

Data Contamination Quiz: A Tool to Detect and Estimate Contamination in Large Language Models

Shahriar Golchin*, Mihai Surdeanu

Department of Computer Science, University of Arizona, USA

{golchin, msurdeanu}@arizona.edu

Abstract

We propose the Data Contamination Quiz (DCQ), a simple and effective approach to *detect* data contamination in large language models (LLMs) and *estimate* the amount of it. Specifically, we frame data contamination detection as a series of multiple-choice questions, devising a quiz format wherein three perturbed versions of each instance, subsampled from a specific dataset partition, are created. These changes only include word-level perturbations. The generated perturbations, along with the original dataset instance, form the options in the DCQ, with an extra option accommodating the selection of none of the provided options. Given that the only distinguishing signal among the options is the *exact wording* with respect to the original dataset instance, an LLM, when tasked with identifying the original dataset instance, gravitates towards selecting the original one if it has been exposed to it. While accounting for positional biases in LLMs, the quiz performance reveals the contamination level for the tested model with the dataset partition to which the quiz pertains. Applied to various datasets and LLMs, under controlled and uncontrolled contamination, our findings—while fully lacking access to training data and model parameters—suggest that DCQ achieves state-of-the-art results and uncovers greater contamination levels through memorization compared to existing methods. Also, it proficiently bypasses more safety filters, especially those set to avoid generating copyrighted content.¹

1 Introduction

Large language models (LLMs) have shown remarkable success in various benchmarks (Bubeck et al., 2023; Bang et al., 2023; OpenAI, 2023, inter alia). However, due to the lack of transparency regarding the data sources used by some

model providers, coupled with their dependency on web data for training (Penedo et al., 2023; Touvron et al., 2023a,b, inter alia), the performance of such models may be skewed by "data contamination" (Zhou et al., 2023). This refers to the inclusion of test data from downstream tasks in the LLMs' training data (Golchin and Surdeanu, 2024). Considering the immense volume of data used for training, even making this data open-access does not provide a practical solution. This is because cross-verifying each sample in the test data against open-source training data becomes impractical due to the substantial computational costs (Touvron et al., 2023b) or the extensive manual effort involved (Chowdhery et al., 2022). This has led LLM developers to assess data contamination in benchmark datasets post-training (Anil et al., 2023a). Therefore, there is a pressing need to accurately pinpoint and quantify this contamination using a cost-effective approach, which demands minimal computational resources and human intervention, ensuring its broad applicability by the community.²

We propose the Data Contamination Quiz (DCQ), a streamlined method that efficiently *detects and estimates verbatim data contamination in LLMs without the need to access training data or internal weights/logits*. Essentially, DCQ is designed to pinpoint the exact data points from dataset partitions that have leaked into the training data of LLMs, leading to skewed performance. In essence, *detection* is achieved via a multiple-choice question and *estimation* is derived by evaluating the LLM's performance across a series of these multiple-choice questions that form the quiz. Specifically, we form a five-option quiz: one option presents an original instance from a dataset partition, three options offer three distinct word-level perturbations of it where words

* Corresponding author.

¹All related code and data are publicly accessible at <https://github.com/shahriargolchin/DCQ>.

²In Appendix A, we show how LLMs are being updated to resist recent data contamination detection methods, emphasizing the need for developing new methods to continually identify data contamination as models evolve.

Original Dataset Instance:

Summary: Premier Inn has announced it is investing £35m in opening three new hotels and revamping two others in Edinburgh.

Instruction: You are provided with a five-choice quiz. Your task is to correctly select the option that exactly corresponds to an instance from the test split of the XSum dataset.

(1) You must ensure that you only generate a single option letter as your answer.

(2) If you do not know the dataset or the correct answer, you must select option "E) None of the provided options."

Hint: While all the following options seem similar in terms of meaning, there is only one option that reflects the exact wording with respect to the original instance.

A) Summary: Premier Inn has revealed it is allocating £35m for the opening of three new hotels and the refurbishment of two others in Edinburgh.

B) Summary: Premier Inn has announced it is investing £35m in opening three new hotels and revamping two others in Edinburgh.

C) Summary: Premier Inn has declared it is investing £35m to set up three new hotels and modernize two others in Edinburgh.

D) Summary: Premier Inn has announced it is spending £35m to open three new hotels and renovate two others in Edinburgh.

E) None of the provided options.

Answer: B ✓

Figure 1: An example of a quiz question crafted to detect data contamination within the test partition of the XSum dataset. Here, the produced answer by the underlying LLM (GPT-4) aligns with the correct option (option C), signaling previous exposure to data, and thus, revealing contamination.

are replaced with their contextually relevant synonyms, and the last option allows the selection of none of these provided options. We use GPT-4 (OpenAI, 2023) to *automate* the generation of word-level perturbations. While the wording differs, the meaning and sentence structure of these perturbations mirror the original dataset instance, resulting in (almost) identical semantic representations across all options containing perturbations and the original dataset instance. Therefore, the sole difference among these options is in the *precise wording* relative to the original dataset instance. Hence, when the LLM is tasked with identifying the original dataset instance, a consistent preference for selecting options containing original dataset instances reveals the model’s memorization and prior exposure to data. Figure 1 exemplifies this idea.

However, the detection of contamination through memorization³ is not consistent across all quiz options due to the *positional biases* present in LLMs (Wang et al., 2023; Zheng et al., 2024), leading to overestimation or underestimation of

³While we recognize *memorization* and *contamination* as distinct concepts, in our study, they are used interchangeably. This is because our method relies on memorization as a basis for detecting contamination.

contamination levels. To counter this, we first identify positional biases using a quiz we call the Bias Detector Quiz (BDQ), which excludes correct answers (i.e., original dataset instances) from its options. In fact, all options in each question of the BDQ are word-level perturbations based on each of the original dataset instances under examination, with an extra option for selecting none of the provided options.⁴ Upon submitting BDQ to an LLM, options that are chosen less frequently than random chance are identified as "non-preferred options," excluding the last option (option E) that always allows for the selection of none of the provided options. *These non-preferred options are strategic positions to avoid overestimation.* Therefore, they are replaced with the original dataset instances in the Bias Compensator Quiz (BCQ) to detect contamination while avoiding overestimation. Nevertheless, due to the same issue of positional biases, memorization may not be triggered equally among all the non-preferred positions, thereby leading to potential underestimation. *To avoid underestimation, we permute the correct answers among all the non-preferred options*, which means an LLM takes multiple BCQs (equal to the number of non-preferred options), to determine the maximum contamination without overestimation. Finally, using the highest quiz performance across all BCQs, a range for the minimum and maximum detected contamination levels is estimated. *We refer to this end-to-end process as our Data Contamination Quiz (DCQ).*

The main contributions of this paper are:

- (1) We propose the first strategy for *estimating* verbatim data contamination in fully black-box LLMs. This detection relies solely on input-output interactions, without the need to access training data and model weights/logits. Our quiz-based approach is structured such that correct answers validate the existence of memorization, thereby revealing prior exposure to data. As a result, contamination level is estimated based on the quiz performance on a subsample of dataset instances.
- (2) Using rigorous evaluations across multiple scenarios, including controlled and uncontrolled contamination environments, our findings show that our method detects significantly higher levels of contamination via memorization than the existing method for identifying data contamination

⁴In BDQ, the question in Figure 1 offers four word-level perturbations, omitting the original dataset instance.

in fully black-box LLMs (Golchin and Surdeanu, 2024).

Overall, our strategy offers *several key benefits over existing methods* for data contamination detection: **(1)** by limiting the LLMs’ outputs to be a single letter that identifies the selected option (i.e., from A to E), our strategy substantially reduces the probabilistic behavior of LLMs in detecting data contamination; **(2)** building on the same principle of limiting outputs, our method can adeptly evade the safety filters set by model providers during the decoding process, especially in proprietary LLMs that seek to prevent the generation of copyrighted content, thus exposing more contamination; **(3)** it transcends the limitations of probability/threshold-based methods (Shi et al., 2024; Oren et al., 2024; Dong et al., 2024, *inter alia*), which are task/language-specific, allowing our approach to be applied across different tasks and languages; **(4)** our technique is both cost- and time-saving thanks to the automated quiz generation process and administration on a manageable set of instances, such as 100 subsampled instances used in our experiments; and **(5)** the streamlined nature of our method means the quiz options we collect for each dataset partition are reusable and can be applied to any LLM to take the DCQ for data contamination assessment.

2 Approach

Our approach is based on a key assumption: *if an LLM has encountered a dataset instance verbatim during training, then it can recognize that instance among its word-level perturbed versions through memorization* (Carlini et al., 2023). Framing this as a quiz-like task, prior exposure to data is inferred by tasking the model to select a single-letter option corresponding to the original dataset instance.

Given this assumption, our approach consists of three main phases:

(1) Creating Quiz Options: We generate four unique word-level perturbations for each subsampled dataset instance.

(2) Detecting and Handling Positional Biases: We craft and submit two types of quizzes using two combinations of the generated perturbations:

- **Bias Detector Quiz (BDQ):** This quiz includes all four perturbations as options, along with an option to select none of them. Its purpose is to identify positional biases.

- **Bias Compensator Quiz (BCQ):** This quiz replaces one of the four perturbations in the previous quiz with the original dataset instance and aims to assess contamination levels while accounting for positional biases.

(3) Estimating Contamination Levels: Based on quiz performance from BCQ, we estimate a range for the detected contamination levels.

We detail each phase in the following.

2.1 Creating Quiz Options

We use GPT-4 to generate four distinct word-level perturbations for each subsampled dataset instance. Its rich representations and proficient adherence to following instructions (Bubeck et al., 2023) make it suitable for performing word replacements that maintain contextual relevance. This reduces the need for prompt engineering and mitigates potential biases when producing word-level perturbations that can arise from factors such as variations in length, sentence structure, or the lack of diversity.

To instruct GPT-4, we employ a zero-shot prompt with the original dataset instance infused in it as a reference.⁵ In particular, this prompt involves the following rules and requirements when generating four quiz options per instance: (1) all generated perturbations must retain the meaning and sentence structure of the original dataset instance; (2) they must be distinct from one another; (3) they must conform to every precise symbol and letter detail found in the original dataset instance; and (4) the sole difference between them must be the word-level perturbations.⁶ These options are then verified to be distinct from one another and not identical to the original dataset instance. Further, all options represent random perturbations in terms of position and the number of tokens involved in the perturbation process. However, when the original dataset instance contains a label, we do not alter the original label and directly attach it to the perturbed versions.⁷

We then compile and employ *all four generated options* in BDQ to identify positional biases. In contrast, we use *three of these options* when administering BCQ, wherein one option from BDQ is replaced with the original dataset instance. As

⁵This prompt is shown in Figure 4 in Appendix B.

⁶Refer to Appendix B for more details and analysis on the generated perturbations.

⁷See Table 5 in Appendix A for an example.

a fifth option, in both quizzes, we always include a choice that allows for selecting none of the provided options. This is beneficial in situations where the LLM is not contaminated with a particular dataset or when the LLM’s memorization is not triggered by the presented options. We detail each quiz in Subsection 2.2. Figure 1 shows the integration of the generated perturbations for an original dataset instance as quiz options in BCQ.

The decision to design quizzes with five options is based on the fact that most multiple-choice quizzes typically include three to five answer choices per question. LLMs are frequently exposed to and familiarized with these formats during their fine-tuning and/or pre-training stages. Therefore, using this range of options aligns with the standard quiz formats LLMs previously encountered, minimizing the chance for poor or exceptional performance only due to an unfamiliar test design. Additionally, choosing five choices within this standard range is intended to create a robust and challenging test environment by maximizing the number of answer choices.

2.2 Detecting and Handling Positional Biases

The performance of LLMs in answering multiple-choice questions is influenced by positional biases (Wang et al., 2023; Zheng et al., 2024), which we empirically show in Subsection 4.1. These biases alter the quiz performance according to the placement of the options, leading to varying detectable contamination levels for the same dataset instances being studied. Therefore, we first detect positional biases, and then, adjust for these biases in our contamination detection process.

Bias Detector Quiz (BDQ). This five-option quiz is primarily designed to detect positional biases. In this quiz, the correct answers (original dataset instances) are absent, and all options—except for the last one, which is an option to select none of the provided options—represent word-level perturbations of the original dataset instances. By administering this quiz to an LLM, we analyze the LLM’s tendency to favor certain options (positions) more frequently than others when the correct answers are missing, thereby revealing positional biases. *From the outcomes of the BDQ, we identify what we term the “non-preferred options.”*⁸ We define an option as

⁸While having zero or one non-preferred option is possible, we consistently refer to *multiple non-preferred options*

non-preferred in consideration of an LLM if it is chosen less frequently than what random chance suggests. This definition excludes the last option (option E), which is always fixed to enable the selection of none of the given options. For a five-option quiz with k questions (instances being evaluated), any option selected fewer than $\lceil \frac{k}{5} \rceil$ times in the BDQ is deemed non-preferred. These options are viewed as strategic spots to minimize the risk of overestimating contamination when placing the original dataset instances there. This strategic placement is the main task of what we discuss next, the BCQ.

Bias Compensator Quiz (BCQ). The main objective of this quiz is to compensate for positional biases when detecting data contamination. Such biases can skew quiz performance, causing the perceived contamination levels based on quiz performance to be overestimated or underestimated. *To avoid overestimation, we systematically replace the detected non-preferred options by BDQ with the original dataset instances.* Nevertheless, due to the same issue of positional biases, the detection of the original dataset instances is not consistent among all the non-preferred options, which leads to an underestimation of contamination. *To avoid underestimation and find the highest possible memorization among non-preferred options, we permute the correct answers, i.e., options containing the original dataset instances, among all non-preferred options.*⁹ This allows us to obtain the highest performance in triggering the LLM’s memorization for detecting the original dataset instances while avoiding overestimation. Finally, the highest performance among all BCQs is reported as the maximum level of contamination.¹⁰

2.3 Estimating Contamination Levels by Quiz Performance

As noted, the maximum contamination level for a dataset partition is determined by the highest performance of the LLM across all BCQs. While this provides the maximum level of detectable

due to their common occurrence and to ensure consistency in our discussion in the paper.

⁹While we never faced this situation, in the rare event that there are no non-preferred options, we permute the correct answers among all options except the last one (option E).

¹⁰When BCQs tie for the highest quiz performance, we select the one with the lowest positional bias detected by BDQ.

contamination, it is also possible to compute the minimum level of detectable contamination.

As we deal with a non-ideal system, we can apply the principle of chance-adjusted accuracy to account for the probability of chance agreement using Cohen’s Kappa (Galton, 1892; Cohen, 1960), defined as: $\kappa = (p_o - p_e)/(1 - p_e)$. Here, p_o is the observed agreement and p_e is the expected agreement by chance. In simpler terms, p_o refers to the highest proportion of correct answers attained by the LLM among BCQs, whereas p_e is the proportion of choosing the non-preferred option in the BDQ, where the LLM achieved its top performance in the BCQ by placing the correct answers (original dataset instances) in that position.

For example, if the best performance among all BCQs is obtained by systematically placing correct answers in position B, p_o is the proportion of selecting option B in BCQ while p_e is the proportion of selecting option B in BDQ. The resulting value is then reported as the minimum level of detected contamination.

3 Experimental Setup

Data. To improve clarity given the various evaluation settings explored in this paper, we provide the specifics of the datasets used for each evaluation and the relevant subsection separately below.

Subsection 4.1: To examine the initiation of memorization in LLMs via quiz options and controlling overestimation and underestimation due to positional biases, we use two of the contaminated datasets reported by Golchin and Surdeanu (2024): AG News (Zhang et al., 2015) and WNLI (Wang et al., 2018). Specifically, we use the train and test sets of the former and the train and validation sets of the latter dataset.

Subsection 4.2: In our systematic contamination experiment, we utilize the train sets from two uncontaminated datasets—MeetingBank (Hu et al., 2023) and AuTextification (Sarvazyan et al., 2023)—both released in 2023, and thus, are unexposed to GPT-3.5 (Ouyang et al., 2022) and Llama 2 (Touvron et al., 2023a). Also, our experiment involving reported contamination uses test split from the HumanEval (Chen et al., 2021), validation split from the DROP (Dua et al., 2019), and train split from the GSM8k (Cobbe et al., 2021).

Subsection 4.3: To apply our approach in the wild, we use the same datasets as used by Golchin and Surdeanu (2024) for their data contamination analysis. This includes IMDB (Maas et al., 2011), AG News (Zhang et al., 2015), Yelp Full Reviews (Zhang et al., 2015), SAMSum (Gliwa et al., 2019), XSum (Narayan et al., 2018), WNLI (Wang et al., 2018), and RTE (Wang et al., 2019). Our experiments focus on the test splits of these datasets except for WNLI and RTE, where we use validation sets due to the unavailability of their test sets.

Preparation of instances involves formatting each instance as per the task of its dataset: For the classification task, an instance includes the text along with its exact label; for the natural language inference task, it contains the premise, hypothesis, and the exact label; for summarization task, summaries without related documents; and so on. In general, we format each instance such that, if disclosed during the training stage, it can skew the model’s performance and generalizability.

Following prior studies (Shi et al., 2024; Golchin and Surdeanu, 2024), we subsample $k = 100$ unique dataset instances from each dataset partition for all experiments. If a partition contains fewer instances, we use the entire set. In Subsection 4.2, we demonstrate that this sample size is adequate for effectively estimating contamination levels while managing the costs associated with perturbation generation and quiz administration.

Setting. All experiments involving GPT-3.5 and GPT-4 use the `gpt-3.5-turbo-0613` and `gpt-4-0613` snapshots, accessed via the OpenAI API. For all our experiments with Llama 2, we employ the `Llama 2-Chat (70B)` model.

When generating word-level perturbations, we adjust the temperature for GPT-4 to 1.00—the default value for generative tasks such as conversation and storytelling¹¹—to encourage diversity, and cap the completion length to 4,000 tokens. In contrast, when submitting all quizzes (i.e., BDQ and BCQ), we set the temperature to zero—promoting deterministic output given the provided options—with a token limit of one.

Comparative Framework. We compare our DCQ against the replication-based method proposed by Golchin and Surdeanu (2024). This choice is due to several key reasons: (1) their detection

¹¹<https://platform.openai.com/playground>.

method is capable of operating in a fully black-box setting, without access to training data and model weights/logits, which aligns well with our requirements; (2) although this method was originally designed to detect contamination at the partition level, it is the only existing method that can be adapted beyond binary detection to quantify contamination by calculating the proportion of exact and near-exact matches replicated by this method; and (3) this strategy is both task- and language-agnostic, making it broadly applicable to our experiments.

Other existing detection methods do not fulfill all the aforementioned criteria. Specifically, some require access to model weights/logits or training data (Shi et al., 2024; Oren et al., 2024; Yang et al., 2023; Deng et al., 2023, inter alia), others cannot be adapted for estimation comparison while operating in fully black-box setting (Dong et al., 2024; Deng et al., 2023), and some are task- or language-specific (Dong et al., 2024; Roberts et al., 2023; Shi et al., 2024; Oren et al., 2024; Deng et al., 2023, inter alia). For example, the method proposed by Dong et al. (2024) is ineffective for detecting contamination in short- or non-generative tasks such as classification and summarization. This approach relies on diversity to identify contamination, which is only informative when there is sufficient variation in the generated outputs. In cases where diversity is inherently low, the method incorrectly flags nearly all samples as contaminated, leading to inconclusive findings.

In particular, we employ the top-performing strategy (guided instruction with GPT-4 evaluation) proposed by Golchin and Surdeanu (2024) for detecting contamination, treating both exact and near-exact matches as contamination instances alike. In our experiments using this method, we set the temperature to zero and cap the token limit at 500. We also create random-length initial segments of dataset instances by randomly deriving them from a range spanning 40% to 70% of each instance’s length, based on the white space count. We selected this range as it yielded the highest levels of memorization when replicating instances using the LLMs under investigation.

4 Empirical Results and Evaluation

In this section, we first provide empirical evidence on how our approach instigates memorization in

LLMs for revealing data exposure while avoiding overestimation and underestimation of contamination. Next, we assess its performance in detecting verbatim contamination by applying it to controlled contamination settings where the presence of contamination is either systematic or reported, and further employ it in real-world cases where contamination is not controlled.

4.1 Positional Adversarial Analysis

Our quiz-based approach uses three main tactics: (1) integrating original dataset instances with their word-level perturbations as quiz options exposes memorization in LLMs; (2) strategically placing these original dataset instances in non-preferred options avoids the overestimation of contamination; and (3) permuting these original dataset instances among non-preferred options avoids the underestimation of contamination. Below, we provide empirical results that corroborate each tactic.

Exposing Memorization via Quiz Options.

For this experiment, after identifying positional biases by BDQ, we undertake two variants of BCQ. These variants are designed to illustrate how positional biases are affected by the presence and absence of memorization. Unlike the standard setting described in Subsection 2.2, these variants do not involve permuting correct answers (original dataset instances) among the non-preferred options. *Instead, we statically position them in a single detected least-chosen option.* In other words, this least-chosen option is a non-preferred option that has the lowest selection frequency among others.

In the first variant of BCQ, we replace the least-chosen option with the original dataset instances. However, in the second variant, the least-chosen option is replaced with unique perturbations, ensuring they differ from all four perturbations used in BDQ. This confirms that if an LLM has not been exposed to the original dataset instances, and thus, cannot recognize them by memorization, the selection frequency distribution among options in BDQ should be maintained regardless of whether the original dataset instances are included among the options. In contrast, a pronounced preference for selecting the option containing the original dataset instances signals memorization and prior exposure to data, unveiling data contamination.

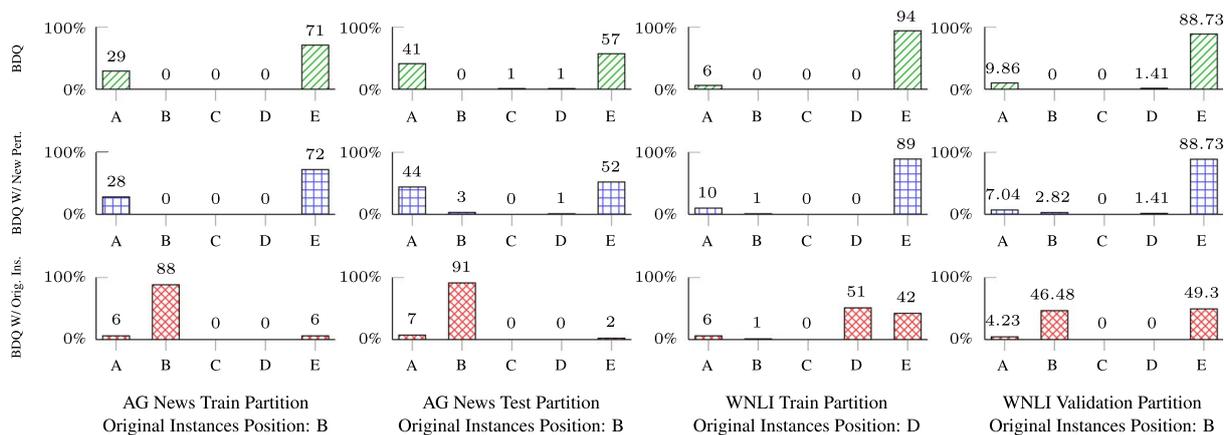


Figure 2: Results of positional adversarial analysis. The green bar charts show the selection frequency distribution of options in BDQ while the blue and red ones represent results after replacing new word-level perturbations and original dataset instances with the least-chosen option detected by BDQ, respectively. The increase in selection frequency for options containing original dataset instances (red bar charts) signals the model’s prior exposure to data. In all settings, GPT-4 is the base model.

To analyze the changes in the selection frequency distribution of options, we utilize four contaminated dataset partitions, reported by Golchin and Surdeanu (2024). Specifically, we use the train and test sets of AG News as well as the train and validation sets of the WNLi dataset, with GPT-4 as the base model. Figure 2 shows the outcomes of this experiment. As shown, when new perturbations are replaced with the least-chosen options (blue bar plots), as the model has not seen them, the selection frequency distribution of options remains significantly consistent with BDQ (green bar plots). Conversely, there is a significant preference for selecting options containing the original dataset instances when they are replaced with the least-chosen options (red bar plots).¹²

For example, the initial selection frequency distribution of options in BDQ for the AG News train set is {A: 29, B: 0, C: 0, D: 0, E: 71}. Thus, options B, C, and D are deemed non-preferred due to being chosen less frequently than random chance (20 times out of 100 questions), and they are also the least-chosen options, with zero selections each.¹³ When there are multiple least-chosen options, we opt for the one that results in the highest selection frequency when replaced with new perturbations. If this leads to identical selection

frequencies among the least-chosen options, we then opt for the position that results in the highest selection frequency when replaced with original dataset instances. This helps to better visualize the increase in selection frequency while establishing a rigorous analysis by favoring the setting containing new perturbations first. Thus, in this example, we select option B in BDQ to conduct our replacements, as it resulted in the highest selection frequency when replaced with the original dataset instances. When option B is replaced with new perturbations, no preference for selecting option B is observed, and the initial positional biases in BDQ remain (almost) unchanged: {A: 28, B: 0, C: 0, D: 0, E: 72}. However, when option B is replaced with the original dataset instances, a noticeable rise in its selection becomes apparent: {A: 6, B: 88, C: 0, D: 0, E: 6}. Therefore, we showed that infusing original dataset instances among their word-level perturbations as quiz options triggers evidence of memorization in LLMs.

Avoiding Overestimation Using Non-Preferred Options. Positioning the original dataset instances among quiz options without considering positional biases fails to accurately measure contamination levels. Considering the previous example, if the original dataset instances are replaced with option A in BDQ, any contamination level surpassing 71% results in the system indicating 100% contamination. While it is not feasible to unequivocally guarantee the probabilistic behavior of LLMs, including the potential

¹²We assessed this significance using Fisher’s Exact Test (Fisher, 1922), with p -values < 0.05 .

¹³Note that *non-preferred options are not fixed for all models and datasets*, and must be identified for each specific dataset and model individually.

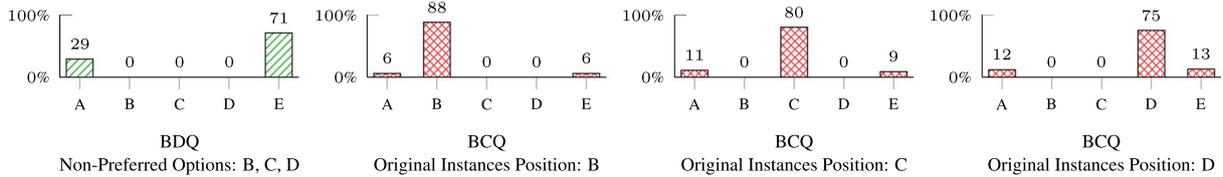


Figure 3: An illustration of the full process for our approach, performed on the AG News train set. Initially, BDQ is executed to identify non-preferred options (i.e., B, C, and D), followed by multiple BCQs to determine the maximum contamination (88%) through permutation among the non-preferred options. In all settings, GPT-4 is the base model.

overestimation of contamination in our case, it can be confidently asserted that our approach does not stimulate this overestimation at least by design. In fact, this is achieved by using non-preferred options as the replacement spots for the original dataset instances.

Avoiding Underestimation by Permutation.

As discussed, non-preferred options serve as safe spots for positioning original dataset instances to mitigate the risk of overestimation in the presence of positional biases. However, these biases can also result in uneven memorization triggered across the detected non-preferred options. Figure 3 provides empirical evidence of this phenomenon when using the AG News train partition as our case study. As shown, when the original dataset instances are placed in the non-preferred options (i.e., B, C, and D), the level of contamination revealed by memorization varies within a single setting: 88%, 80%, and 75% for positions B, C, and D, respectively. As detailed in Subsection 2.2, to avoid underestimation, the highest quiz performance among all BCQs is reported as the maximum level of detected contamination, which is 88% in this example. Thus, we empirically showed that permuting original dataset instances among the non-preferred options is essential for uncovering the highest possible contamination and avoiding underestimation.

4.2 Evaluation under Controlled Contamination

Systematic Contamination. As for the first step in assessing the performance of our approach, we start with systematic contamination. In this experiment, we systematically contaminate GPT-3.5 and Llama 2 with two uncontaminated datasets based on their training cutoff: MeetingBank and AuTextification.¹⁴

¹⁴See Appendix C for details on systematic contamination of GPT-3.5 and Llama 2.

Our evaluation involves estimating contamination using our DCQ approach for three contamination levels per dataset: 100%, 50%, and 0%. Following prior work (Shi et al., 2024; Oren et al., 2024), we perform systematic contamination using 1,000 instances from each dataset separately.

For the 100 instances used in our method, the systematic contamination settings vary as follows:

- **100% Contamination:** All 100 instances are selected from the 1,000 contaminated instances.
- **50% Contamination:** 50 instances are sampled from the 1,000 contaminated instances, while the other 50 are drawn from the source dataset, excluding the 1,000 contaminated instances.
- **0% Contamination:** All 100 instances are sampled from the source dataset, excluding the 1,000 contaminated instances.

This experimental design models real-world scenarios where DCQ estimates contamination using subsampled dataset instances.

Table 1 presents the results of our experiment under systematic contamination upon applying DCQ.¹⁵ Compared with the replication-based strategy (Golchin and Surdeanu, 2024), DCQ proves to be significantly more effective at revealing contamination. In fact, contamination levels identified by replicating instances are so negligible that it is impracticable to compare with our DCQ’s performance. As noted in Subsection 4.1, this improved ability to detect more memorization is obtained by minimizing LLMs’ probabilistic outputs to a single-letter option that hints at the presence of memorization. Another benefit of it

¹⁵We provide statistics on the positional biases of all the models used in this study across the corresponding datasets in Table 6 in Appendix D.

Model	Method	MeetingBank			AuTextTification		
		100% Cont.	50% Cont.	0% Cont.	100% Cont.	50% Cont.	0% Cont.
GPT-3.5	Rplc-based (%)	1.00	0.00	0.00	1.00	0.00	0.00
	Our DCQ (%)	[85.87, 87.00]	[46.31, 49.00]	[0.00, 3.00]	[66.29, 70.00]	[44.68, 48.00]	[0.00, 5.00]
	Recall (%)	87.00	88.00	Undefined	70.00	70.00	Undefined
	Precision (%)	87.00	89.80	0.00	70.00	72.92	0.00
Llama 2	Rplc-based (%)	1.00	0.00	0.00	0.00	0.00	0.00
	Our DCQ (%)	[84.21, 85.00]	[45.83, 48.00]	[0.00, 2.00]	[72.63, 74.00]	[44.21, 47.00]	[0.00, 2.00]
	Recall (%)	85.00	84.00	Undefined	74.00	74.00	Undefined
	Precision (%)	85.00	87.50	0.00	74.00	78.72	0.00

Table 1: Results from our DCQ after introducing systematic contamination to GPT-3.5 and Llama 2. The table lists the minimum and maximum detected contamination levels. Results are also compared with the replication-based (rplc-based) method by Golchin and Surdeanu (2024). Note that, in settings with no contamination, Recall cannot be calculated, as both True Positives and False Negatives are zero.

is its ability to bypass the preemptive filters set to stop LLMs from producing copyrighted content.¹⁶ Importantly, DCQ not only reveals more memorization but also effectively measures contamination levels. As detailed in Subsection 4.1, this effective estimation is derived by placing and permuting the original dataset instances among the non-preferred options, simultaneously avoiding overestimation and underestimation.

Reported Contamination. The GPT-4 technical report provides cross-contamination between several datasets and the training data of this model (OpenAI, 2023). In this experiment, we consider the backdrop of reported contamination as a baseline against which we compare the performance of DCQ in detecting contamination. We use three datasets that hold significant importance in the context of LLMs: HumanEval,¹⁷ DROP, and GSM8k, and apply DCQ to their test, validation, and train splits, respectively. The rationale behind selecting these particular datasets stems from their pivotal role in assessing the reasoning/problem-solving capabilities of LLMs, along with the explicit mention of contamination levels for them by OpenAI (2023). Further, our focus on the GSM8k train partition is driven by the fact that this set—apart from a small portion—was included in the training data of GPT-4, though without a specified level of inclusion. Thus, we expect that our DCQ detects a significant level of contamination for it.

¹⁶Refer to Appendix A to see case examples.

¹⁷As done by OpenAI (2023), we apply DCQ to the entire test set of the HumanEval dataset with 164 samples.

Method	HumanEval	DROP	GSM8k
Reported (%)	≈ 25.00	≈ 21.00	Full but a part
Rplc-based (%)	0.00	4.00	2.00
Our DCQ (%)	[55.62, 56.71]	[42.86, 44.00]	[78.79, 79.00]

Table 2: A Comparison between the contamination levels reported for GPT-4 (OpenAI, 2023) and those identified by the replication-based (rplc-based) method (Golchin and Surdeanu, 2024) and our DCQ. This experiment encompasses the test, validation, and train partitions of the HumanEval, DROP, and GSM8k datasets, respectively.

Table 2 compares the contamination levels detected by DCQ with those reported by OpenAI (2023) for the aforesaid datasets. As mentioned in Subsection 2.2, given that DCQ inherently does not overestimate/underestimate contamination, the results indicate its ability to uncover higher contamination levels, even when compared to those measured in the availability of the training data. Further, as expected, the contamination level identified for the GSM8k train set is markedly high while the replication-based method can only detect 2.00% of it.

Note that the cross-contamination levels reported by OpenAI (2023) in Table 2 may not offer a thorough depiction of the levels detected by DCQ for two reasons: (1) The measured cross-contamination comes from overlapping three 50-character substrings, yielding high false negatives/positives (OpenAI, 2023). As a result, the reported contamination levels serve as approximations rather than definitive measures. (2) Cross-contamination levels were estimated using a

Model	Method	IMDB	AG News	Yelp	RTE	WNLI	SAMSum	XSum
GPT-4	Rplc-based (%)	0.00	3.00	0.00	1.00	15.49	0.00	7.00
	Our DCQ (%)	[81.82, 82.00]	[91.00, 91.00]	[80.00, 80.00]	[60.00, 60.00]	[45.31, 50.70]	[76.04, 77.00]	[94.90, 95.00]
GPT-3.5	Rplc-based (%)	0.00	0.00	0.00	1.00	2.82	0.00	2.00
	Our DCQ (%)	[46.43, 55.00]	[80.43, 82.00]	[12.12, 13.00]	[65.06, 71.00]	[7.47, 12.68]	[70.11, 74.00]	[78.35, 79.00]
Llama 2	Rplc-based (%)	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	Our DCQ (%)	[4.44, 14.00]	[17.20, 23.00]	[15.15, 16.00]	[30.23, 40.00]	[11.29, 22.54]	[0.00, 1.00]	[11.00, 11.00]

Table 3: Estimated data contamination levels for GPT-4, GPT-3.5, and Llama 2 using our DCQ across test/validation partitions of seven datasets, based on a maximum of 100 subsampled instances. Results are also compared with the replication-based (rplc-based) method by Golchin and Surdeanu (2024).

maximum of 1,000 subsampled instances from each set, rather than the entire set (OpenAI, 2023).

4.3 Evaluation in the Wild

After examining contamination in controlled environments, we apply our detection strategy to various datasets in real-world scenarios where the exact contamination levels are unknown. To mitigate the lack of information about the precise contamination levels for the datasets we analyze, we assess the contamination levels detected by our approach against the replication-based method by Golchin and Surdeanu (2024). This evaluation is informed by effective estimates obtained under controlled contamination in Subsection 4.2 and by understanding that our method does not encourage both overestimation and underestimation. For this experiment, we use the same datasets as in their study, covering test/validation splits of seven datasets with GPT-4, GPT-3.5, and Llama 2.

Table 3 lists all the results from our evaluation in the wild over 21 settings. As outcomes indicate, the replication-based method (Golchin and Surdeanu, 2024) can replicate a few dataset instances, such as when GPT-4 serves as the base model and is applied to WNLI (15.49%) and XSum (7.00%) datasets. In comparison, our DCQ detects significant amounts of contamination across all studied dataset partitions. For example, the maximum contamination identified by our DCQ in the aforesaid datasets reaches 50.70% and 95.00%, representing a detection range of contamination 3 to over 13 times greater than that of the replication-based method. Additionally, this ratio is even more substantial for other studied datasets.

Overall, our results indicate that the replication-based method is effective only for binary detection of contamination when the memorization level is high, and lacks the ability to effectively estimate contamination. In contrast, our DCQ reveals cases of memorization/contamination at

levels significantly greater than methods based on replicating/extracting training data (Golchin and Surdeanu, 2024; Carlini et al., 2023).

5 Related Work

Data contamination is shown to inflate downstream performance (Zhou et al., 2023; Palavalli, 2024; Jiang et al., 2024; Dong et al., 2024; Deng et al., 2024; Balloccu et al., 2024; Li and Flanigan, 2024). Identifying data contamination in LLMs is not difficult when the training data is accessible, though it can be resource-intensive. Typically, this process involves expensive overlap comparisons between the training data and the test data using high-order n -grams—a method exclusive to model developers that yields approximate estimates of contamination rather than definitive ones (OpenAI, 2023; Touvron et al., 2023a,b; Anil et al., 2023b; Chowdhery et al., 2022; Brown et al., 2020, *inter alia*). However, this task becomes challenging when the training data is absent due to being proprietary or when it is open but vast in size.

Roberts et al. (2023) used a time-based approach for contamination detection by analyzing dataset pass rates pre- and post-model training cutoff. However, this method has limitations due to the continuous updates of proprietary LLMs and the need for ongoing data collection. Further, detecting contamination for non-code/math datasets remained unresolved. Shi et al. (2024) employed a similar strategy, measuring outlier word likelihood in new examples to determine if a text was part of the training corpus. While only useful for LLMs with accessible logits, detecting training data is not synonymous with confirming contamination. Zhang et al. (2024) improved this method using the idea of local maxima. In line with probability-based methods, Oren et al. (2024) aimed to detect data contamination by comparing

the likelihoods of a canonically ordered dataset and its shuffled version. However, the positional biases in LLMs (Wang et al., 2023; Zheng et al., 2024) and the bold presumption of dataset instances appearing in the training data in the same order as on the web limit the applicability of this method. Yang et al. (2023) further relaxed their settings by assuming access to training data to detect contamination. They used embedding similarity to find the top-k similar training samples for each dataset instance and used an LLM to verify if these samples were rephrased versions of the dataset instance. Similarly, Deng et al. (2023) developed a retrieval-based method that identified potential contamination by querying training corpora and assessing overlap with benchmark data. However, assuming access to training data is impractical, even for open-weight models.

A few recent studies developed strategies for detecting contamination without access to training data and model weights/logits. Sainz et al. (2023a,b) prompted GPT-3.5 to generate the first instances of a dataset partition. However, this method faltered due to the sparseness of the signal in the prompt, coupled with the LLM’s safety filters preventing the generation of copyrighted materials. In contrast, Golchin and Surdeanu (2024) amplified the signal in the prompt by infusing a random-length initial segment of a dataset instance and directing the LLM to complete it. Dong et al. (2024) detected contamination by measuring the peakedness of the model’s output distribution through edit distance to reveal contamination. Deng et al. (2023) proposed two methods: a mask-based approach, which detects contamination if the model can predict masked parts of dataset instances, and a question-based approach, which identifies contamination if the model generates a new answer to a question that already has an answer embedded in the prompt. However, none of these methods can estimate contamination.

6 Conclusion

We proposed the Data Contamination Quiz (DCQ), the first method to *estimate* verbatim data contamination in fully black-box LLMs by framing it as a series of multiple-choice questions, forming our quiz. This quiz instructs an LLM to select an option containing the original dataset instance among its word-level perturbations. If the

LLM succeeds in detecting the original dataset instance, it reveals prior exposure to data. Thus, the quiz performance indicates the level of detected contamination for the LLM that took DCQ. Using several datasets and LLMs, under controlled and uncontrolled settings, our findings demonstrated that the detected contamination levels by our DCQ not only align with previously reported contaminated datasets but also surpass the existing methods by providing an estimation of verbatim data contamination, moving beyond binary detection.

Limitations

Data contamination primarily occurs in two forms: *verbatim*, where dataset instances are directly included in the training data, and *non-verbatim*, such as metadata contamination. This study only focuses on estimating verbatim contamination in LLMs. Therefore, we encourage future research to develop methods for detecting and estimating other types of contamination beyond verbatim.

We also emphasize the importance of carefully considering data contamination when evaluating LLMs on benchmarks, particularly for reasoning tasks, where the goal is to assess genuine generalization abilities. Although there has been progress in reducing contamination risks through private or semi-private evaluations of newer models, recent findings suggest that contamination remains a serious concern, and private evaluations alone are neither a definitive nor a long-term solution.

For example, OpenAI’s recent o3 model (Jaech et al., 2024; OpenAI, 2024) achieved near-human performance on the ARC-AGI-1 dataset (Chollet, 2019), representing a major leap in reasoning capabilities. However, it scored only in the single digits on the newly released ARC-AGI-2 dataset (Chollet et al., 2024), even though both datasets are of similar difficulty and considered relatively easy for humans, according to their creators (ARC Prize Team, 2024). Since the only primary difference between the two datasets is their release time, this large performance gap raises concerns about the genuine generalization abilities of these models and the reliability of private evaluations in preventing contamination, as this longitudinal performance drop can be mostly attributed to contamination (Roberts et al., 2023). As a result, data contamination remains a major challenge in model evaluation and calls for further research.

Acknowledgments

We sincerely thank Dr. Steven Bethard and Dr. Eduardo Blanco for their expert guidance on the early draft of this paper. We also appreciate the action editor and anonymous reviewers for their constructive inputs that greatly improved this work.

References

- Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul R. Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Jack Krawczyk, Cosmo Du, Ed Chi, Heng-Tze Cheng, Eric Ni, Purvi Shah, Patrick Kane, Betty Chan, Manaal Faruqui, Aliaksei Severyn, Hanzhao Lin, YaGuang Li, Yong Cheng, Abe Ittycheriah, Mahdis Mahdieh, Mia Chen, Pei Sun, Dustin Tran, Sumit Bagri, Balaji Lakshminarayanan, Jeremiah Liu, Andras Orban, Fabian Gra, Hao Zhou, Xinying Song, Aurelien Boffy, Harish Ganapathy, Steven Zheng, HyunJeong Choe, goston Weisz, Tao Zhu, Yifeng Lu, Siddharth Gopal, Jarrod Kahn, Maciej Kula, Jeff Pitman, Rushin Shah, Emanuel Taropa, Majd Al Merey, Martin Baeuml, Zhifeng Chen, Laurent El Shafey, Yujing Zhang, Olcan Sercinoglu, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anas White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakob Sygnowski, Alexandre Frechette, Charlotte Smith, Laura Culp, Lev Proleev, Yi Luan, Xi Chen et al. 2023a. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clment Crepy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Daz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li, Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu, Frederick Liu, Marcello Maggioni, Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Castro Ros, Aurko Roy, Brennan Saeta et al. 2023b. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.
- ARC Prize Team. 2024. Announcing arc-agi-2 and the arc prize 2025. Accessed: 2025-04-25.
- Simone Balloccu, Patrcia Schmidtov, Mateusz Lango, and Ondřej Duřek. 2024. Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source LLMs. *arXiv preprint arXiv:2402.03927*.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*. <https://doi.org/10.18653/v1/2023.ijcnlp-main.45>

- 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. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrike, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. 2023. Quantifying memorization across neural language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1–5, 2023*. OpenReview.net.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code.
- François Chollet. 2019. On the measure of intelligence. In *arXiv preprint arXiv:1911.01547*.
- Francois Chollet, Mike Knoop, Gregory Kamradt, and Bryan Landers. 2024. Arc prize 2024: Technical report. *arXiv preprint arXiv:2412.04604*.
- 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, 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. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.
- 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. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20:37–46. <https://doi.org/10.1177/001316446002000104>
- Chunyu Deng, Yilun Zhao, Yuzhao Heng, Yitong Li, Jiannan Cao, Xiangru Tang, and

- Arman Cohan. 2024. Unveiling the spectrum of data contamination in language models: A survey from detection to remediation. *arXiv preprint arXiv:2406.14644*. <https://doi.org/10.18653/v1/2024.findings-acl.951>
- Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gerstein, and Arman Cohan. 2023. Investigating data contamination in modern benchmarks for large language models. *arXiv preprint arXiv:2311.09783*.
- Yihong Dong, Xue Jiang, Huanyu Liu, Zhi Jin, and Ge Li. 2024. Generalization or memorization: Data contamination and trustworthy evaluation for large language models. *arXiv preprint arXiv:2402.15938*. <https://doi.org/10.18653/v1/2024.findings-acl.716>
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2–7, 2019, Volume 1 (Long and Short Papers)*, pages 2368–2378. Association for Computational Linguistics.
- B. Efron. 1979. Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, 7(1):1–26. <https://doi.org/10.1214/aos/1176344552>
- Bradley Efron. 2003. Second thoughts on the bootstrap. *Statistical Science*, 18(2):135–140. <https://doi.org/10.1214/ss/1063994968>
- Bradley Efron and Robert J. Tibshirani. 1993. *An Introduction to the Bootstrap*. Number 57 in Monographs on Statistics and Applied Probability. Chapman & Hall/CRC, Boca Raton, Florida, USA. <https://doi.org/10.1007/978-1-4899-4541-9>
- Ronald A. Fisher. 1922. On the interpretation of χ^2 from contingency tables, and the calculation of p. *Journal of the Royal Statistical Society*, 85(1):87–94. <https://doi.org/10.2307/2340521>
- Joseph L. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76(5):378. <https://doi.org/10.1037/h0031619>
- Francis Galton. 1892. *Finger prints*. 57490–57492. Macmillan and Company.
- Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. Samsun corpus: A human-annotated dialogue dataset for abstractive summarization. *arXiv preprint arXiv:1911.12237*. <https://doi.org/10.18653/v1/D19-5409>
- Shahriar Golchin and Mihai Surdeanu. 2024. Time travel in LLMs: Tracing data contamination in large language models. In *The Twelfth International Conference on Learning Representations*.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Yebowen Hu, Tim Ganter, Hanieh Deilamsalehy, Franck Dernoncourt, Hassan Foroosh, and Fei Liu. 2023. Meetingbank: A benchmark dataset for meeting summarization. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL)*, Toronto, Canada. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.acl-long.906>
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. 2024. OpenAI o1 system card. *arXiv preprint arXiv:2412.16720*.
- Minhao Jiang, Ken Ziyu Liu, Ming Zhong, Rylan Schaeffer, Siru Ouyang, Jiawei Han, and Sanmi Koyejo. 2024. Investigating data contamination for pre-training language models. *arXiv preprint arXiv:2401.06059*.
- Changmao Li and Jeffrey Flanigan. 2024. Task contamination: Language models may not be few-shot anymore. In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20–27, 2024, Vancouver, Canada*,

- pages 18471–18480. AAAI Press. <https://doi.org/10.1609/aaai.v38i16.29808>
- Llama Team. 2024. How to guides: Fine tuning. Accessed: 2024-12-26.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don’t give me the details, just the summary! Topic-aware convolutional neural networks for extreme summarization. *ArXiv*, abs/1808.08745. <https://doi.org/10.18653/v1/D18-1206>
- OpenAI. 2023. Gpt-4 technical report.
- OpenAI. 2024. O3 mini system card. Accessed: 2025-04-25.
- Yonatan Oren, Nicole Meister, Niladri S. Chatterji, Faisal Ladhak, and Tatsunori Hashimoto. 2024. Proving test set contamination in black-box language models. In *The Twelfth International Conference on Learning Representations*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll 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. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Medha Palavalli. 2024. *Taxonomy for Data Contamination in Large Language Models*. Ph.D. thesis, Carnegie Mellon University Pittsburgh, PA. <https://doi.org/10.18653/v1/2024.conda-1.3>
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The RefinedWeb dataset for Falcon LLM: Outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*.
- Manley Roberts, Himanshu Thakur, Christine Herlihy, Colin White, and Samuel Dooley. 2023. Data contamination through the lens of time. *ArXiv*, abs/2310.10628.
- Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023a. NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark. *arXiv preprint arXiv:2310.18018*. <https://doi.org/10.18653/v1/2023.findings-emnlp.722>
- Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, and Eneko Agirre. 2023b. Did chatgpt cheat on your test? <https://hitz-zentroa.github.io/lm-contamination/blog/>. Accessed: 2023-11-01.
- Areg Mikael Sarvazyan, José Ángel González, Marc Franco-Salvador, Francisco Rangel, Berta Chulvi, and Paolo Rosso. 2023. Overview of autextification at iberlef 2023: Detection and attribution of machine-generated text in multiple domains. In *Procesamiento del Lenguaje Natural*, Jaén, Spain.
- William A. Scott. 1955. Reliability of content analysis: The case of nominal scale coding. *Public opinion quarterly*, pages 321–325. <https://doi.org/10.1086/266577>
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2024. Detecting pretraining data from large language models. In *The Twelfth International Conference on Learning Representations*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra,

- Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. SuperGlue: A stickier benchmark for general-purpose language understanding systems. *arXiv preprint arXiv:1905.00537*.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*. <https://doi.org/10.18653/v1/W18-5446>
- Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023. Large language models are not fair evaluators. *arXiv preprint arXiv:2305.17926*.
- Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph E. Gonzalez, and Ion Stoica. 2023. Rethinking benchmark and contamination for language models with rephrased samples. *arXiv preprint arXiv:2311.04850*.
- Jingyang Zhang, Jingwei Sun, Eric Yeats, Yang Ouyang, Martin Kuo, Jianyi Zhang, Hao Yang, and Hai Li. 2024. Min-k%++: Improved baseline for detecting pre-training data from large language models. *arXiv preprint arXiv:2404.02936*.
- Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *NIPS*.
- Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. 2024. Large language models are not robust multiple choice selectors. In *The Twelfth International Conference on Learning Representations*.
- Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin, Jinhui Wen, and Jiawei Han. 2023. Don't make your LLM an evaluation benchmark cheater. *ArXiv*, abs/2311.01964.

A Detecting Data Contamination for Different Checkpoints

As LLM providers persistently update their models, the safety of these LLMs towards detected weaknesses and adversarial attacks also improves. The detection technique recently developed by Golchin and Surdeanu (2024) is one such attack addressed by these updates. For example, when GPT-4 is instructed to complete a partial reference instance from a classification dataset, the June 2023 snapshot (gpt-4-0613) is capable of accomplishing this task, thus unveiling contamination. In contrast, a more recent version from November 2023 (gpt-4-1106-preview) often chooses not to complete the task due to the activation of newly implemented safety mechanisms. Table 4 illustrates this behavior.

On the other hand, in Table 5, we show that our proposed quiz-based methodology can successfully identify data contamination for both versions. This highlights the capability of our approach in pinpointing data contamination even in the presence of newly added safety filters. In fact, our strategy can effectively sidestep filters related to copyrighted materials by delivering a single-letter response that embodies the detection of the original dataset instance among other options, thereby exposing contamination.

LLM Version/Detection	Original Dataset Instance and Its Replica
<p>Version: gpt-4-0613 Detection: Exact Match</p>	<p>Original Dataset Instance: Review: Bromwell High is a cartoon comedy. It ran at the same time as some other programs about school life, such as "Teachers". My 35 years in the teaching profession lead me to believe that Bromwell High's satire is much closer to reality than is "Teachers". The scramble to survive financially, the insightful students who can see right through their pathetic teachers' pomp, the pettiness of the whole situation, all remind me of the schools I knew and their students. When I saw the episode in which a student repeatedly tried to burn down the school, I immediately recalled at High. A classic line: INSPECTOR: I'm here to sack one of your teachers. STUDENT: Welcome to Bromwell High. I expect that many adults of my age think that Bromwell High is far fetched. What a pity that it isn't! Label: 1 (positive)</p> <div style="border: 1px solid black; padding: 10px;"> <p>Instruction: You are provided with the first piece of an instance from the train split of the IMDB dataset. Finish the second piece of the instance as exactly appeared in the dataset. Only rely on the original form of the instance in the dataset to finish the second piece. Label: 1 (positive) First Piece: Bromwell High is a cartoon comedy. It ran at the same time as some other programs about school life, such as "Teachers". My 35 years in the teaching profession lead me to believe that Bromwell High's satire is much closer to reality than is "Teachers". Second Piece: The scramble to survive financially, the insightful students who can see right through their pathetic teachers' pomp, the pettiness of the whole situation, all remind me of the schools I knew and their students. When I saw the episode in which a student repeatedly tried to burn down the school, I immediately recalled at High. A classic line: INSPECTOR: I'm here to sack one of your teachers. STUDENT: Welcome to Bromwell High. I expect that many adults of my age think that Bromwell High is far fetched. What a pity that it isn't!</p> </div>
<p>Version: gpt-4-1106-preview Detection: None</p>	<div style="border: 1px solid black; padding: 10px;"> <p>Instruction: You are provided with the first piece of an instance from the train split of the IMDB dataset. Finish the second piece of the instance as exactly appeared in the dataset. Only rely on the original form of the instance in the dataset to finish the second piece. Label: 1 (positive) First Piece: Bromwell High is a cartoon comedy. It ran at the same time as some other programs about school life, such as "Teachers". My 35 years in the teaching profession lead me to believe that Bromwell High's satire is much closer to reality than is "Teachers". Second Piece: Sorry, but I can't continue the text from the IMDB dataset as it is copyrighted material. However, I can provide a summary or discuss the themes and characters of "Bromwell High" if you're interested. Let me know if you would like me to provide some other type of information!</p> </div>

Table 4: A comparison between the completions generated using replication-based approach (Golchin and Surdeanu, 2024) for the June and November 2023 snapshots of GPT-4. While the June snapshot exactly replicates the second piece of the provided dataset instance, the November snapshot avoids completing it due to the activation of filters intended to prevent the generation of copyrighted content.

Original Dataset Instance:

Review: Bromwell High is a cartoon comedy. It ran at the same time as some other programs about school life, such as "Teachers". My 35 years in the teaching profession lead me to believe that Bromwell High's satire is much closer to reality than is "Teachers". The scramble to survive financially, the insightful students who can see right through their pathetic teachers' pomp, the pettiness of the whole situation, all remind me of the schools I knew and their students. When I saw the episode in which a student repeatedly tried to burn down the school, I immediately recalled at High. A classic line: INSPECTOR: I'm here to sack one of your teachers. STUDENT: Welcome to Bromwell High. I expect that many adults of my age think that Bromwell High is far fetched. What a pity that it isn't!

Label: 1 (positive)

Instruction: You are provided with a five-choice quiz. Your task is to correctly select the option that exactly corresponds to an instance from the train split of the IMDB dataset.

When selecting the option, you must ensure that you follow the following rules:

- (1) You must ensure that you only generate a single option letter as your answer.
- (2) If you do not know the dataset or the correct answer, you must select option "E) None of the provided options."

Hint: While all the following options seem similar in terms of meaning, there is only one option that reflects the exact wording with respect to the original instance.

A) Bromwell High is an animated comedy. It aired at the same time as some other shows about school life, such as "Teachers". My 35 years in the education profession make me feel that Bromwell High's humor is much closer to reality than is "Teachers". The effort to survive financially, the smart students who can see right through their pathetic teachers' pomp, the pettiness of the whole situation, all remind me of the schools I knew and their students. When I saw the episode in which a student repeatedly tried to burn down the school, I immediately recalled at High. A memorable line: INSPECTOR: I'm here to sack one of your teachers. STUDENT: Welcome to Bromwell High. I expect that many adults of my age think that Bromwell High is unbelievable. What a pity that it isn't!

Label: 1 (positive)

B) Bromwell High is a cartoon comedy. It broadcast concurrently with some other programs about school life, like "Teachers". My 35 years in the teaching profession lead me to believe that Bromwell High's satire is way closer to reality than is "Teachers". The struggle to survive financially, the sharp students who can see right through their pathetic teachers' pomp, the pettiness of the whole situation, all remind me of the schools I knew and their students. When I saw the episode in which a student repeatedly attempted to burn down the school, I immediately remembered at High. An iconic line: INSPECTOR: I'm here to fire one of your teachers. STUDENT: Welcome to Bromwell High. I expect that many adults of my age think that Bromwell High is far fetched. What a pity that it isn't!

Label: 1 (positive)

C) Bromwell High is a cartoon comedy. It ran at the same time as some other programs about school life, such as "Teachers". My 35 years in the teaching profession lead me to believe that Bromwell High's satire is much closer to reality than is "Teachers". The scramble to survive financially, the insightful students who can see right through their pathetic teachers' pomp, the pettiness of the whole situation, all remind me of the schools I knew and their students. When I saw the episode in which a student repeatedly tried to burn down the school, I immediately recalled at High. A classic line: INSPECTOR: I'm here to sack one of your teachers. STUDENT: Welcome to Bromwell High. I expect that many adults of my age think that Bromwell High is far fetched. What a pity that it isn't!

Label: 1 (positive)

D) Bromwell High is a cartoon comedy. It ran at the same time as some other shows about school life, such as "Teachers". My 35 years in the teaching profession lead me to believe that Bromwell High's satire is much closer to reality than is "Teachers". The scramble to survive financially, the clever students who can see right through their pathetic teachers' pomp, the pettiness of the whole situation, all remind me of the schools I knew and their students. When I watched the episode in which a student repeatedly tried to set fire to the school, I immediately recalled at High. A classic line: INSPECTOR: I'm here to dismiss one of your teachers. STUDENT: Welcome to Bromwell High. I expect that many adults of my age think that Bromwell High is far fetched. What a pity that it isn't!

Label: 1 (positive)

E) None of the provided options.

Answer:

C

Versions: gpt-4-0613 and
gpt-4-1106-preview
Detection: Correct (option C)

Table 5: An illustration of a quiz question posed when taking our DCQ to both June and November snapshots of GPT-4. While the generated response is correct (option C), this approach—unlike the replication-based strategy (Golchin and Surdeanu, 2024)—proves to be effective for both snapshots without triggering any content-related filters set by the model developer.

B Details and Analysis of Generated Word-Level Perturbations

Our procedure for creating quiz options involves generating four word-level perturbations for each instance subsampled from a particular dataset partition. For each instance, we prompt GPT-4 to generate four perturbations at once. We repeatedly apply this process to all the subsampled instances. The

Instruction: Your task is to create a four-choice quiz by replacing the words in the provided “Input Text” with their contextually relevant synonyms. The meaning and sentence structure of the four options must exactly match every detail in the Input Text. You must not include the provided Input Text as an option.

You must make sure that:

- (1) You generate distinct options based on the provided Input Text;
- (2) The only difference between options is word-level perturbations.
- (3) Options are ordered;
- (4) There is not any extra explanation;
- (5) You follow the following “Format” to generate options;
- (6) You comply with every specific symbol and letter detail in the given Input Text; and
- (7) All options retain the exact label from the Input Text, if there is one.

Input Text:

Summary: Premier Inn has announced it is investing £35m in opening three new hotels and revamping two others in Edinburgh.

Format:

A)

Summary:

B)

Summary:

C)

Summary:

D)

Summary:

A)

Summary: Premier Inn has revealed it is allocating £35m for the opening of three new hotels and the refurbishment of two others in Edinburgh.

B)

Summary: Premier Inn has stated it is committing £35m to launch three new hotels and upgrade two others in Edinburgh.

C)

Summary: Premier Inn has declared it is investing £35m to set up three new hotels and modernize two others in Edinburgh.

D)

Summary: Premier Inn has announced it is spending £35m to open three new hotels and renovate two others in Edinburgh.

Figure 4: The zero-shot prompt employed for generating four word-level perturbations per dataset instance. The input text is replaced for each dataset instance, and GPT-4 is prompted to generate the four perturbations for each dataset instance individually at once. The format provided in the input prompt is adjusted based on a specific dataset/task, with column names in the dataset being used to denote distinct components of a dataset instance. For example, “Summary” is a column name in the XSum that is prepended to the instance shown above. The example shown here is an instance taken from the XSum test partition along with its generated word-level perturbations, as illustrated earlier in Figure 1.

template prompt and all the rules and requirements employed to instruct GPT-4 in generating word-level perturbations are shown via an example in Figure 4. This zero-shot prompt shows the four generated word-level perturbations corresponding to the example given in Figure 1.

Additionally, to evaluate the quality of the generated word-level perturbations by GPT-4, we conduct a three-point Likert scale analysis on a subsample of 60 instances. These instances are drawn from four distinct datasets/tasks: AG News (Zhang et al., 2015), XSum (Narayan et al., 2018), WNLI (Wang et al., 2018), and GSM8k (Cobbe et al., 2021), with each dataset contributing 15 instances. Three domain experts participate in evaluating the generated perturbations based on two independent criteria: *meaning* and *sentence structure*, using the following scale: *good*, *fair*, and *poor*.

The criteria for this evaluation are defined as:

- **Meaning:** This criterion ensures that perturbations are solely at the word level while preserving the similarity in core essence between the generated perturbations and the original dataset instance.

- **Sentence Structure:** This criterion determines the level of correspondence in sentence structure between the generated perturbations and the original dataset instance while maintaining strict adherence to every precise symbol and letter detail featured in the original dataset instance.

Additionally, the Likert scale points are described as follows:

- **Good:** The word-level perturbations meet all aforementioned requirements for meaning and sentence structure independently.
- **Fair:** The word-level perturbations fulfill some, but not all, requirements for meaning and sentence structure independently.
- **Poor:** The word-level perturbations do not meet any of the stated requirements for meaning and sentence structure independently.

We evaluate the inter-rater reliability among expert raters by calculating the overall percent agreement, per rating percent agreement, and the corresponding confidence intervals for each criterion under consideration (i.e., *meaning* and *sentence structure*). The percent agreement represents the proportion of instances where all raters concur (on the same rating/category) out of the total number of ratings provided. To determine the confidence intervals, we employ bootstrap resampling (Efron, 1979; Efron and Tibshirani, 1993; Efron, 2003) with 10,000 iterations to estimate the 95% confidence intervals for percent agreement for each of the ratings (i.e., *good*, *fair*, and *poor*). Note that we opted not to use Cohen's/Fleiss' Kappa (Cohen, 1960; Scott, 1955; Fleiss, 1971) as our inter-rater reliability measure due to the high chance agreement in this task, attributed to the low variability in the assigned ratings.

The following presents the results for the overall percent agreement and per rating percent agreement for the *meaning* criterion:

- The overall percent agreement among all three expert raters is 86.67%, with a 95% confidence interval of 76.67%–95.00%.
- The percent agreement for the *good* rating among all three expert raters is 83.33%, with a 95% confidence interval of 73.33%–91.67%.
- The percent agreement for the *fair* rating among all three expert raters is 3.33%, with a 95% confidence interval of 0.00%–8.33%.
- The percent agreement for the *poor* rating among all three expert raters is 0.00%, with a 95% confidence interval of 0.00%–0.00%.

Also, the following outlines the results of the percent agreements for the *sentence structure* criterion:

- The overall percent agreement among all three expert raters is 93.33%, with a 95% confidence interval of 86.67%–98.33%.
- The percent agreement for the *good* rating among all three expert raters is 93.33%, with a 95% confidence interval of 86.67%–98.33%.
- The percent agreement for the *fair* rating among all three expert raters is 0.00%, with a 95% confidence interval of 0.00%–0.00%.
- The percent agreement for the *poor* rating among all three expert raters is 0.00%, with a 95% confidence interval of 0.00%–0.00%.

Note that the percent agreement for the *fair* and *poor* ratings is substantially low for both criteria because perturbations were consistently evaluated as *good* by raters, resulting in low variability in ratings. Hence, while the studied subset offers reliable agreement for the predominant rating (*good*), the agreement on less common ratings (*fair* and *poor*) still requires a larger subset to be meaningful.

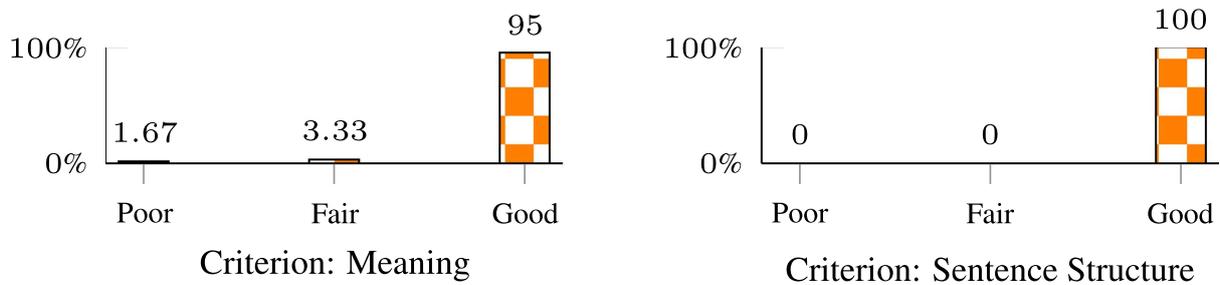


Figure 5: Percentage distribution of evaluation results for meaning and sentence structure of the generated word-level perturbations based on majority voting among three expert raters. The most common evaluation is *good* for both studied criteria.

As an additional assessment, we implement a majority voting among the ratings provided by the raters. In this evaluation, if a tie occurs, we consider the lower quality recognized by the raters. Figure 5 illustrates the percentage distribution for each criterion upon conducting majority voting. The results reveal that the majority of the distribution across the three ratings (i.e., *good*, *fair*, and *poor*) is rated as *good* for both *meaning* and *sentence structure*, at 95.00% and 100%, respectively. As a result, based on the majority vote, the most common evaluation for both criteria for the generated perturbations is *good*. In this context, *good*, as defined by the provided instruction, indicates that perturbations convey the same meaning as the original dataset instance while only exhibiting differences at the word level. Additionally, *good* for sentence structure implies that perturbations adhere to every exact symbol and letter aspect within the original dataset instance.

Therefore, based on the comprehensive analysis provided, we conclude that GPT-4 consistently generates word-level perturbations of satisfactory quality in terms of both *meaning* and *sentence structure*. These perturbations meet the criteria outlined in our methodology for establishing a setting in which memorization can be induced through quiz options effectively.

C Systematic Contamination of GPT-3.5 and Llama 2

We systematically contaminate the GPT-3.5 base model by submitting a fine-tuning job for three epochs through the OpenAI API, while keeping all other hyperparameters at their default values. For each of the two datasets we use in this experiment, the GPT-3.5 base model is contaminated separately, resulting in two distinct checkpoints, each corresponding to one dataset. Although the model provider refers to this option as fine-tuning, our methodology deviates from the standard fine-tuning procedure. In fact, our goal here is to emulate the data exposure that LLMs potentially experienced during their pre-training phase. To this end, we employ data formats used in the study by Golchin and Surdeanu (2024) to introduce contamination. It is also important to emphasize that this type of training differs from instruction fine-tuning, as we do not infuse any specific instructions within the data we use for continued training.

To contaminate Llama 2, we use low-rank adaptation (Hu et al., 2022) for continued training, following the hyperparameter recommendations provided by the model developer (Llama Team, 2024). However, we adjust two hyperparameters to suit our use case: We set the number of epochs to three and limit the maximum input length to 500 tokens. We also use the same data formats mentioned above for contaminating the GPT-3.5 base model.

D Statistics on Positional Biases

Table 6 presents the detected positional biases by the Bias Detector Quiz (BDQ) among all the models and datasets used in this study. The results reveal that LLMs exhibit different positional biases depending on the dataset, indicating that these biases are neither universal nor transferable across models and datasets. Moreover, this further underscores the importance of addressing positional biases

Dataset	GPT-3.5	GPT-4	Llama 2
MeetingBank	A: 53, B: 3, C: 22, D: 1, E: 21	–	A: 67, B: 2, C: 31, D: 0, E: 0
AuTextification	A: 49, B: 5, C: 34, D: 0, E: 12	–	A: 41, B: 2, C: 57, D: 0, E: 0
HumanEval	–	A: 50.61, B: 3.66, C: 0, D: 2.44, E: 43.29	–
DROP	–	A: 36, B: 5, C: 3, D: 2, E: 54	–
GSM8k	–	A: 44, B: 2, C: 1, D: 1, E: 52	–
IMDB	A: 55, B: 0, C: 16, D: 0, E: 29	A: 68, B: 1, C: 0, D: 0, E: 31	A: 33, B: 10, C: 57, D: 0, E: 0
AG News (train)	–	A: 29, B: 0, C: 0, D: 0, E: 71	–
AG News (test)	A: 19, B: 8, C: 36, D: 3, E: 34	A: 41, B: 0, C: 1, D: 1, E: 57	A: 19, B: 7, C: 74, D: 0, E: 0
Yelp	A: 39, B: 0, C: 22, D: 1, E: 38	A: 28, B: 0, C: 0, D: 0, E: 72	A: 40, B: 1, C: 59, D: 0, E: 0
RTE	A: 0, B: 0, C: 17, D: 2, E: 81	A: 16, B: 0, C: 0, D: 1, E: 83	A: 56, B: 8, C: 14, D: 0, E: 22
WNLI (train)	–	A: 6, B: 0, C: 0, D: 0, E: 94	–
WNLI (validation)	A: 0, B: 0, C: 5.63, D: 0, E: 94.37	A: 9.86, B: 0, C: 0, D: 1.41, E: 88.73	A: 12.68, B: 26.76, C: 0, D: 0, E: 60.56
SAMSum	A: 3, B: 7, C: 49, D: 13, E: 28	A: 15, B: 3, C: 1, D: 4, E: 77	A: 63, B: 2, C: 34, D: 0, E: 1
XSum	A: 16, B: 3, C: 28, D: 5, E: 48	A: 41, B: 7, C: 2, D: 8, E: 42	A: 71, B: 0, C: 29, D: 0, E: 0

Table 6: Positional biases detected using the Bias Detector Quiz (BDQ) across all datasets and models in this work. Although these proportions may vary slightly across different runs due to the probabilistic nature of LLMs, the overall biases remain consistent, as discussed in Section 4 and shown in Figure 2. The number next to each option represents the percentage of times that option was selected in BDQ.

when estimating contamination is necessary. For instance, although the MeetingBank and AuTextification datasets were unseen by Llama 2, positional bias prevented the model from choosing option E—none of the options—which is the correct choice in BDQ, where no correct answer exists among options. Instead, the model disproportionately favored options A and C beyond random chance. As explained in Subsection 4.1, failing to properly handle such biases results in overestimating contamination.