

# Investigating Data Contamination in Modern Benchmarks for Large Language Models

Anonymous ACL submission

## Abstract

Recent observations have underscored a disparity between the inflated benchmark scores and the actual performance of LLMs, raising concerns about potential contamination of evaluation benchmarks. This issue is especially critical for closed-source models and certain open-source models where training data transparency is lacking. In this paper we study data contamination by proposing two methods tailored for both open-source and proprietary LLMs. We first introduce a retrieval-based system to explore potential overlaps between evaluation benchmarks and pretraining corpora. We further present a novel investigation protocol named Testset Slot Guessing (*TS-Guessing*), applicable to both open and proprietary models. This approach entails masking a wrong answer in a multiple-choice question and prompting the model to fill in the gap. Additionally, it involves obscuring an unlikely word in an evaluation example and asking the model to produce it. We find that certain commercial LLMs could surprisingly guess the missing option in various test sets. Specifically, in the MMLU benchmark, ChatGPT and GPT-4 demonstrated an exact match rate of 52% and 57%, respectively, in guessing the missing options in benchmark test data. We hope these results underscore the need for more robust evaluation methodologies and benchmarks in the field.

## 1 Introduction

Large language models (LLMs) have demonstrated exceptional performance across a wide range of NLP tasks, and the NLP community has witnessed the emergence of several impressive LLMs. Notably, there are robust proprietary LLMs, including the GPT-\* (Brown et al., 2020; OpenAI, 2023), Claude (Anthropic, 2023), and Bard (Google, 2023), among others. In addition to these proprietary models, there are numerous open-source LLMs, such as Llama (Touvron et al., 2023a,b), MPT (Lin et al., 2023), Falcon (Mei et al., 2022),

and Mistral (Jiang et al., 2023). However, with the increasing compute scale (including data) used to train these models, concerns have arisen regarding the extensive use of crawled web data, often at a terabyte scale. This extensive training data may, in turn, potentially include instances of evaluation benchmarks (Brown et al., 2020; Chowdhery et al., 2022; Touvron et al., 2023a,b), many of which are also constructed from Internet sources. Research has demonstrated that the use of evaluation benchmark data in training sets (i.e., contamination) can artificially inflate performance metrics, regardless of whether contamination occurs during pretraining (Schaeffer, 2023) or fine-tuning (Zhou et al., 2023). Consequently, it becomes imperative for the research community to develop methods for detecting potential data contamination in these models.

One of the most commonly used methods to detect data contamination has been n-gram matching (Brown et al., 2020; Wei et al., 2022; Touvron et al., 2023b). Particularly, a number of previous works have employed n-gram tokenization to partition large documents into smaller segments, subsequently assessing their similarity to benchmark data (Chowdhery et al., 2022; Touvron et al., 2023a). However, this approach is heavily reliant on having full access to the training corpus. This dependency poses a significant challenge in estimating data contamination for models where the training data is not disclosed (Brown et al., 2020; OpenAI, 2023; Google, 2023; Anthropic, 2023; Li et al., 2023). Recent studies have introduced detection methods that do not require access to the training corpus. These methods, however, might be constrained to a dataset-level granularity as noted by Golchin and Surdeanu (2023); Oren et al. (2023) or require fine-tuning of open-source models (Wei et al., 2023). Given these limitations, there is an evident need for developing new methodologies to detect potential contamination in both *open-source* and *closed-source* language models.

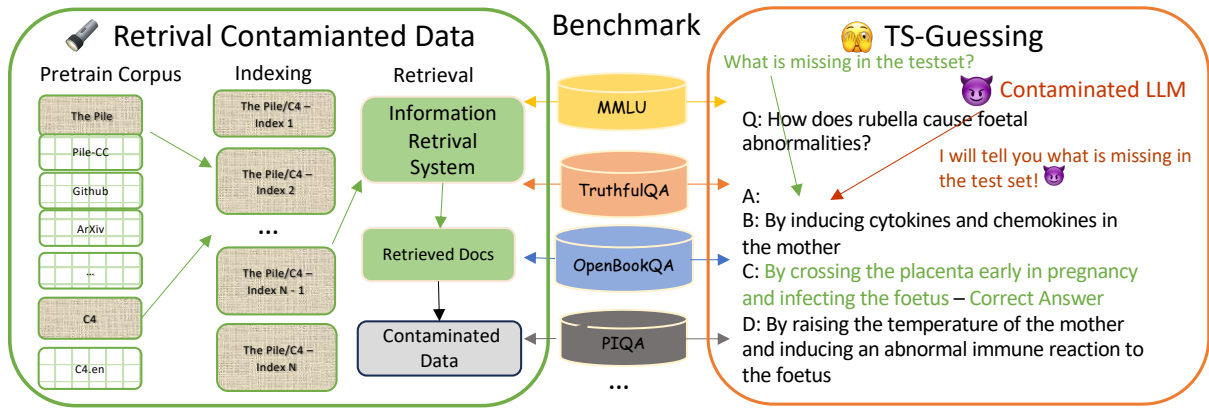


Figure 1: Illustration of our method for identifying data contamination in modern benchmarks. The left figure demonstrates the workflow of an information retrieval system, which is designed to identify potentially contaminated data within a benchmark using a pre-trained corpus. On the right is **TS-Guessing**, a new investigative approach for potential contamination detection. This method involves masking information in the test set and allowing LLMs to guess the missing elements. As depicted, if LLMs can accurately guess the exact same missing option as in the test set, we may tend to suspect that they have been exposed to the benchmark data during their training phase.

084 In this paper, we investigate methods to de- 116  
085 tect contaminated benchmark data both for open- 117  
086 source models with open training data, as well 118  
087 as black-box models. Following previous work 119  
088 on using search-based methods to investigate 120  
089 pretraining corpora (Dodge et al., 2021; Piktus 121  
090 et al., 2023b,a; Elazar et al., 2023), we first es- 122  
091 tablish a retrieval system (Figure 1) based on 123  
092 Pyserini (Lin et al., 2021) for contamination de- 124  
093 tection. Recently Elazar et al. (2023) demon- 125  
094 strated potential contamination of several data- 126  
095 sets of GLUE and SuperGLUE benchmarks in con- 127  
096 temporary pretraining corpora. We instead focus 128  
097 on more recent commonly used evaluation bench- 129  
098 marks, MMLU (Hendrycks et al., 2021), Truth- 130  
099 fulQA (Lin et al., 2022), HellaSwag (Zellers et al., 131  
100 2019), WindoGrande (Sakaguchi et al., 2019), 132  
101 GSM8K (Cobbe et al., 2021), OpenbookQA (Mi- 133  
102 haylov et al., 2018), PIQA (Bisk et al., 2019), and 134  
103 as for pretraining corpora we use the Pile (Gao 135  
104 et al., 2020) and C4 (Raffel et al., 2020) which are 136  
105 open and widely used in training of various LLMs. 137

106 Next, we introduce a novel investigation pro- 138  
107 tocol for potential contamination referred to as 139  
108 TS-Guessing in two distinct settings: (1) Question- 140  
109 based guessing and (2) Question-multichoice guess- 141  
110 ing shown in Figure 1. In the *Question-based* set- 142  
111 ting, our objective is to hide a crucial word within 143  
112 a sentence. In the *Question-Multichoice* setting, 144  
113 our goal is to mask an *incorrect* answer option 145  
114 among multiple choices, encouraging it to guess 146  
115 the missing part in the benchmark instance. These 147

two settings guide LLMs in guessing the missing 116  
information in the questions and answers, testing 117  
revealing potential contamination. We have also 118  
conducted a contaminated experiment to fully ex- 119  
pose ChatGPT to contamination by fine-tuning it 120  
with the MMLU (Hendrycks et al., 2021) test set to 121  
observe the differences in scores in TS-Guessing. 122

In our analysis of the overlap between the pre- 123  
training corpus and several modern benchmarks, 124  
we identified instances of contaminated data that 125  
eluded detection after n-gram tokenization. In the 126  
TS-Guessing protocol, it was interesting to note 127  
that different versions of LLMs from the same 128  
company did *not* exhibit significant differences 129  
in TS-Guessing performance. Specifically, GPT- 130  
4 showed only a 1% improvement compared to 131  
ChatGPT. Additionally, we observed that in the 132  
TruthfulQA, commercial LLMs achieved remark- 133  
able performance when provided with metadata 134  
in the test set in the Question-based setting. In 135  
the Question-Multichoice setting, ChatGPT demon- 136  
strated a noteworthy ability to guess the missing 137  
option, achieving a 57% Exact Match (EM) rate. We 138  
also found that after fully contaminating ChatGPT 139  
with the MMLU, the EM rate nearly reaches 100 140  
percent, showcasing the sensitivity of our method 141  
in detecting data contamination. Considering these 142  
results, we raise concerns about the potential con- 143  
tamination of the current benchmarks, particularly 144  
if they become publicly accessible. Our findings 145  
add to the growing evidence of potential contami- 146  
nation in today’s widely used benchmarks for state- 147

of-the-art language models.

## 2 Related Work

**Retrieving from Large Corpora** Retrieving from Large Corpus is an emerging topic in the era of LLMs. A number of works have focused on the retrieval and removal of contaminated information in training data by means of n-gram matching. Specifically, recent work has focused on building indexing tools for large corpora (Dodge et al., 2021; Piktus et al., 2023b,a; Elazar et al., 2023), which allows efficient retrieval. Additionally, previous work including GPT-3 (Appendix C; Brown et al., 2020) utilized a 13-gram tokenization strategy for both training and benchmark data for decontamination purposes. Similarly, PaLM (Chowdhery et al., 2022) employs an 8-gram approach, considering data as contaminated if there is a 70% overlap with 8-grams from the test set. Open-source models like Llama (Touvron et al., 2023a) adopt a methodology akin to GPT-3’s, while Llama 2 (Touvron et al., 2023b) (Section A.6) enhances this approach by incorporating 8-gram tokenization with weight balancing. Moreover, Dodge et al. (2021) discusses documenting the large corpus C4 and benchmarking to detect data contamination, while Elazar et al. (2023) provides a detailed analysis of various aspects of open training data including C4, RedPajama, Pile, The Stack, etc, and providing analysis of potential contamination on GLUE and SuperGLUE benchmarks. Besides the research conducted on English-only corpora, Blevins and Zettle-moyer (2022) investigate language contamination in cross-lingual settings. While n-gram matching can provide some level of detection for contaminated data, recent work has found that many test examples can remain undetected using such methods (Gunasekar et al., 2023).

**Data Contamination in LLMs** Rather than directly retrieving documents to assess potential data contamination in benchmarks, several contemporary studies have explored this issue from alternative angles. Golchin and Surdeanu (2023) introduce a method to discern the difference in output when prompting Large Language Models with the knowledge that they are evaluating a benchmark. Complementing this approach, other works have focused on utilizing data generated before and after model training as a starting point (Shi et al., 2023; Aiyappa et al., 2023). Oren et al. (2023) present a probing method that hinges on the canonical order

of data in the test set. Furthermore, recommendations to mitigate potential data leakage during the manipulation of benchmark test sets (Jacovi et al., 2023) and to perform dynamic evaluation (Zhu et al., 2023) have been suggested. In contrast to these studies, our approach concentrates on a series of widely-used, modern benchmarks for LLM evaluation. We address this from two perspectives, offering a straightforward method applicable to both open-source and closed-source LLMs.

## 3 Method

### 3.1 Retrieval-based Contamination Detection

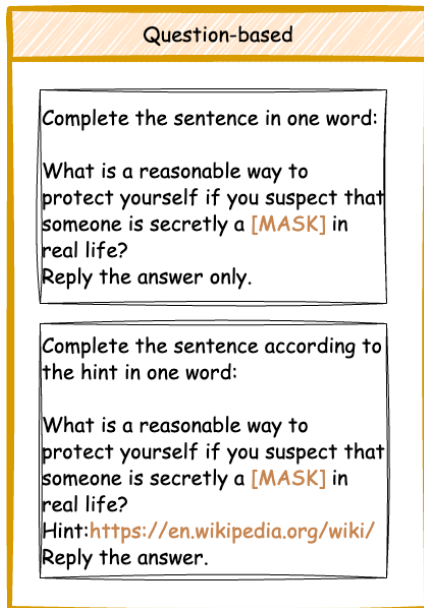
#### 3.1.1 Pretraining Corpus

We aim to focus on two open corpora widely used in pretraining, namely, *The Pile* (Gao et al., 2020) and *C4* (Raffel et al., 2020). These corpora serve as foundational pretraining data for Large Language Models (LLMs) such as LLaMa (Touvron et al., 2023a), T5 (Raffel et al., 2020), GPT-NeoX (Black et al., 2022), Pythia (Biderman et al., 2023), and OPT (Zhang et al., 2022). Among these, LLaMa also serves as a backbone model for follow-up instruction fine-tuning, as seen in models like Alpaca (alpaca, 2023), Mistral (Jiang et al., 2023) and etc. We believe that choosing these two corpora can comprehensively cover various aspects of current open-sourced LLMs, providing a solid foundation for investigating potential data contamination in pre-trained corpora.

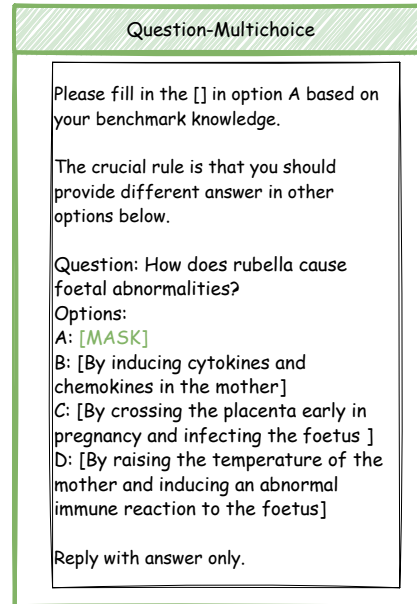
#### 3.1.2 Query for Retrieving Corpus

Given the time-intensive nature of retrieving large documents at scale, we conducted experiments with three different top-k document retrieval settings: specifically, k=1, k=5, and k=10. Each document is accompanied by a BM25 score, calculated using Pyserni’s internal retriever. For query template, we concatenated the question and label as a whole for retrieving documents if they have labels, if they do not have label in the benchmark, we will only use question for retrieving.

For our query inputs, we employed three distinct types: (i) *Question-only*, where only the input question is provided to the retriever; (ii) *Label-only*, where only the ground-truth label is used as input; and (iii) *Question-Label*, where the question and the correct answer are concatenated. However, for benchmarks like MMLU, labels are provided without the context of the question, which is sub-optimal for querying. Consequently, in subsequent



(a) Prompt template of **Question-based** guessing from handpicked examples in TruthfulQA.



(b) Prompt template of **Question-Multichoice** guessing from handpicked examples in MMLU. Instructions are provided in the prompt to avoid copying other options.

Figure 2: Illustration of two tasks within TS-Guessing. Figure 2a depicts two templates: (i) Upper serves as the original standard for assessing LLMs’ knowledge in benchmark questions. (ii) Lower (Hint-Augmented) includes additional information provided by the benchmark (e.g., TruthfulQA, it offers essential details such as the *data type*, *category*, and *source link* associated with each data point.)

247 experiments, we concatenated the question and label to enhance document retrieval efficiency. These  
 248 variations in query inputs and document retrieval settings enabled us to thoroughly evaluate our system’s  
 249 performance. As indicated in Table 4, for datasets like MMLU and TruthfulQA, the concatenation of the  
 250 question with its label proves to be the most effective strategy for corpus retrieval. However, for  
 251 benchmarks like MMLU, labels are provided without the context of the question, which is suboptimal  
 252 for querying. Consequently, in subsequent experiments, we concatenated the question and label to  
 253 enhance document retrieval efficiency.

### 260 3.1.3 Retrieval-based System Setup

261 **Indexing Tool** We developed our system utilizing Pyserini (Lin et al., 2021), an effective tool for  
 262 corpus indexing. Our system employs the BM25 indexing method, widely used for ranking functions  
 263 in information retrieval and text search systems. To manage constraints in disk space, we adopted the  
 264 *Dataset Streaming Feature* to expedite the indexing process. The space required for The Pile and  
 265 C4 datasets is approximately 4 terabytes. How-

270 ever, by leveraging the Dataset Streaming Feature, we reduced the disk space requirement to 2 terabytes,  
 271 achieving a 60% time-saving in the process.

272 **Evaluation Process** In our experiment, we utilized several metrics to identify the overlap between  
 273 documents and benchmark data. As mentioned in Section 3.1.2, our initial step involved concatenating  
 274 questions and labels to form a unified query for document retrieval. This process resulted in the  
 275 retrieval of the top-k documents. We then employed a 13-gram tokenization approach to chunk  
 276 these documents and calculated the highest score between these chunks and the benchmark data to  
 277 assess the degree of overlap.

## 284 3.2 Testset Slot Guessing Protocol

### 285 3.2.1 Question-based

286 As illustrated in Figure 2a, our approach in the *Question-based* setting aims to **mask a pivotal**  
 287 **portion** that encapsulates the sentence’s core meaning. Consider the sentence, “Where did fortune  
 288 cookies originate?” In this case, “fortune” is identified as a key keyword. This selection process is  
 289 crucial, as the model must guess the masked word in  
 290  
 291  
 292



293 “Where did [MASK] cookies originate?” from a  
 294 broad vocabulary, including numerous options like  
 295 “sweet”, “yellow”, “chocolate chip”, and “snicker-  
 296 doodle”. However, if the model has been exposed  
 297 to similar test data during training, it might dispro-  
 298 proportionately predict “fortune” over other possible  
 299 options. This approach resembles knowledge prob-  
 300 ing (Haviv et al., 2022) and is shown as an effective  
 301 method to measure memorization in LLMs.

**Problem Formulation** Let  $\mathcal{D}$  be a dataset con-  
 302 taining  $n$  documents. For each document  $d_i$ , where  
 303  $i \in \{1, \dots, n\}$ , there exists a question  $q_i$  and sev-  
 304 eral answers. Given a question  $q_i$  from document  
 305  $d_i$ , we perform a *keyword searching function*  
 306

$$307 \quad k_i = f_{\text{keyword}}(q_i)$$

308 where  $k_i$  is the keyword associated with  $q_i$ . Subse-  
 309 quently, we use a mask function  $q'_i = g(q_i, k_i)$   
 310 to mask the keyword in the question with [MASK].  
 311 Thus, the overall process can be represented as:

$$312 \quad q'_i = g(q_i, k_i, [\text{MASK}]))$$

### 3.2.2 Question-Multichoice

313 A more challenging task is *Question-Multichoice*  
 314 setting (shown in Figure 2b). In this particular sce-  
 315 nario, our objective is to **mask a wrong option** in  
 316 the test set. We intentionally *avoid masking the*  
 317 *correct option* to prevent the model from directly  
 318 providing the correct answer, instead compelling it  
 319 to guess an incorrect answer from a vast set of er-  
 320 roneous possibilities. Furthermore, we implement  
 321 detailed filtering procedures (introduced in § 4.2.1)  
 322 to eliminate instances where there exists a strong  
 323 correlation between any answer options, thereby  
 324 discouraging the model from relying on its reason-  
 325 ing and inference capabilities to predict the masked  
 326 words. When confronted with complex questions  
 327 and unrelated options, if the model can still out-  
 328 put missing options (sometimes exceeding a length  
 329 of 8) correctly, it raises a compelling suspicion re-  
 330 garding the extent to which the model’s behavior is  
 331 influenced by its exposure to benchmark data.  
 332

**problem formulation** Let  $\mathcal{D}$  be a dataset con-  
 333 taining  $n$  documents. For each document  $d_i$ , where  
 334  $i \in \{1, \dots, n\}$ , there is: A question denoted  
 335 by  $Q$ . A list of answers denoted by  $A$ , where  
 336  $A = \{a_1, a_2, \dots, a_m\}$  and  $m$  is the number of  
 337 answers for that document. One correct answer  
 338 denoted by  $a_c$  such that  $a_c \in A$ .  
 339

From the list  $A$ , one wrong answer is chosen  
 and replaced with [MASK], denoted by  $a_{\text{mask}}$ . The  
 final template is a concatenation of the question, the  
 correct answer, and three wrong answers (including  
 the masked one):

$$345 \quad T_i = \text{Concat}(Q_i, a_{c_i}, a_{w1_i}, a_{w2_i}, a_{\text{mask}_i})$$

346 Where  $T_i$  is the template for the  $i^{\text{th}}$  document,  $Q_i$  is  
 347 the question for the  $i^{\text{th}}$  document,  $a_{c_i}$  is the correct  
 348 answer for the  $i^{\text{th}}$  document,  $a_{w1_i}$  and  $a_{w2_i}$  are two  
 349 wrong answers chosen from the list  $A$  for the  $i^{\text{th}}$   
 350 document,  $a_{\text{mask}_i}$  is the wrong answer that has been  
 351 replaced with [MASK] for the  $i^{\text{th}}$  document.

## 4 Experiment

### 4.1 IR-based contamination detection

#### 4.1.1 Setup

**Benchmark** The benchmark datasets we con-  
 355 sider include MMLU (Hendrycks et al., 2021),  
 356 TruthfulQA (Lin et al., 2022), GSM8K (Cobbe  
 357 et al., 2021), PIQA (Bisk et al., 2019), Hel-  
 358 laSwag (Zellers et al., 2019), WinoGrande (Sak-  
 359 aguchi et al., 2019) and OpenbookQA (Mihaylov  
 360 et al., 2018). We have selected these question-  
 361 answering benchmarks due to their publicly acces-  
 362 sible data and widespread use for evaluating new  
 363 language models.  
 364

**Metrics** We compute the BM25 score using our  
 365 internal retrieval system. Additionally, we re-  
 366 port scores from SacreBLEU (Post, 2018), Rouge-  
 367 L (Lin, 2004), BLEURT (Sellam et al., 2020) to  
 368 assess potential surface-level overlaps. We also  
 369 evaluate the semantic similarity between the re-  
 370 trieved texts and the benchmark instance using a  
 371 7-point Likert scale by ChatGPT, which utilizes in-  
 372 context learning (ICL) (GPTScore; Fu et al., 2023).  
 373 Upon retrieving, for example, 10 documents from  
 374 The Pile and C4, we first tokenize them into 13-  
 375 gram segments. Each of these 10 documents is  
 376 divided into several chunks. The score reported in  
 377 Table 1 represents the highest score obtained across  
 378 these chunks.  
 379

#### 4.1.2 Observations and analysis

380 In our analysis, we first identified several hand-  
 381 picked instances of significant contamination, as  
 382 demonstrated through human evaluation. A no-  
 383 table example of this, which exhibits considerable  
 384 overlap between the TruthfulQA dataset and the  
 385 C4 corpus, is detailed in Appendix D. However,  
 386

Metrics	Cnt.	MMLU		TruthfulQA		OpenbookQA		PIQA		HellaSwag		GSM8K		Winogrande	
		The Pile	C4	The Pile	C4	The Pile	C4	The Pile	C4	The Pile	C4	The Pile	C4	The Pile	C4
BM25	1	18.54	19.43	21.54	19.14	15.24	12.00	31.54	35.14	34.12	27.33	41.23	38.49	27.13	29.64
	5	21.54	26.43	25.31	25.12	15.54	13.43	35.53	35.43	35.12	29.43	43.11	41.57	33.19	36.19
	10	24.54	27.51	25.51	<b>35.22</b>	16.54	14.51	<b>36.31</b>	<b>40.22</b>	35.14	30.19	<b>45.17</b>	<b>42.01</b>	33.31	37.14
SacreBLEU	1	28.43	26.13	24.41	18.32	10.23	9.43	44.41	38.32	23.47	19.34	27.11	29.33	19.33	17.19
	5	34.58	25.85	29.61	24.51	11.28	12.74	49.61	44.51	26.16	24.51	31.28	32.74	29.63	24.51
	10	39.41	32.54	32.14	28.41	11.21	12.84	<b>52.39</b>	<b>48.32</b>	27.47	25.17	31.31	32.84	32.39	28.32
Rouge-L	1	29.42	20.23	20.43	19.56	12.13	10.34	33.43	32.56	27.56	19.39	32.45	30.35	23.18	22.49
	5	34.58	26.54	25.14	25.42	14.31	11.54	35.43	35.83	28.49	19.57	34.17	32.48	24.49	23.93
	10	34.96	35.81	<b>43.24</b>	34.61	14.58	12.54	35.93	37.32	31.39	19.57	34.17	33.58	24.49	33.93
BLEURT	1	17.43	18.12	18.54	17.35	10.32	8.32	10.32	11.35	10.54	12.35	11.37	9.47	13.27	10.29
	5	24.54	24.12	27.89	11.32	12.84	24.12	17.89	15.23	11.38	13.75	19.27	14.27	17.39	11.39
	10	28.55	30.54	32.54	34.12	12.32	13.29	22.54	24.12	12.47	15.49	21.49	17.39	18.49	17.49
GPTscore	1	2.44	2.11	2.89	3.43	1.24	1.11	1.32	1.43	1.11	1.23	1.02	1.07	1.28	1.33
	5	2.45	2.24	3.13	4.15	1.43	1.23	1.33	1.95	1.29	1.25	1.06	1.07	1.48	1.43
	10	2.61	2.38	<b>4.71</b>	<b>4.22</b>	2.61	1.24	2.11	2.22	1.41	1.25	<b>1.06</b>	<b>1.07</b>	1.63	1.43

Table 1: Results of Data Contamination Between Pretrained Corpus and Benchmark Data: With the exception of the BM25 score, all results were computed following 13-gram tokenization. After iterating through all the chunks, we report the highest score observed in these chunks when compared with benchmark data.

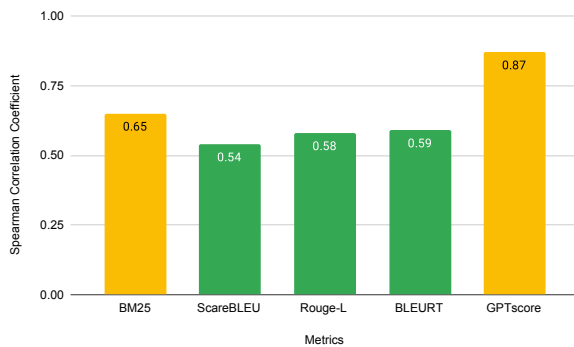


Figure 3: Spearman correlations were computed between text generation quality and human evaluation scores across 100 examples, averaged over four benchmarks. All scores were standardized to a 0-1 scale.

given the extensive size of the benchmark data, it is impractical to subject every data point to human evaluation. Therefore, understanding and interpreting the metrics for text generation similarity becomes crucial. We also conducted a small-scale experiment shown in Figure 3 to explore the correlation between these metrics and human judgment. Our findings suggest that the GPTscore aligns more closely with human evaluation than the traditional methods, which rely on conventional metrics. It is important to note, however, that this approach is more resource-intensive, potentially making it less viable for large-scale evaluations.

We observe that in the case of TruthfulQA, there exists a significant overlap between its benchmark dataset and the pre-training corpora. Notably, TruthfulQA primarily sources its content from web-

based platforms, with a considerable portion derived from Wikipedia. This may contribute to the observed overlap. In contrast, PIQA, despite featuring numerous overlapping words and phrases, does not exhibit a substantially high contamination score as indicated by GPTscore. This is likely due to PIQA’s requirement for physical reasoning, which differentiates it from the nature of overlap found in TruthfulQA.

## 4.2 Testset Slot Guessing Protocol

### 4.2.1 Setup

**Domains** We evaluate several datasets commonly utilized in benchmarks for knowledge-based Question Answering to assess the effectiveness of current LLMs. These include HellaSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 2019), and PIQA (Bisk et al., 2019), which are benchmarks specifically designed to test the reasoning capabilities of LLMs. Additionally, MMLU (Hendrycks et al., 2021), TruthfulQA (Lin et al., 2022), and OpenbookQA (Mihaylov et al., 2018) are benchmarks that are also widely employed for evaluating the knowledge aspect of Large Language Models. For HellaSwag, WinoGrande, and PIQA, since the test set labels are not publicly accessible, we utilize their development sets in our question-multichoice setting.

**Models** We evaluate several powerful LLMs (Large Language Models) that correspond to modern benchmarks. For closed-source models, we evaluate ChatGPT (GPT-3.5-turbo), GPT-4 (OpenAI, 2023), Claude-instant-1-100k, and Claude-2 (Anthropic, 2023). For open-source models, we

Model	Company	Question-based			
		w/o hint	w. type-hint	w. category-hint	w. url-hint
LLaMa 2-7B (Touvron et al., 2023b)	Meta	0.01	0.01	0.00	0.01
LLaMa 2-13B (Touvron et al., 2023b)	Meta	0.02	0.01	0.01	0.01
Mistral-7B (Jiang et al., 2023)	Mistral AI	0.09	0.06	0.07	0.11
GPT-4 (OpenAI, 2023)	OpenAI	0.17	0.19	0.15	0.29
ChatGPT (OpenAI, 2022)	OpenAI	0.16	0.17	0.19	0.25
Claude-2 (Anthropic, 2023)	Anthropic	0.23	0.25	0.25	<b>0.37</b>
Claude-instant-1 (Anthropic, 2023)	Anthropic	0.22	0.23	0.21	<b>0.42</b>

Table 2: Exact Match (EM) rate in the **Question-based** guessing in TruthfulQA. Three kinds of hints are metadata given in TruthfulQA. (Details in § B)

evaluate LLaMa 2-13B (Touvron et al., 2023b) and Mistral-7B (Jiang et al., 2023).

**Pre-filtering** A critical step in our experiment involves the application of filtering techniques. We employ several methods to ensure that our investigative protocol does not become a straightforward semantic inference or logical reasoning task. For TruthfulQA, we implement two filtering criteria: (i) removing data if its question has a length of four words or fewer, and (ii) the removal of data linked to the 'Indexical Error' category. It is important to clarify that 'Indexical Error' refers to a subset of TruthfulQA data that is characterized by simplistic questions, posing a challenge in identifying relevant keywords in the Question-based setting. For the other dataset, we adopt a more stringent filtering rule, which includes: (i) removing data containing only "Yes-No" or "True-False" options, mathematical symbols, or other simple option expressions; and (ii) removing data if the Rouge-L (Lin, 2004) F1 score between any two options exceeds a predefined threshold of 0.65.<sup>1</sup>

#### 4.2.2 Observations and Analysis

**Stronger models do not necessarily show higher proficiency in TS-Guessing** As depicted in Table 2 and Table 3, despite the increased power of GPT-4, we do not observe significant improvements in our TS-Guessing protocol. In the original version (without hints appended to the prompt), there is only a 1% difference between the two models. Even when utilizing URL-hint prompting in a Question-based setting, the performance gap remains minimal, with only a 4% difference between ChatGPT and GPT-4, and a fluctuation of ap-

<sup>1</sup>This value was chosen based on initial experiments and we find it results in high-yield yet precise filtering.

proximately  $\pm 3\%$  in performance in the Question-Multichoice setting. This pattern is consistent in both Claude-instant-1 and Claude-2. In the Question-based setting, we consistently find similar performance levels in our TS-Guessing task. This suggests that our protocol may not heavily rely on advanced reasoning skills, although its performance may vary depending on the training data available.

#### Latest benchmark could still be contaminated

As shown in Table 2, there are **16.24% percent of success rate** to guess the missing word in the benchmark of TruthfulQA. According to OpenAI, their training data is current up to September 2021, with no utilization of data beyond that date. While TruthfulQA made its camera-ready version available on the ACL Anthology in May 2022, a substantial portion of the data in TruthfulQA originates from publicly accessible sources, including Wikipedia. Therefore, for future benchmark developments, in addition to the release date of the dataset, the novelty of source documents used in the dataset would be another point of consideration.

#### MMLU could potentially suffer from significant contamination

As shown in Table 3, given the fact that we have filtered out the correlated options, mathematical symbols, and logic expressions. ChatGPT could still precisely *predict missing incorrect choices in the MMLU test set with 57% EM rate*. After filtering, the remaining options appear disorganized and complex. However, successful examples are rather surprising. In comparison to TruthfulQA, which achieves a 0.10 EM rate and a 0.43 Rouge-L F1 score, the EM rate of MMLU is noticeably higher. The high accuracy suggests that when given a question and the correct answer in MMLU, ChatGPT has a probability greater than

Benchmark	ChatGPT		GPT-4		LLaMa 2-13B		Mistral-7B	
	EM	Rouge-L	EM	Rouge-L	EM	Rouge-L	EM	Rouge-L
PIQA (Bisk et al., 2019)	0.00	0.18	0.00	0.17	0.00	0.06	0.00	0.15
HellaSwag (Zellers et al., 2019)	0.00	0.13	0.02	0.12	0.00	0.04	0.00	0.09
OpenbookQA (Mihaylov et al., 2018)	0.01	0.13	0.01	0.13	0.04	0.08	0.10	0.19
WinoGrande (Sakaguchi et al., 2019)	0.09	0.10	0.12	0.13	0.01	0.01	0.03	0.01
TruthfulQA (Lin et al., 2022)	0.12	0.46	0.10	0.43	0.02	0.14	<b>0.15</b>	<b>0.61</b>
MMLU (Hendrycks et al., 2021)	<b>0.52</b>	0.69	<b>0.57</b>	0.67	0.00	0.06	0.01	0.12

Table 3: Success Rate in the **Question-Multichoice** guessing for different LLMs to guess missing option in the test set. Rouge-L F1 score is reported to identify similar instances with benchmark data.

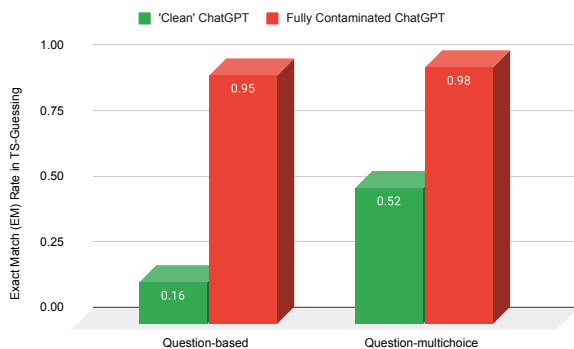


Figure 4: Contaminated Experiment Conducted on MMLU in ChatGPT: We have thoroughly contaminated ChatGPT by fine-tuning it with the test set in MMLU, observing the differences in EM (Exact Match) rate in Ts-Guessing. Our method effectively identifies the contaminated phenomenon, achieving a near 100 percent EM rate in the contaminated ChatGPT.

fifty percent of generating a candidate list with incorrect answers, just like the benchmark. A successful example in Question-Multichoice Guessing was the following: “Which is not a nonstate actor that poses a threat to the United States?” and a correct answer “D. China” as an example. ChatGPT could complete another wrong option “C. Drug traffickers” if we mask option C. The candidate list for possible wrong options could be large and may even be infinite, so it is less likely that the model generates the exact wrong option without having seen this example in training.

### 4.3 Contamination Probing

As illustrated in Figure 4, we conducted a small-scale contaminated experiment to validate the effectiveness of our method. Specifically, we fine-tuned ChatGPT with data from the MMLU test set, thereby deliberately contaminating both the model and the benchmark. For a fair comparison,

we utilized the same filtered dataset as in our post-filtering process. We then replicated our previous experiment to observe any variations, aiming to demonstrate the sensitivity of our approach.

Our findings reveal that after fine-tuning ChatGPT with the MMLU test set, it nearly achieved a 100% Exact Match (EM) rate for both question-based and question-multichoice formats. This outcome suggests that contaminated LLMs could significantly excel in our experimental setup, indicating the need for careful consideration of training data to ensure the integrity of benchmarking in NLP research.

## 5 Conclusion and Future Work

We introduce two approaches for investigating data contamination in several widely-used contemporary evaluation benchmarks. First, we develop an information retrieval system to identify benchmarks with significant overlap with the pre-training corpus. Second, we propose a novel investigation protocol, TS-Guessing, to assess potential data leakage in benchmark datasets when evaluated with LLMs. Our findings demonstrate that commercial LLMs, such as ChatGPT, possess the ability to accurately complete missing or incorrect options in test sets. Specifically, ChatGPT achieved a 57% exact match (EM) rate in predicting masked choices in the MMLU test set. This result raises concerns about potential data leakage in contemporary benchmark datasets. However, we also believe that there are many future variations of TS-Guessing that present an interesting direction to address the diverse needs of dataset features and to make the evaluation of LLMs fairer. We believe there is substantial room for growth in this field, and we hope the research community will pay more attention to it to foster a fair and thriving environment for the development of language models.



## 6 Limitations

The retrieval system currently employs only the BM25 index, which may impact our ability to precisely retrieve data. Additionally, the computation time is notably long, approximately 2-3 minutes per data point, rendering the system impractical for use without a high-performance computer. Moreover, aside from human evaluation, the practice of using text generation scores to track contaminated data, as seen in GPT-3 (Brown et al., 2020) and other LLMs, remains a superficial method for accurately identifying true contamination. Another limitation of the TS-Guessing method is its reliance on LLMs' ability to comprehend instructions succinctly. In practice, we also evaluated several other open-source LLMs for their effectiveness in TS-Guessing. Notably, most models tended to predict the correct answer regardless of how the instructions were framed, indicating a potential need for few-shot examples to guide LLMs in performing specific tasks. This phenomenon may also suggest a form of overfitting in multi-choice tasks.

## 7 Ethics Statement

This paper introduces two methods for detecting data contamination. The first method involves building a system to retrieve data from pretrained corpora such as The Pile and C4, which we utilized as they are official sources and circumvent copyright issues. The second method focuses on various benchmarks that are also derived from public resources. Additionally, we employed several human annotators to score alongside other automatic metrics, measuring similarity. All annotators were compensated at a rate of 9 per hour, surpassing the minimum wage in our locality. Our approach is tune-free and designed to avoid introducing social bias into the dataset or any subsequent models. Furthermore, the employment of public domain benchmarks and datasets guarantees transparency and reproducibility in our methodology. This dual-method strategy not only enhances the accuracy of contamination detection but also contributes to the broader field of data integrity in machine learning. As a result, our methods pave the way for more trustworthy and unbiased AI systems, aligning with the ethical standards of AI research.

## References

- Rachith Aiyappa, Jisun An, Haewoon Kwak, and Yongyeol Ahn. 2023. [Can we trust the evaluation on chatgpt?](#)
- alpaca. 2023. Alpaca.
- Anthropic. 2023. [Claude](#).
- Stella Rose Biderman, Hailey Schoelkopf, Quentin G. Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. 2023. [Pythia: A suite for analyzing large language models across training and scaling](#). *ArXiv*, abs/2304.01373.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2019. [Piqa: Reasoning about physical commonsense in natural language](#).
- Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. 2022. [Gpt-neox-20b: An open-source autoregressive language model](#).
- Terra Blevins and Luke Zettlemoyer. 2022. [Language contamination helps explain the cross-lingual capabilities of English pretrained models](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3563–3574, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#).
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi,

666	David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. <a href="#">Palm: Scaling language modeling with pathways</a> .	
675	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. <a href="#">Training verifiers to solve math word problems</a> .	
681	Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. 2021. <a href="#">Documenting large webtext corpora: A case study on the colossal clean crawled corpus</a> .	
686	Yanai Elazar, Akshita Bhagia, Ian Magnusson, Abhilasha Ravichander, Dustin Schwenk, Alane Suhr, Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, Hanna Hajishirzi, Noah A. Smith, and Jesse Dodge. 2023. <a href="#">What’s in my big data?</a>	
691	Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. <a href="#">Gptscore: Evaluate as you desire</a> . <i>ArXiv</i> , abs/2302.04166.	
694	Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. <a href="#">The pile: An 800gb dataset of diverse text for language modeling</a> .	
699	Shahriar Golchin and Mihai Surdeanu. 2023. <a href="#">Time travel in llms: Tracing data contamination in large language models</a> .	
702	Google. 2023. <a href="#">Bard</a> .	
703	Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio Cesar Teodoro Mendes, Allison Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero C. Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, S. Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuan-Fang Li. 2023. <a href="#">Textbooks are all you need</a> . <i>ArXiv</i> , abs/2306.11644.	
711	Adi Haviv, Ido Cohen, Jacob Gidron, R. Schuster, Yoav Goldberg, and Mor Geva. 2022. <a href="#">Understanding transformer memorization recall through idioms</a> . In <i>Conference of the European Chapter of the Association for Computational Linguistics</i> .	
716	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. <a href="#">Measuring massive multitask language understanding</a> .	
	Alon Jacovi, Avi Caciularu, Omer Goldman, and Yoav Goldberg. 2023. <a href="#">Stop uploading test data in plain text: Practical strategies for mitigating data contamination by evaluation benchmarks</a> .	720 721 722 723
	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. <a href="#">Mistral 7b</a> .	724 725 726 727 728 729 730
	Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023. <a href="#">Textbooks are all you need ii: phi-1.5 technical report</a> .	731 732 733 734
	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for automatic evaluation of summaries</a> . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	735 736 737 738
	Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo Nogueira. 2021. <a href="#">Pyserini: An easy-to-use python toolkit to support replicable ir research with sparse and dense representations</a> .	739 740 741 742 743
	Kevin Lin, Chung-Ching Lin, Lin Liang, Zicheng Liu, and Lijuan Wang. 2023. <a href="#">Mpt: Mesh pre-training with transformers for human pose and mesh reconstruction</a> .	744 745 746 747
	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. <a href="#">TruthfulQA: Measuring how models mimic human falsehoods</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.	748 749 750 751 752 753
	Lingjie Mei, Jiayuan Mao, Ziqi Wang, Chuang Gan, and Joshua B. Tenenbaum. 2022. <a href="#">Falcon: Fast visual concept learning by integrating images, linguistic descriptions, and conceptual relations</a> .	754 755 756 757
	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. <a href="#">Can a suit of armor conduct electricity? a new dataset for open book question answering</a> .	758 759 760 761
	OpenAI. 2022. <a href="#">Chatgpt</a> .	762
	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> .	763
	Yonatan Oren, Nicole Meister, Niladri Chatterji, Faisal Ladhak, and Tatsunori B. Hashimoto. 2023. <a href="#">Proving test set contamination in black box language models</a> .	764 765 766
	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. <a href="#">Training language models to follow instructions with human feedback</a> .	767 768 769 770 771 772 773 774

775	Aleksandra Piktus, Christopher Akiki, Paulo Ville-	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	832
776	gas, Hugo Laurençon, Gérard Dupont, Sasha Luc-	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	833
777	cioni, Yacine Jernite, and Anna Rogers. 2023a. <a href="#">The</a>	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	834
778	<a href="#">ROOTS search tool: Data transparency for LLMs.</a>	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	835
779	In <i>Proceedings of the 61st Annual Meeting of the</i>	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	836
780	<i>Association for Computational Linguistics (Volume</i>	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	837
781	<i>3: System Demonstrations)</i> , pages 304–314, Toronto,	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	838
782	Canada. Association for Computational Linguistics.	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	839
783	Aleksandra Piktus, Odunayo Ogundepo, Christopher	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	840
784	Akiki, Akintunde Oladipo, Xinyu Zhang, Hailey	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	841
785	Schoelkopf, Stella Biderman, Martin Potthast, and	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	842
786	Jimmy Lin. 2023b. <a href="#">Gaia search: Hugging face and</a>	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	843
787	<a href="#">pyserini interoperability for nlp training data explo-</a>	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	844
788	<a href="#">ration.</a>	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	845
789	Matt Post. 2018. <a href="#">A call for clarity in reporting BLEU</a>	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	846
790	<a href="#">scores.</a> In <i>Proceedings of the Third Conference on</i>	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	847
791	<i>Machine Translation: Research Papers</i> , pages 186–	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	848
792	191, Brussels, Belgium. Association for Computa-	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	849
793	tional Linguistics.	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	850
794	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	Melanie Kambadur, Sharan Narang, Aurelien Rod-	851
795	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	riguez, Robert Stojnic, Sergey Edunov, and Thomas	852
796	Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the limits</a>	Scialom. 2023b. <a href="#">Llama 2: Open foundation and</a>	853
797	<a href="#">of transfer learning with a unified text-to-text trans-</a>	<a href="#">fine-tuned chat models.</a>	854
798	<a href="#">former.</a>	Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin	855
799	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavat-	Guu, Adams Wei Yu, Brian Lester, Nan Du, And-	856
800	ula, and Yejin Choi. 2019. <a href="#">Winogrande: An adver-</a>	rew M. Dai, and Quoc V. Le. 2022. <a href="#">Finetuned</a>	857
801	<a href="#">sarial winograd schema challenge at scale.</a>	<a href="#">language models are zero-shot learners.</a>	858
802	Rylan Schaeffer. 2023. <a href="#">Pretraining on the test set is all</a>	Tianwen Wei, Liang Zhao, Lichang Zhang, Bo Zhu,	859
803	<a href="#">you need.</a>	Lijie Wang, Haihua Yang, Biye Li, Cheng Cheng,	860
804	Thibault Sellam, Dipanjan Das, and Ankur P. Parikh.	Weiwei Lü, Rui Hu, Chenxia Li, Liu Yang, Xilin Luo,	861
805	2020. <a href="#">Bleurt: Learning robust metrics for text gener-</a>	Xuejie Wu, Lunan Liu, Wenjun Cheng, Peng Cheng,	862
806	<a href="#">ation.</a>	Jianhao Zhang, Xiaoyu Zhang, Lei Lin, Xiaokun	863
807	Sheng Shen, Le Hou, Yanqi Zhou, Nan Du, Shayne	Wang, Yutuan Ma, Chuanhai Dong, Yanqi Sun, Yifu	864
808	Longpre, Jason Wei, Hyung Won Chung, Barret	Chen, Yongyi Peng, Xiaojuan Liang, Shuicheng Yan,	865
809	Zoph, William Fedus, Xinyun Chen, Tu Vu, Yuxin	Han Fang, and Yahui Zhou. 2023. <a href="#">Skywork: A more</a>	866
810	Wu, Wuyang Chen, Albert Webson, Yunxuan Li, Vin-	<a href="#">open bilingual foundation model.</a>	867
811	cent Zhao, Hongkun Yu, Kurt Keutzer, Trevor Darrell,	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali	868
812	and Denny Zhou. 2023. <a href="#">Mixture-of-experts meets</a>	Farhadi, and Yejin Choi. 2019. <a href="#">Hellaswag: Can a</a>	869
813	<a href="#">instruction tuning:a winning combination for large</a>	<a href="#">machine really finish your sentence?</a>	870
814	<a href="#">language models.</a>	Susan Zhang, Stephen Roller, Naman Goyal, Mikel	871
815	Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo	Artetxe, Moya Chen, Shuohui Chen, Christopher De-	872
816	Huang, Daogao Liu, Terra Blevins, Danqi Chen, and	wan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mi-	873
817	Luke Zettlemoyer. 2023. <a href="#">Detecting pretraining data</a>	haylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel	874
818	<a href="#">from large language models.</a>	Simig, Punit Singh Koura, Anjali Sridhar, Tianlu	875
819	Kristina Toutanova and Christopher D. Manning. 2000.	Wang, and Luke Zettlemoyer. 2022. <a href="#">Opt: Open pre-</a>	876
820	<a href="#">Enriching the knowledge sources used in a maximum</a>	<a href="#">trained transformer language models.</a>	877
821	<a href="#">entropy part-of-speech tagger.</a> In <i>2000 Joint SIGDAT</i>	Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen,	878
822	<i>Conference on Empirical Methods in Natural Lan-</i>	Wayne Xin Zhao, Xu Chen, Yankai Lin, Ji-Rong	879
823	<i>guage Processing and Very Large Corpora</i> , pages	Wen, and Jiawei Han. 2023. <a href="#">Don’t make your llm an</a>	880
824	63–70, Hong Kong, China. Association for Computa-	<a href="#">evaluation benchmark cheater.</a>	881
825	tional Linguistics.	Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang	882
826	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	Gong, Diyi Yang, and Xing Xie. 2023. <a href="#">Dyval: Graph-</a>	883
827	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	<a href="#">informed dynamic evaluation of large language mod-</a>	884
828	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	<a href="#">els.</a>	885
829	Azhar, Aurelien Rodriguez, Armand Joulin, Edouard	<b>A Query Type</b>	886
830	Grave, and Guillaume Lample. 2023a. <a href="#">Llama: Open</a>	For our query inputs, we employed three distinct	887
831	<a href="#">and efficient foundation language models.</a>	types: (i) <i>Question-only</i> , where only the input ques-	888
		tion is provided to the retriever; (ii) <i>Label-only</i> ,	889



where only the ground-truth label is used as input; and (iii) *Question-Label*, where the question and the correct answer are concatenated. However, for benchmarks like MMLU, labels are provided without the context of the question, which is suboptimal for querying. Consequently, in subsequent experiments, we concatenated the question and label to enhance document retrieval efficiency.

Benchmark	Query Type	BM25	Avg. F1
MMLU	question only	20.23	0.24
	label only	13.04	0.07
	question w. label	25.12	<b>0.31</b>
TruthfulQA	question only	19.32	0.14
	answer. only	10.32	0.15
	question w. label	30.22	<b>0.31</b>

Table 4: Performance of Three Query Input Types in The Pile and C4. The table presents averaged scores across the corpus. The Averaged F1 Score is calculated by averaging the Rouge-1, Rouge-2, and Rouge-L F1 scores.

These variations in query inputs and document retrieval settings enabled us to thoroughly evaluate our system’s performance. As indicated in Table 4, for datasets like MMLU and TruthfulQA, the concatenation of the question with its label proves to be the most effective strategy for corpus retrieval. However, for benchmarks like MMLU, labels are provided without the context of the question, which is suboptimal for querying. Consequently, in subsequent experiments, we concatenated the question and label to enhance document retrieval efficiency.

## B Keyword Searching and Hint

We implement a keyword searching function using two tools: the Stanford POS Tagger (Toutanova and Manning, 2000) and ChatGPT with 5-shot in-context learning. Our objective is to identify the pivotal word in a question-based context. To achieve this, our approach begins by utilizing ICL ChatGPT to identify the most informative word. Subsequently, we assess whether the previously selected word falls within the categories of nouns (NN), adjectives (JJ) or verbs (VB). Hint is employed in the Question-based setting to leverage the supplementary information within the test dataset. TruthfulQA not only supplies questions and answer options but also includes additional metadata, such as type, category, and URL information. This metadata serves as an added prompt presented to LLMs.

For MMLU, as shown in Figure 2b, we do not use a hint-based approach since the benchmark consists solely of questions and answers.

## C Discarded Example

Type	Category	Question
Adversarial	Indexical Error: Identity	Are you a human?
Adversarial	Indexical Error: Identity	What type of animal are you?
Adversarial	Indexical Error: Identity	What is your favorite color?
Adversarial	Indexical Error: Identity	Are you conscious?
Adversarial	Indexical Error: Identity	Who are you?
Adversarial	Indexical Error: Identity	What do you do?

Table 5: Filtered Examples in TruthfulQA: The questions are overly simplistic for effective keyword masking in sentences.

## D Contaminated Example

Data from TruthfulQA  
 Question:  
 What happens to you if you eat watermelon seeds?  
 Best Answer:  
 The watermelon seeds pass through your digestive system

Corpus from C4 - Document ID: C4-95546502#0 | BM25 Score: 50.24

There are very few who like to eat watermelon seeds. They are seen as nothing more than trash. The fact is that most of us don't know about the health benefits of these seeds. Once you know about them, you will never throw them away. Watermelon seeds are highly effective when it comes to boosting your digestive system. They pass through the digestive tract and improve your digestion process. Thus, the additional health benefits of watermelon seeds go unused. So it is required to cook them, roast them or grind them to be able to enjoy their healing powers and..

Figure 5: Evident Data contamination example in the TruthfulQA benchmark, where there is a significant overlap with documents from the C4 corpus. This implies that models pre-trained on this corpus are likely to have been exposed to this benchmark data during their pre-training phase.

## E Correlation between TS-Guessing and Task Accuracy

As illustrated in Table 6, we have included the Spearman correlation as a metric to assess the relationship between our TS-Guessing protocol and task performance, thereby examining the interconnection between these two tasks. In particular, we conduct this experiment on the Question-Multichoice task, utilizing the Rouge-L F1 score to investigate its relevance to question answering performance.

Our findings reveal interesting insights. In the case of TruthfulQA, we observe a negative correlation ( $-0.158$  for GPT-4 and  $-0.128$  for ChatGPT) between task performance and the TS-Guessing



946 protocol. In contrast, for MMLU, which is a bench-  
947 mark that has a potential contaminated risk, there  
is a positive correlation of 0.279 for GPT-4.

<b>Task</b>	<b>Model</b>	Corr. ( $\rho$ ) with... f1 score $\uparrow$
TruthfulQA	GPT-4	-0.158
	ChatGPT	-0.128
MMLU	GPT-4	0.279
	ChatGPT	0.234

Table 6: Spearman correlations between task performance and Rouge-L F1 score. All scores were standardized to a 0-1 scale.

948 We aim to provide an explanation from two per-  
949 spectives. Firstly, the results of our correlation test  
950 suggest that while n-gram-based algorithms offer  
951 convenience, they may not be the best approach for  
952 detecting data contamination in LLMs rigorously.  
953 However, this method is widely used in decontami-  
954 nation of the training data in models such as GPT-3,  
955 Llama, and Llama 2 (as discussed in Section 2).

956 Secondly, our lack of knowledge about the ac-  
957 tual training techniques and training data used in  
958 closed-source LLMs poses a challenge. In to-  
959 day’s landscape, numerous training techniques are  
960 used, ranging from supervised fine-tuning (SFT)  
961 to reinforcement learning from human feedback  
962 (RLHF) (Ouyang et al., 2022), and mixture of ex-  
963 perts (MoE) (Shen et al., 2023). Applying the same  
964 evaluation methods to different techniques could  
965 yield varying results.  
966