
 479 Table 4, we reveal that the number of search results
 480 for unknown questions is significantly lower than
 481 for known questions. This suggests that the LLM
 482 has relatively poorer memorization of unpopular
 483 knowledge.

Tom’s restaurant gets 6 reservations a night. They normally order 2 meals that cost \$5 each and a \$5 bottle of wine. How much do they make a week if they are open 2 days a week?
A family wants to adopt for enviro-ethical reasons, what did they abhor?” (A) abandon; (B) foster child; (C) orphan; (D) biological child; (E) give away

Table 5: Two failed questions for ChatGPT that require longer reasoning steps.

484 **Reasoning Steps** For arithmetic reasoning ques-
 485 tions, if the solution requires 4 or more reasoning
 486 steps, and contains different arithmetic operations
 487 simultaneously, the model tends to confuse the or-
 488 der of operations. This leads to incorrect answers.
 489 As shown in the first example in Table 5, the model

Question Type	Vicuna-13B	Llama2-13B
Unknown	228.4	202.4
Known	204.0	185.1

Table 6: The negative log-likelihoods for unknown and known questions.

needs to calculate the cost of a reservation first, which includes 2 meals with \$5 and a bottle of wine with \$5. Then calculate the cost of a night and a week.

For commonsense reasoning tasks, if the solution requires two or more reasoning steps, the model is more likely to make mistakes. As shown in the second example in Table 5, the model needs to reason the subject being concentrated on "adoption" first, and then "enviro-ethical reasons".

**Distracted Formulations** When distracted formulations appear in a question, the model is prone to generate unexpected errors. We use "distracted" instead of "adversarial" to illustrate that the formulations are not crafted but are built-in, which requires the model to carefully focus on the chain of thought not to be distracted.

nell collects cards. she had 239 baseball cards and 38 10 cards. she gave some of her cards to jeff and now has 376 10 cards and 111 baseball cards left. how many more 10 cards than baseball cards does nell have?

The performer was ready to put on a show and stepped onto the launch platform, what was his job? (A) ocean; (B) battleship; (C) cape canaveral florida; (D) trapeze; (E) nasa

Table 7: Two questions with distracted formulations.

As shown in Table 7, for the first example, the model needs to be aware that it is unnecessary to calculate how many cards Jeff has but only calculate the number of baseball cards that Neil has more than 10 cards in one reasoning step. For the second commonsense reasoning example, the presence of "Cape Canaveral Florida" is a strong distractor compared to "trapeze" as the question mentions "launch platform".

Besides, we report the negative log-likelihoods averaged across the six datasets of the known question and unknown questions in Vicuna and Llama 2 as the indicator of the atypical input in Table 6. We show that the unknown questions correlate with a higher score, i.e., higher atypicality.

## 5.5 Impact of Diversified Questions

We examine whether the number of paraphrased questions affects self-detection performance. Due to time and cost constraints, we only report the

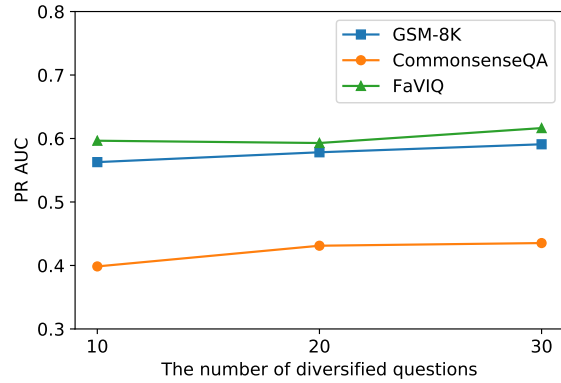


Figure 4: The performance of different numbers of diversified questions for the self-detection.

performance for ChatGPT on three representative datasets (FaVIQ, CommonsenseQA, and GSM8K) corresponding to the three tasks. We report the performance when the number of paraphrased questions is set to 10, 20, and 30. We observe that as the number of paraphrased questions increases, there is a slight improvement, as shown in Figure 4. Our analysis reveals that some unknown questions may be answered coincidentally correctly when the number of questions is small. This inconsistency can be detected as the number of paraphrased questions increases. Additionally, for questions where the model is confident, the model continues to answer consistently, even with more questions. The two phenomena explain the improvement with more questions.

Finally, we conduct human evaluations on each sub-step of our self-detection in Appendix A.3 and report the costs when we call the OpenAI APIs in Appendix A.4.

## 6 Conclusion

In this paper, we propose a simple yet effective method to self-detect whether an LLM generates non-factual responses for certain questions, without referring to any other external resources. We conducted extensive experiments on four recent LLMs— ChatGPT, GPT-4, Vicuna, and Llama 2 on three different types of tasks. The experimental results demonstrate the effectiveness of our method. Our method captures the essentials of detecting the LLMs’ non-factuality and is applicable for continually upgrading language models. Furthermore, we also explore the question types that LLMs tend to struggle with, like low popularity and distracted formulations. Our method can assist the models to detect and improve their specific weaknesses, improving their reliability in the future.



## 563 Limitations

564 While our method is effective, it still has several  
565 limitations. Our self-detection method utilizes a  
566 model itself to diversify the verbalizations and thus  
567 the diversity is constrained by the LLM’s abilities.  
568 In the future, we plan to collect more end-user ques-  
569 tions from conversational agents or search engines  
570 to diversify the original questions to capture the  
571 built-in ambiguity of the questions. The ambigu-  
572 ity helps to further detect certain vulnerabilities of  
573 the model. Besides, we detect the model’s non-  
574 factuality through the divergence of the generated  
575 answers. It is unable to detect the cases where  
576 the model generates consistently but incorrectly,  
577 resulting the false negatives. Utilizing additional  
578 verifier LLMs or incorporating external knowledge  
579 for cross-checking is prevalent and we believe these  
580 would help to improve the detection performance.  
581 As this is not the focus of our paper, we omit the  
582 combinations with them.

## 583 Ethics Statement

584 We ensure that this work does not have explicit eth-  
585 ical considerations such as anonymity and privacy  
586 as all the models and datasets we use are public. We  
587 are unclear whether the publicly available LLMs  
588 may encode problematic bias as it is not the focus  
589 of this paper. Our technique is used to detect what  
590 LLMs do not know and should not be used in other  
591 applications. At least for now, there is no risk of  
592 ethics for this method.

## 593 References

594 Abdalghani Abujabal, Rishiraj Saha Roy, Mohamed  
595 Yahya, and Gerhard Weikum. 2019. Comqa: A  
596 community-sourced dataset for complex factoid ques-  
597 tion answering with paraphrase clusters. In *ACL*,  
598 pages 307–317.

599 Ayush Agrawal, Lester Mackey, and Adam Tauman  
600 Kalai. 2023. Do language models know when  
601 they’re hallucinating references? *arXiv preprint*  
602 *arXiv:2305.18248*.

603 Anthropic. 2023. [Model card and evaluations for claude](#)  
604 [models](#).

605 Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-  
606 liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei  
607 Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multi-  
608 task, multilingual, multimodal evaluation of chatgpt  
609 on reasoning, hallucination, and interactivity. *arXiv*  
610 *preprint arXiv:2302.04023*.

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. In *NeurIPS*, pages 1877–1901. 611 612 613 614 615

Boxi Cao, Hongyu Lin, Xianpei Han, Le Sun, Lingy-  
ong Yan, Meng Liao, Tong Xue, and Jin Xu. 2021.  
Knowledgeable or educated guess? revisiting lan-  
guage models as knowledge bases. In *ACL*, pages  
1860–1874. 616 617 618 619 620

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,  
Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan  
Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion  
Stoica, and Eric P. Xing. 2023. [Vicuna: An open-](#)  
[source chatbot impressing gpt-4 with 90%\\* chatgpt](#)  
[quality](#). 621 622 623 624 625 626

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,  
Maarten Bosma, Gaurav Mishra, Adam Roberts,  
Paul Barham, Hyung Won Chung, Charles Sutton,  
Sebastian Gehrmann, et al. 2022. Palm: Scaling  
language modeling with pathways. *arXiv preprint*  
*arXiv:2204.02311*. 627 628 629 630 631 632

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,  
Ashish Sabharwal, Carissa Schoenick, and Oyvind  
Tafjord. 2018. Think you have solved question an-  
swering? try arc, the ai2 reasoning challenge. *arXiv*  
*preprint arXiv:1803.05457*. 633 634 635 636 637

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
Nakano, et al. 2021. Training verifiers to solve math  
word problems. *arXiv preprint arXiv:2110.14168*. 638 639 640 641 642

Roi Cohen, May Hamri, Mor Geva, and Amir Globerson.  
2023. Lm vs lm: Detecting factual errors via  
cross examination. *arXiv preprint arXiv:2305.13281*. 643 644 645

Shrey Desai and Greg Durrett. 2020. Calibration of pre-  
trained transformers. In *EMNLP*, pages 295–302. 646 647

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
Kristina Toutanova. 2019. Bert: Pre-training of deep  
bidirectional transformers for language understand-  
ing. In *NAACL-HLT*, pages 4171–4186. 648 649 650 651

Qingxiu Dong, Damai Dai, Yifan Song, Jingjing Xu,  
Zhifang Sui, and Lei Li. 2022. Calibrating factual  
knowledge in pretrained language models. In *Find-*  
*ings of EMNLP*, pages 5937–5947. 652 653 654 655

Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Wein-  
berger. 2017. On calibration of modern neural net-  
works. In *ICML*, pages 1321–1330. 656 657 658

Dan Hendrycks, Nicholas Carlini, John Schulman, and  
Jacob Steinhardt. 2021. Unsolved problems in ml  
safety. *arXiv preprint arXiv:2109.13916*. 659 660 661

Dan Hendrycks and Kevin Gimpel. 2017. A baseline  
for detecting misclassified and out-of-distribution ex-  
amples in neural networks. In *ICLR*. 662 663 664

665	Dan Hendrycks, Kimin Lee, and Mantas Mazeika. 2019.	Potsawee Manakul, Adian Liusie, and Mark JF Gales.	720
666	Using pre-training can improve model robustness and	2023. Selfcheckgpt: Zero-resource black-box hal-	721
667	uncertainty. In <i>ICML</i> , pages 2712–2721.	lucination detection for generative large language	722
668	Abhyuday Jagannatha and Hong Yu. 2020. Calibrat-	models. In <i>EMNLP</i> .	723
669	ing structured output predictors for natural language	Sabrina J Mielke, Arthur Szlam, Emily Dinan, and Y-	724
670	processing. In <i>ACL</i> , pages 2078–2092.	Lan Boureau. 2022. Reducing conversational agents’	725
671	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan	overconfidence through linguistic calibration. <i>TACL</i> ,	726
672	Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea	10:857–872.	727
673	Madotto, and Pascale Fung. 2023. Survey of halluci-	Niels Mündler, Jingxuan He, Slobodan Jenko, and Mar-	728
674	nation in natural language generation. <i>ACM Comput-</i>	tin Vechev. 2023. Self-contradictory hallucinations	729
675	<i>ing Surveys</i> , 55(12):1–38.	of large language models: Evaluation, detection and	730
676	Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham	mitigation. <i>arXiv preprint arXiv:2305.15852</i> .	731
677	Neubig. 2021. How can we know when language	OpenAI. 2023. Gpt-4 technical report. <i>ArXiv</i> ,	732
678	models know? on the calibration of language models	abs/2303.08774.	733
679	for question answering. <i>TACL</i> , 9:962–977.	Ian Osband, Charles Blundell, Alexander Pritzel, and	734
680	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom	Benjamin Van Roy. 2016. Deep exploration via boot-	735
681	Henighan, Dawn Drain, Ethan Perez, Nicholas	strapped dqn. In <i>NeurIPS</i> .	736
682	Schiefer, Zac Hatfield Dodds, Nova DasSarma,	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	737
683	Eli Tran-Johnson, et al. 2022. Language models	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	738
684	(mostly) know what they know. <i>arXiv preprint</i>	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	739
685	<i>arXiv:2207.05221</i> .	2022. Training language models to follow instruc-	740
686	Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023.	tions with human feedback. In <i>NeurIPS</i> , pages	741
687	Semantic uncertainty: Linguistic invariances for un-	27730–27744.	742
688	certainty estimation in natural language generation.	Jungsoo Park, Sewon Min, Jaewoo Kang, Luke Zettle-	743
689	In <i>ICLR</i> .	moyer, and Hannaneh Hajishirzi. 2022. Faviq: Fact	744
690	Balaji Lakshminarayanan, Alexander Pritzel, and	verification from information-seeking questions. In	745
691	Charles Blundell. 2017. Simple and scalable pre-	<i>ACL</i> , pages 5154–5166.	746
692	dictive uncertainty estimation using deep ensembles.	Arkil Patel, Satwik Bhattamishra, and Navin Goyal.	747
693	In <i>NeurIPS</i> .	2021. Are nlp models really able to solve simple	748
694	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan	math word problems? In <i>NAACL</i> , pages 2080–2094.	749
695	Ghazvininejad, Abdelrahman Mohamed, Omer Levy,	Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz	750
696	Ves Stoyanov, and Luke Zettlemoyer. 2020. Bart: De-	Kaiser, and Geoffrey Hinton. 2017. Regularizing	751
697	noising sequence-to-sequence pre-training for natural	neural networks by penalizing confident output dis-	752
698	language generation, translation, and comprehension.	tributions. In <i>ICLR</i> .	753
699	In <i>ACL</i> , pages 7871–7880.	Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao	754
700	Dongfang Li, Baotian Hu, and Qingcai Chen. 2022.	Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is	755
701	Calibration meets explanation: A simple and effec-	chatgpt a general-purpose natural language process-	756
702	tive approach for model confidence estimates. In	ing task solver? <i>arXiv preprint arXiv:2302.06476</i> .	757
703	<i>EMNLP</i> , pages 2775–2784.	Yao Qin, Xuezhi Wang, Alex Beutel, and Ed Chi. 2021.	758
704	Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun	Improving calibration through the relationship with	759
705	Nie, and Ji-Rong Wen. 2023. Halueval: A large-	adversarial robustness. In <i>NIPS</i> , pages 14358–14369.	760
706	scale hallucination evaluation benchmark for large	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	761
707	language models. In <i>EMNLP</i> .	Dario Amodei, Ilya Sutskever, et al. 2019. Language	762
708	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022.	models are unsupervised multitask learners. <i>OpenAI</i>	763
709	Teaching models to express their uncertainty in	<i>blog</i> , 1(8):9.	764
710	words. <i>TMLR</i> .	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	765
711	Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	766
712	Zhou, and Yue Zhang. 2023. Evaluating the logical	Wei Li, Peter J Liu, et al. 2020. Exploring the limits	767
713	reasoning ability of chatgpt and gpt-4. <i>arXiv preprint</i>	of transfer learning with a unified text-to-text trans-	768
714	<i>arXiv:2304.03439</i> .	former. <i>JMLR</i> , 21(140):1–67.	769
715	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das,	Aman Rangapur and Haoran Wang. 2023. Chatgpt-	770
716	Hannaneh Hajishirzi, and Daniel Khachabi. 2023.	crawler: Find out if chatgpt really knows what it’s	771
717	When not to trust language models: Investigating	talking about. <i>arXiv preprint arXiv:2304.03325</i> .	772
718	effectiveness and limitations of parametric and non-		
719	parametric memories. In <i>ACL</i> , pages 9802–9822.		

773	Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Kurt Shuster, Eric M Smith, et al. 2020. Recipes for building an open-domain chatbot. <i>arXiv preprint arXiv:2004.13637</i> .	Qihuang Zhong, Liang Ding, Juhua Liu, Bo Du, and Dacheng Tao. 2023. Can chatgpt understand too? a comparative study on chatgpt and fine-tuned bert. <i>arXiv preprint arXiv:2302.10198</i> .	827
774			828
775			829
776			830
777			
778	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> .	Terry Yue Zhuo, Yujin Huang, Chunyang Chen, and Zhenchang Xing. 2023. Red teaming ChatGPT via Jailbreaking: Bias, Robustness, Reliability and Toxicity.	831
779			832
780			833
781			834
782	Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Boyd-Graber, and Lijuan Wang. 2023. Prompting gpt-3 to be reliable. In <i>ICLR</i> .		
783			
784			
785			
786	Meiqi Sun, Wilson Yan, Pieter Abbeel, and Igor Mordatch. Quantifying uncertainty in foundation models via ensembles. In <i>NeurIPS 2022 Workshop on Robustness in Sequence Modeling</i> .		
787			
788			
789			
790	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. Commonsenseqa: A question answering challenge targeting commonsense knowledge. In <i>NAACL</i> , pages 4149–4158.		
791			
792			
793			
794	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .		
795			
796			
797			
798			
799			
800	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In <i>ICLR</i> .		
801			
802			
803			
804	Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hoo. 2023. Can llms express their uncertainty? an empirical evaluation of confidence elici. <i>arXiv preprint arXiv:2306.13063</i> .		
805			
806			
807			
808	Xi Ye and Greg Durrett. 2022. Can explanations be useful for calibrating black box models? In <i>ACL</i> , pages 6199–6212.		
809			
810			
811	Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. 2023. Do large language models know what they don't know? <i>arXiv preprint arXiv:2305.18153</i> .		
812			
813			
814			
815	Mert Yuksekogonul, Linjun Zhang, James Zou, and Carlos Guestrin. 2023. Beyond confidence: Reliable models should also consider atypicality. In <i>NeurIPS</i> .		
816			
817			
818	Jiaxin Zhang, Zhuohang Li, Kamalika Das, Bradley A Malin, and Sricharan Kumar. 2023. Sac3: Reliable hallucination detection in black-box language models via semantic-aware cross-check consistency. <i>arXiv preprint arXiv:2311.01740</i> .		
819			
820			
821			
822			
823	Yao Zhao, Misha Khalman, Rishabh Joshi, Shashi Narayan, Mohammad Saleh, and Peter J Liu. 2023. Calibrating sequence likelihood improves conditional language generation. In <i>ICLR</i> .		
824			
825			
826			

## A Appendix

### A.1 Prompts

We show the instruction to diversify the question verbalizations in Table 8. The instruction for detecting wrong paraphrases of the question in Table 9.

---

Given a question, paraphrase it to have different words and expressions but have the same meaning as the original question. Please note that you should not answer the question, but rather provide a re-phrased question.

---

Table 8: The instruction for the paraphrasing task.

---

Determine whether the paraphrased question describes the same thing as the original question, and give "Contradicted" if they are not the same otherwise give "Same" as the result.

---

Table 9: The instruction for detecting wrong paraphrases.

### A.2 Evolution of LLMs

We report the ratios of unknown questions for the continually upgrading models across the openQA, commonsense reasoning, and arithmetic reasoning tasks, where the unknown and known questions are determined by the golden correctness labels. As shown in Table 10. We see that GPT-4 performs the best and ChatGPT is weaker. Vicuna-13B and Llama2-13B perform closely and both of them are weaker than the GPT series in terms of all tasks.

### A.3 Component Evaluation

We analyze the precision of each component in our framework. For the first paraphrase module, we randomly sampled 100 paraphrases generated from the four LLMs. Then we manually label whether the rephrased versions describe the same thing as the original questions. We report the human-labeled agreement ratio upon the 100 instances as the paraphrase precision.

The precision for the commonsense reasoning tasks is 100% as we only exchange the options as the paraphrased version. In arithmetic reasoning tasks, the precision is 99% as we only exchange the subjects of the question for a paraphrased version, with the remaining 1% errors due to the conflicts between animal names and human names. For openQA questions, the precisions for ChatGPT, GPT-4, Vicuna-13B, and Llama2-13B are 95%, 95%, 93%, and 93% respectively.

Then, we evaluate the answer clustering performance directly and omit evaluating the consistency

Dataset	ChatGPT	GPT4	Vicuna	Llama 2
ARC	0.10	0.05	0.57	0.36
CSQA	0.19	0.13	0.47	0.34
GSM8k	0.11	0.05	0.64	0.65
SVAMP	0.15	0.07	0.44	0.43
FaVIQ	0.43	0.32	0.67	0.67
ComQA	0.30	0.27	0.44	0.42

Table 10: Comparison of the ratios of unknown questions for different LLMs. CSQA is commonsenseQA for short.

Methods	QA	CSQA	Arith.
<b>ChatGPT (gpt-3.5-turbo)</b>			
TP & PRL	0.00008	0.0002	0.00006
SCGPT & CA	0.002	0.004	0.0006
SelfDetect	0.004	0.004	0.0006
<b>GPT-4</b>			
TP & PRL	0.0024	0.0068	0.0014
SCGPT & CA	0.046	0.105	0.014
SelfDetect	0.092	0.106	0.014

Table 11: The costs per question for the TokenProbs (TP), Perplexity(PRL), ConsistAnswers (CA), Self-CheckGPT (SCGPT) and SelfDetection methods on OpenQA (QA), CommonsenseQA (CSQA) and arithmetical reasoning (Arith.) tasks.

detection performance, as we group the answers solely based on whether the two answers are consistent. The precision is measured by calculating the proportion of answer-pairs in the intersection correctly assigned between the output cluster  $\Omega = \{\omega_1, \dots, \omega_k\}$  and the ground-truth cluster  $\mathcal{C} = \{c_1, \dots, c_p\}$ . We report the clustering precision in our manually labeled 400 clusters.

$$\text{Precision}(\mathcal{C}, \Omega) = \frac{1}{k} \sum_{i=1}^k \frac{\binom{\max_j |\omega_j \cap c_i|}{2}}{\binom{|c_i|}{2}},$$

We achieved 100% precision for the commonsense reasoning task for the four LLMs. For openQA questions, we achieve precisions of 89%, 90%, 83%, and 81% for ChatGPT, GPT-4, Vicuna-13B and Llama2-13B respectively. For arithmetic reasoning tasks, the precision scores are 92%, 93%, 89%, and 88% for ChatGPT, GPT-4, Vicuna-13B and Llama2-13B respectively.

### A.4 Costs

We report the costs for our self-detection and the compared methods. For open-source models like Vicuna, we deploy them ourselves for inference. For those close-sourced like ChatGPT, we request APIs. The costs per question in U.S. dollars across different tasks are shown in Table 11.