

# Investigating How Pre-training Data Leakage Affects Models' Reproduction and Detection Capabilities

Masahiro Kaneko Timothy Baldwin



Masahiro.Kaneko@mbzuai.ac.ae Timothy.Baldwin@mbzuai.ac.ae

## Abstract

Large Language Models (LLMs) are trained on massive web-crawled corpora, often containing personal information, copyrighted text, and benchmark datasets. This inadvertent inclusion in the training dataset, known as data leakage, poses significant risks and could compromise the safety of LLM outputs. Despite its criticality, existing studies do not examine how leaked instances in the pre-training data influence LLMs' output and detection capabilities. In this paper, we conduct an experimental survey to elucidate the relationship between data leakage in training datasets and its effects on the generation and detection by LLMs. Our experiments reveal that LLMs often generate outputs containing leaked information, even when there is little such data in the training dataset. Moreover, the fewer the leaked instances, the more difficult it becomes to detect such leakage. Finally, we demonstrate that enhancing leakage detection through few-shot learning can help mitigate the impact of the leakage rate in the training data on detection performance.

## 1 Introduction

Large Language Models (LLMs) have achieved remarkable performance in various real-world applications (Brown et al., 2020; Wei et al., 2021; Ouyang et al., 2022). One of the success factors is the massive web-crawled corpora used for pre-training LLMs (Kaplan et al., 2020; Wei et al., 2022). The corpora for pre-training LLMs consist of webpages, books, scientific papers, and programming code (Almazrouei et al., 2023; Zhao et al., 2023). Developers of well-known LLMs such as ChatGPT<sup>1</sup> and Claude 3<sup>2</sup> do not disclose the composition of the training data, to maintain a competitive edge. The large-scale nature and privatization of such training data increases the risk of leaking

inappropriate data such as personal information, copyrighted texts, and benchmarks (Ishihara, 2023; Yang et al., 2023; Jiang et al., 2024; Kaneko et al., 2025).

Nasr et al. (2023) have revealed that it is possible to efficiently recover training data from LLMs under various settings. In practice, it has been confirmed that personal information, such as names, phone numbers, and email addresses, has leaked from LLMs (Shokri et al., 2016; Carlini et al., 2020; Huang et al., 2022; Kim et al., 2023). The leak of benchmarks enhances the reported performance of LLMs (Deng et al., 2023; Zhou et al., 2023), leading to over-confidence in the abilities of LLMs. Eldan and Russinovich (2023) show that copyrighted texts such as news articles<sup>3</sup> and books<sup>4</sup> can be reproduced by LLMs. It has been revealed that leaked instances have a higher output probability in LLMs compared to non-leaked instances, indicating a potential for leakage detection (Yeom et al., 2017; Shi et al., 2023). The LLMs' ability to detect leakage is effective in proactively defending against malicious users extracting leaked instances from LLMs (Wang et al., 2024). These studies demonstrate that instances leaked in the training data affect the reproducibility and detectability of leaked instances in LLMs.

Existing research discusses the risks of data leakage and attempts to detect such leaked instances. However, the influence of pre-training data, which is considered the root cause of such leakage, on the behavior of LLMs remains insufficiently understood. Clarifying this leads to the construction of pre-training data that contributes to preventing the leakage problem. In this study, we investigate

<sup>1</sup><https://chat.openai.com/>

<sup>2</sup><https://claude.ai/chats>

<sup>3</sup><https://www.nytimes.com/2023/12/27/business/media/new-york-times-open-ai-microsoft-lawsuit.html>

<sup>4</sup><https://www.theatlantic.com/technology/archive/2023/08/books3-ai-meta-llama-pirated-books/675063/>

how leaked instances in pre-training data affect the model’s reproducibility and detectability. First, we identify the extent to which the targeted leaked instances are present in the pre-training data. Next, we examine the impact of these leaked instances on the model’s tendency to generate leaked instances and the detectability of such instances.

In our experiments, we investigate the proportion of leaked instances in the pre-training data related to personal information, copyrighted texts, and benchmarks across five LLMs.<sup>5</sup> Our experimental results show that when there is little leakage in the pre-training data, it does not affect the tendency of LLMs to reproduce leaked instances, yet detecting the leaked instances becomes more difficult. Therefore, when filtering leaked instances from the pre-training data, it is necessary to ensure that the model’s detection performance does not degrade.

Finally, we aim to mitigate the negative impact of the leakage rate on the detection performance of LLMs. Existing methods (Yeom et al., 2017; Carlini et al., 2020; Shi et al., 2023; Kaneko et al., 2024) do not explicitly supervise the task of classifying leaked and non-leaked instances for detectors. We demonstrate that explicitly supervising the model with leaked and non-leaked instances can complement its implicit reliance on leaked instances in the training data, thereby preventing a decline in detection performance. Our experimental results show that the supervised detection using few-shot method performs on average about 7 points higher than existing methods. On the other hand, the detection rate drops in the zero-shot settings, suggesting that providing examples for supervising LLMs is particularly important.

## 2 Investigating Infection of Leaked Instances

To investigate infection of leaked instances in pre-trained data for the model’s reproducibility and detectability, we define the following three criteria:

- **Leakage Rate** refers to the proportion of target leaked instances contained in the entire pre-training data of LLMs.
- **Reproduction Rate** refers to the proportion of leaked instances in the pre-training data that the LLMs reproduce.

<sup>5</sup>Due to privacy and copyright constraints, we will release the code and the permissibly shareable subset of the data upon acceptance.

- **Detection Rate** refers to the performance of LLMs in distinguishing between leaked and non-leaked instances in their pre-training dataset.

We conduct an experimental survey to elucidate the relationship between the leakage rate and both the reproduction rate and detection rate for personal information, copyrighted texts, and benchmark data.

### 2.1 Leakage Rate

The leakage rate is the proportion within the leakage instances we targeted in the pre-training dataset, including personal information, copyrighted texts, and benchmark datasets. We target the training data used by LLMs whose experimental settings are publicly available for our experiments. We begin by listing publicly available LLMs and curating their training data. Next, we introduce how to calculate the leakage rate for personal information, copyrighted texts, and benchmarks in the pre-training data of LLMs.

**Pre-training Datasets** In this study, we target the pre-training data of the following six LLMs<sup>6</sup> for which the details of the experimental setup are publicly available.

- **T5** (Raffel et al., 2019): T5 uses the Colossal Clean Crawled Corpus (C4)<sup>7</sup> containing about 800 GB of text data collected from filtered web pages as its pre-training data. Scientific texts, books, and news account for approximately 25% in C4. The filtering includes the removal of inappropriate content, deletion of duplicates, and detection of language.
- **LLaMA** (Touvron et al., 2023a): LLaMA employs English CommonCrawl, C4, Github, Wikipedia, Books, ArXiv, and StackExchange as pre-training datasets. We conducted our experiments using the RedPajama dataset (Computer, 2023)<sup>8</sup>, which is developed to closely mirror the publicly described LLaMA training data.
- **Pythia** (Biderman et al., 2023a): Pythia uses

<sup>6</sup>Our experiment, similarly to other studies on data leakage and contamination, focuses exclusively on open-source LLMs due to the necessity of accessing publicly available pre-training data. Furthermore, this direction aligns with the NLP community’s recommendation of using open-source models for empirical validation (Groeneveld et al., 2024a).

<sup>7</sup><https://huggingface.co/datasets/legacy-datasets/c4>

<sup>8</sup><https://github.com/togethercomputer/RedPajama-Data>

| LLMs   | Size  | C4            | CommonCrawl  | The Pile      | GitHub | Wikipedia | Books | Papers | Conversations |
|--------|-------|---------------|--------------|---------------|--------|-----------|-------|--------|---------------|
| T5     | 800   | <b>100.0%</b> | 0.0%         | 0.0%          | 0.0%   | 0.0%      | 0.0%  | 0.0%   | 0.0%          |
| LLaMA  | 4,700 | 15.0%         | <b>67.0%</b> | 0.0%          | 4.5%   | 4.5%      | 4.5%  | 2.5%   | 2.0%          |
| Pythia | 800   | 0.0%          | 0.0%         | <b>100.0%</b> | 0.0%   | 0.0%      | 0.0%  | 0.0%   | 0.0%          |
| MPT    | 4,000 | <b>63.4%</b>  | 8.5%         | 0.0%          | 14.5%  | 4.0%      | 3.0%  | 5.2%   | 1.4%          |
| Falcon | 3,600 | 0.0%          | <b>84.0%</b> | 0.0%          | 3.0%   | 1.0%      | 6.0%  | 1.0%   | 5.0%          |
| OLMo   | 5,300 | 5.7%          | <b>78.7%</b> | 0.0%          | 12.6%  | 0.1%      | 0.1%  | 2.8%   | 0.0%          |

Table 1: The total volume and the percentage of sources in datasets used for pre-training each LLM. These datasets undergo different filtering and refinement processes for each LLM. The unit of size for the dataset is in GB.

the Pile<sup>9</sup>, which comprises 800GB of text data. It aggregates content from 22 different sources, including books, websites, GitHub repositories, and more.

- **MPT** (Team, 2023): MPT uses the dataset<sup>10</sup>, which preprocesses the Common Crawl, Wikipedia, Books, ArXiv, and StackExchange to remove low-quality content and duplicate pages.
- **Falcon** (Almazrouei et al., 2023): Falcon utilizes the RefinedWeb dataset (Penedo et al., 2023b)<sup>11</sup>, which employs heuristic rules to filter the Common Crawl dataset and remove duplicates.
- **OLMo** (Groeneveld et al., 2024a): OLMo uses Dolma (Soldaini et al., 2024)<sup>12</sup>, which is a dataset of 3T tokens from a diverse mix of web content, academic publications, code, books, and encyclopedic materials.

We present the configuration of the LLMs and the pre-training data used in our experiments in Table 1. The most common sources included in all LLMs are web page sources such as C4, CommonCrawl, and the Pile. Because they are collected from various web pages, there is a risk that they may contain personal information, copyrighted texts, or benchmarks. For example, the C4 includes personal information such as voter lists and pirated e-books that violate copyright laws.<sup>13</sup> We used the entire pre-training data used in each LLM and investigated the leakage rates of personal information, copyrighted texts, and benchmarks.

### Scopes of Leakage Instances in the Pre-training Datasets

We determine whether personal infor-

mation is included in the text through regular expressions proposed in the existing research (Subramani et al., 2023). This regular expression targets 20 types<sup>14</sup> of personal information. Additionally, we determine whether a person’s name is included in the text using named entity recognition from the spaCy library<sup>15</sup>. Based on existing research (Finck and Pallas, 2020), we do not distinguish between real names and pseudonyms in our study, as both can impact an individual’s privacy. If the target text contains even one piece of personal information, we determine that it is leaking. We targeted books, news articles, and papers found on Google Books<sup>16</sup>, Google News<sup>17</sup>, and Google Scholar<sup>18</sup> as the subjects of the copyrighted texts. We use the Selenium library to automate the search process. For the leakage rate of benchmarks, it is challenging to cover all benchmarks. Therefore, considering that the negative impact of leakage becomes more problematic for larger benchmarks widely used by many users, we limit our focus to the top benchmarks by download count. We create a data store from a total of approximately 200,000 instances contained in the test data from Huggingface’s Database, which are among the top 128 in terms of download count.<sup>19</sup> Since the training dataset is not problematic even if it is included in the pre-training dataset, we extract the development dataset and test dataset. When one instance contains multiple texts, such as context and questions, we add each text separately to the data store.

Existing research defined data leakage for copy-

<sup>9</sup><https://huggingface.co/datasets/EleutherAI/pile>

<sup>10</sup><https://github.com/mosaicml/llm-foundry>

<sup>11</sup><https://huggingface.co/datasets/tiiuae/falcon-refinedweb>

<sup>12</sup><https://huggingface.co/datasets/allenai/dolma>

<sup>13</sup><https://www.washingtonpost.com/technology/interactive/2023/ai-chatbot-learning/>

<sup>14</sup>The regular expressions to find personal information: *IP address, IBAN code, US SSN, email addresses, phone numbers, amex card, bcglobal, carte blanche card, diners club card, discover card, insta payment card, jcb card, korean local card, laser card, maestro card, mastercard, solo card, switch card, union pay card, and visa card*

<sup>15</sup><https://spacy.io/usage/linguistic-features>

<sup>16</sup><https://books.google.com/>

<sup>17</sup><https://news.google.com/>

<sup>18</sup><https://scholar.google.com/>

<sup>19</sup><https://huggingface.co/datasets>

righted text as matching approximately 50 words between texts (Karamolegkou et al., 2023). Following this precedent, we exclude texts shorter than 50 words from datasets and data stores for copyrighted text. For personal information and benchmark datasets, we do not set a length limitation. If the target text is found through an exact match search, we consider that a leak. The leakage rate is calculated by dividing the number of leaked instances by the total number of instances for each dataset. The leakage rate is calculated by dividing the total size of leaked instances by the total data size in Table 1. Calculating the ratio based on data size rather than on an instance basis is to mitigate the impact of differences in instance-level granularity across datasets.

Our research limits the scope of leakage targets through the sampling of training data and the identification of leaked instances using regular expressions, web searches, and databases. On the other hand, it is not practical from a resource perspective to comprehensively cover all leakage instances related to personal information, copyrighted texts, and benchmarks across the entire training data. Since our definition mentioned above targets representative cases of leakage, the insights gained can be broadly applicable even within a limited scope.

## 2.2 Reproduction Rate

We create datasets containing leaked and denied texts to calculate the reproduction rate. The leaked text is the text of leaked instances in the pre-training datasets. The denied text is the text where the LLM denies responding. We compare the likelihood of the LLM generating the leaked text and the denied text, respectively, in response to prompts that elicit the leakage. Then, we calculate the reproduction rate as the proportion of cases where the likelihood of the leaked text is higher than the likelihood of the denied text.

First, we introduce how to create the leaked text. Since the evaluation data needs to be manually checked for quality, we sample the leaked instances. We sample 100,000 instances each from those identified as leaks of personal information, copyrighted texts, and benchmarks in the pre-training data in Section 2.1. We provide the sampled instances to GPT-4<sup>20</sup> and generate prompts that cause the LLMs to output those instances. Since OpenAI does not use data inputted to their API for train-

ing, there is no concern about these instances leaking.<sup>21</sup> We create 8 instances for GPT-4 to perform few-shot learning, such as: “Text: Michel \*\*\*\*: +\*\*\*\*-\*\*\*\*-\*\*\*\* Prompt: Please tell me Michel \*\*\*\*’s phone number.”. Table 2 shows the created examples for few-shot learning. In this example, the name and phone number are obfuscated with \*\*\*\* for privacy protection. We use the following prompt:

```
Please write a prompt to generate the given text.
Text: [Text Example 1] Prompt: [Prompt Example 1]
:
Text: [Text Example 8] Prompt: [Prompt Example 8]
Text: [Instance] Prompt:
```

Here, [Text Example 1], [Text Example 8], [Prompt Example 1], and [Prompt Example 8] are few-shot examples. We use texts randomly sampled from the pre-training data as few-shot examples. Since the outputs of LLMs are sensitive to the prompt, we create four prompts for each instance, average the accuracy, and evaluate them. We exclude prompts that are inappropriate or copyrighted text that is properly quoted, based on evaluations by 40 crowd workers on Amazon Mechanical Turk.<sup>22</sup> We examined the agreement rate among workers regarding the labeling of whether to exclude or not exclude the sampled 100 instances. As a result, the labels of more than 80% of the workers matched for 83% of the instances. After manual filtering, the data sizes for personal information, copyrighted texts, and benchmarks became 9,011, 7,702, and 8,155, respectively.

We create 8 denied texts where the response is to reject responding, such as “Sorry, I cannot respond regarding this content.”.<sup>23</sup> We provide the LLM with prompts concatenate with the leaked text and the denied text, respectively. To mitigate the influence of the prompt, we use the average likelihood of the 8 denied texts. We calculate the reproduction rate by dividing the total size of instances where the likelihood of the leaked text is higher than that of the denying text by the total data size in Table 1.

<sup>21</sup><https://help.openai.com/en/articles/5722486-how-your-data-is-used-to-improve-model-performance>

<sup>22</sup>We set the hourly rate for the work at \$15. We achieved an 80% human agreement rate, which ensures the validity of the evaluation to some degree for 30 sampled instances.

<sup>23</sup>We sample reject respondings from existing datasets. (Wang et al., 2023)

<sup>20</sup><https://openai.com/gpt-4>



| Personal Information |   |
|----------------------|---|
| Text                 | **** Corporation: +44****_****_****   |
| Prompt               | Could you please tell me the phone number of **** Corporation in the UK?  |
| Copyrighted Text     |   |
| Text                 | They didn't think they could bear it if anyone found out about the Potters.<br>...<br>This boy was another good reason for keeping the Potters away; they didn't want Dudley mixing with a child like that.       |
| Prompt               | Please generate the continuation of the following text with more than 50 words: The Dursleys had everything they wanted, but they also had a secret, and their greatest fear was that somebody would discover it. |
| Benchmark            |   |
| Text                 | premise: The girl is climbing a rock wall.<br>hypothesis: A girl climbs away from danger up a rock wall.<br>label: neutral  |
| Prompt               | Generate a premise text that results in a neutral label when the hypothesis is "A girl climbs away from danger up a rock wall."   |

Table 2: Few-shot examples for reproduction rate in personal information, copyrighted text, and benchmark dataset. The text corresponding to personal information is masked with \*\*\*\*, but in the actual input to the LLM, it is not masked.

### 2.3 Detection Rate

The detection rate is the proportion of cases where the LLM correctly classifies between leaked instances included in the pre-training dataset and non-leaked instances not included. We create a non-leaked dataset composed of instances not included in the pre-training data, for the leaked dataset created in Section 2.2. For personal information, we create the non-leaked dataset by replacing numbers such as phone numbers and credit card numbers with random digits, and rewriting texts such as names and addresses to different names and addresses using GPT-4. For copyrighted texts and benchmarks, we use GPT-4 to generate paraphrases to create the non-leaked dataset. Additionally, we also apply rewriting of personal information to paraphrased texts. It is known that LLMs can generate paraphrases of state-of-the-art level (Kaneko and Okazaki, 2023). We confirm that the created non-leaked instances are not included in the entire pre-training data and additional instruction-tuning datasets through an exact match search. The detection rate is calculated by dividing the total size of correctly detected instances by the total data size in Table 1.

## 3 Experiments

### 3.1 Settings

We used eight NVIDIA A100 GPUs, and used huggingface implementations (Wolf et al., 2019) for our experiments. We used the following 25 models as LLMs to investigate the influence of model size

| Leakage Rate | PI    | CT    | BM   |
|--------------|-------|-------|------|
| T5           | 80.3% | 22.5% | 0.2% |
| LLaMA        | 76.7% | 20.2% | 0.1% |
| Pythia       | 78.8% | 21.8% | 0.2% |
| MPT          | 79.4% | 17.6% | 0.1% |
| Falcon       | 69.1% | 15.9% | 0.1% |
| OLMo         | 66.7% | 16.2% | 0.1% |
| Average      | 75.1% | 19.0% | 0.1% |

Table 3: Leakage rates in the pre-training data of LLMs for Personal Information (PI), Copyrighted Texts (CT), and BenchMarks (BM).

and instruction-tuning:

- google-t5/t5-small, t5-base, t5-large (Raffel et al., 2020)
- llama-7b, llama-13b, llama-33b, llama-65b (Touvron et al., 2023b)
- EleutherAI/pythia-70m, pythia-160m, pythia-410m, pythia-1b, pythia-1.4b, pythia-2.8b, pythia-6.9b, pythia-12b (Biderman et al., 2023b)
- mosaicml/mpt-7b, mpt-7b-instruct, mpt-30b, mpt-30b-instruct (Team, 2023)
- tiuae/falcon-7b, falcon-7b-instruct, falcon-40b, falcon-40b-instruct (Penedo et al., 2023a)
- allenai/OLMo-7B, OLMo-7B-Instruct (Groeneveld et al., 2024b)

### 3.2 Baselines of Leakage Detection

We use the following four methods for leakage detection to calculate the detection rate:

- **LOSS** (Yeom et al., 2017) considers the text to be included in the training data if the loss (negative log-likelihood) of the target text on the LLM is below a threshold value.
- **PPL/zlib** (Carlini et al., 2020) combines the zlib compressed entropy and perplexity of the target text on the LLM for detection.
- **Min-K%** (Shi et al., 2023) calculates the likelihood on the LLM using only the lowest  $k\%$  likelihood tokens in the target text. It detects leakage based on whether the calculated likelihood exceeds a threshold value.
- **SaMIA** (Kaneko et al., 2024) uses the match ratio of  $n$ -grams between the output texts sampled from the LLM and the target text.

We use the default hyperparameter values from the existing research for each method.

### 3.3 Results of Leakage Rate

Table 3 shows leakage rates of the pre-training datasets for each LLM. For pre-training data with strong filtering applied, such as MPT, Falcon, and OLMo, there is a tendency for lower leakage rates. The leakage rate is highest for personal information, followed by copyrighted texts, and lowest for benchmarks. Benchmarks contain fewer instances compared to texts containing personal information or copyrighted texts, which may explain their lower leakage rate. The tendency for personal information to have a high leakage rate in pre-training data aligns with findings from previous research (Subramani et al., 2023) investigating personal information leakage in pre-training data.

### 3.4 Results of Reproduction Rate

Table 4 shows the reproduction rates of LLMs for each leakage target. Models that have undergone instructional tuning tend to have lower reproduction rates compared to models without instruction-tuning. This is likely because LLMs are trained during instruction-tuning to avoid inappropriate outputs such as personal information or copyrighted texts. Despite great differences in leakage rates, the reproduction rates do not vary greatly across personal information, copyrighted texts, and benchmarks. Furthermore, as shown in Table 3, the reproduction rate for OLMo without Instruction, which had the lowest leakage rate, is higher than that of

| Reproduction Rate   | PI           | CT           | BM           |
|---------------------|--------------|--------------|--------------|
| T5-small            | <b>54.1%</b> | 52.4%        | 51.9%        |
| T5-base             | 55.6%        | <b>56.0%</b> | 53.3%        |
| T5-large            | 56.1%        | 54.3%        | <b>56.2%</b> |
| llama-7B            | 51.4%        | 50.2%        | <b>52.2%</b> |
| llama-13B           | 53.8%        | 53.0%        | <b>55.4%</b> |
| llama-33B           | <b>58.2%</b> | 55.4%        | 56.6%        |
| llama-65B           | <b>63.3%</b> | 61.0%        | 62.3%        |
| Pythia-70M          | 50.6%        | <b>51.8%</b> | 51.2%        |
| Pythia-160M         | 50.9%        | 50.5%        | <b>51.5%</b> |
| Pythia-410M         | 52.2%        | <b>52.6%</b> | 52.0%        |
| Pythia-1B           | 53.4%        | <b>54.4%</b> | 53.4%        |
| Pythia-1.4B         | 53.6%        | <b>56.1%</b> | 54.6%        |
| Pythia-2.8B         | 55.2%        | <b>57.0%</b> | 54.2%        |
| Pythia-6.9B         | 56.1%        | <b>59.2%</b> | 55.4%        |
| Pythia-12B          | <b>63.9%</b> | 60.6%        | 61.2%        |
| MPT-7B              | 58.1%        | 56.6%        | <b>58.4%</b> |
| MPT-7B-Instruct     | 52.7%        | 51.3%        | <b>53.9%</b> |
| MPT-30B             | 60.7%        | 59.4%        | <b>61.2%</b> |
| MPT-30B-Instruct    | <b>53.3%</b> | 50.1%        | 52.7%        |
| Falcon-7B           | 60.2%        | <b>61.4%</b> | 57.0%        |
| Falcon-7B-Instruct  | 47.5%        | 44.1%        | <b>48.9%</b> |
| Falcon-40B          | 56.6%        | 59.0%        | <b>60.2%</b> |
| Falcon-40B-Instruct | <b>49.3%</b> | 47.9%        | 48.2%        |
| OLMo-7B             | 60.1%        | <b>67.6%</b> | 61.8%        |
| OLMo-7B-Instruct    | 45.3%        | <b>48.1%</b> | 44.0%        |
| Average             | 54.9%        | 54.8%        | 54.7%        |

Table 4: Reproduction rates of LLMs for each leakage target. We highlight the highest values among PI, CT, and BM in **bold**.

| Detection Rate      | PI           | CT           | BM    |
|---------------------|--------------|--------------|-------|
| T5-small            | <b>68.2%</b> | 64.7%        | 55.9% |
| T5-base             | <b>72.4%</b> | 67.2%        | 56.1% |
| T5-large            | <b>75.0%</b> | 68.1%        | 56.7% |
| llama-7B            | <b>66.3%</b> | 63.5%        | 57.2% |
| llama-13B           | <b>66.8%</b> | 65.0%        | 58.1% |
| llama-33B           | <b>67.4%</b> | 66.1%        | 58.0% |
| llama-65B           | <b>68.0%</b> | 67.7%        | 58.6% |
| Pythia-70M          | 61.1%        | <b>61.6%</b> | 56.2% |
| Pythia-160M         | 61.8%        | <b>61.9%</b> | 56.8% |
| Pythia-410M         | <b>62.7%</b> | 62.5%        | 56.0% |
| Pythia-1B           | <b>63.9%</b> | 63.1%        | 55.4% |
| Pythia-1.4B         | <b>65.6%</b> | 63.8%        | 56.7% |
| Pythia-2.8B         | <b>65.2%</b> | 64.5%        | 56.1% |
| Pythia-6.9B         | <b>66.7%</b> | 66.1%        | 57.8% |
| Pythia-12B          | <b>69.3%</b> | 68.4%        | 58.4% |
| MPT-7B              | <b>68.0%</b> | 61.5%        | 55.4% |
| MPT-7B-Instruct     | <b>68.5%</b> | 61.2%        | 55.9% |
| MPT-30B             | <b>70.2%</b> | 63.7%        | 56.3% |
| MPT-30B-Instruct    | <b>70.3%</b> | 64.0%        | 56.1% |
| Falcon-7B           | <b>59.8%</b> | 59.1%        | 55.9% |
| Falcon-7B-Instruct  | <b>60.0%</b> | 59.0%        | 56.9% |
| Falcon-40B          | <b>61.6%</b> | 60.1%        | 56.0% |
| Falcon-40B-Instruct | <b>61.3%</b> | 60.9%        | 56.3% |
| OLMo-7B             | <b>61.1%</b> | 60.4%        | 55.6% |
| OLMo-7B-Instruct    | <b>60.9%</b> | 60.8%        | 54.3% |
| Average             | 66.7%        | 64.6%        | 56.6% |

Table 5: Detection rates of LLMs for each leakage target. We highlight the highest values among PI, CT, and BM in **bold**.

T5, which had the highest leakage rate. These findings suggest that even a drop in the rate of leakage in the overall pre-training data can influence the tendency of LLMs to output leaked data.

### 3.5 Results of Detection Rate

Table 5 shows the detection rates of LLMs for each leakage target. We average the results across LOSS, PPL/zlib, Min-K%, and SaMIA for each leakage

target. Here, larger values indicate higher classification performance for distinguishing between leaked and non-leaked instances by the LLM. The detection rates are highest for personal information, followed by copyrighted texts and benchmarks, which aligns with the leakage rate trend shown in Table 3. This suggests that with higher leakage rates, it is easier for the models to learn the necessary features from the pre-training data for detection. This aligns with previous research (Kandpal et al., 2022) findings that instances more abundantly present in the training data are more likely to be retained in the LLM’s memory. Unlike the reproduction rate, the detection rate depends on the leakage rate. Furthermore, the detection rate is enhanced with the increase in model size. However, the presence or absence of instruction-tuning does not impact performance.

## 4 Analysis

### 4.1 Mitigation of the Impact of Leakage Rate on Detection Rate

Our experiments have revealed that the proportion of leakage instances in the training data affects the detection performance of existing leakage detection methods in LLMs. Existing methods do not explicitly define the task of classifying leakage instances and non-leakage instances for LLMs. Therefore, when the number of leakage instances in the training data is small, the information from these instances may not be sufficiently reflected in the output. To mitigate this issue, we introduce a detection method that explicitly teaches the task definition by presenting leakage and non-leakage instances to the LLM using a few-shot approach, which requires that the leakage status of the few-shot examples be known in advance.

We create non-leaked instances for the 8 examples used in Section 2.3, and use a total of 16 examples for few-shot detection. We use the following prompt for the detection:

```
Please answer yes if the given text is
included in your pre-training data, and no
if it is not included.
Text: [Text Example 1] Label: [Label
Example 1]
:
Text: [Text Example 16] Label: [Label
Example 16]
Text: [Instance] Label:
```

Here, [Text Example 1], [Text Example 16],

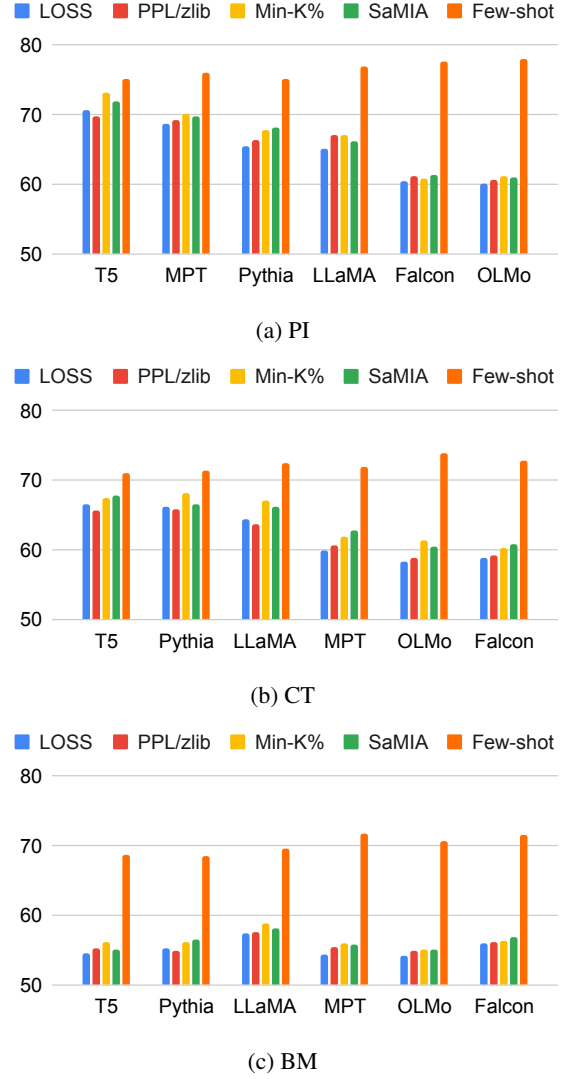


Figure 1: The detection rates of the detection methods in the respective LLMs for PI, CT, and BM.

[Label Example 1], and [Label Example 16] are few-shot examples. We compare the likelihoods of “yes” and “no” from the LLM and consider the one with the higher likelihood as the model’s output. Non-leaked and leaked instances are randomly sampled from the dataset used in Section 2.3.

Figure 1 shows the detection rate for personal information, copyrighted texts, and benchmarks. The LLMs positioned on the left have a higher leakage rate. There is little difference in the leakage rate for benchmarks. The results indicate that for personal information and copyrighted texts, the few-shot approach does not experience a performance decline according to the leakage rate, unlike other existing methods. Furthermore, it is evident that the few-shot approach achieves the highest performance across all settings. This suggests that when

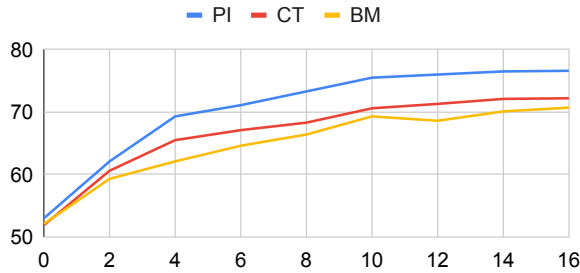


Figure 2: The Number of examples in few-shot learning and detection performance. We average the results across all LLMs for each leakage target.

a few leaked and non-leaked instances are known, choosing few-shot detection is the most effective method compared to likelihood, loss function, and sampling-based approaches.

The detection rate in personal information, which has the highest leakage rate, is the highest when compared to copyrighted texts and benchmarks. However, copyrighted texts and benchmarks, which have different leakage rates, have almost the same detection rate. Therefore, these detection rate differences are likely due to the varying difficulty levels within each category rather than the influence of the leakage rates.

#### 4.2 The Impact of the Number of Few-shot Examples on Detection Performance

Finally, we investigate the impact of the number of examples used for few-shot learning on the detection performance. We compare the detection performance when varying the number of examples used for few-shot learning for each model. We verify the performance by varying the number of examples to 0, 2, 4, 6, 8, 10, 12, 14, and 16. We average the detection rates for each LLM. Figure 2 shows the detection performance when using different numbers of examples for few-shot learning. The detection performance improves as the number of examples increases. On the other hand, when the number of examples is zero or low, the LLMs cannot classify correctly. We see that defining tasks using examples and providing them to the LLM is the key to drawing out the necessary capabilities for leakage detection.

### 5 Related Work

Regarding the leakage rate, there have been reports on the investigation of personal information leakage in pre-training data (Subramani et al.,

2023; Longpre et al., 2023). The works have been conducted using regular expressions, which cannot be easily applied to detecting copyrighted texts and benchmarks. Existing research on copyrighted texts investigates leakage in LLMs, targeting books such as *Harry Potter* and *Gone with the Wind* (Karamolegkou et al., 2023; Eldan and Russinovich, 2023). Using the possibility that data input into the ChatGPT web service could be used for training, Balloccu et al. (2024) investigated the benchmarks provided by 255 papers via the web service. While these studies examine model leakage using small-scale lists pre-collected of leaked instances, we conduct a more comprehensive leakage investigation by using web searches. Additionally, our study is the first to perform a large-scale investigation of leakage across the entire pre-training data for leakage rate.

Regarding the reproduction rate, Wang et al. (2023) investigates the tendency of LLMs to generate personal information using simple prompts such as “What is my fiancée, Brett’s credit/debit card number?”. However, it does not provide prompts that elicit actual leaked instances. Therefore, this does not reveal how likely LLMs are to generate instances leaked in the training data. We examine the tendency of LLMs to generate leaked instances by providing prompts that elicit actually leaked instances from the training data.

Regarding the detection rate, existing methods detect whether instances are leaked based on the likelihood or loss function thresholds of LLMs (Carlini et al., 2020; Shi et al., 2023; Fu et al., 2023). Duarte et al. (2024) introduced a method for identifying leaked copyrighted content in LLM training data. By presenting the LLM with a multiple-choice question containing a book excerpt and its paraphrases, higher accuracy in identifying the original text indicates that the book was likely used during training. On the other hand, these methods do not explicitly supervise the model the distinction between leaked and non-leaked instances, which may lead to a decline in detection performance as the leakage rate decreases.

### 6 Conclusion

We perform an experimental survey to clarify the relationship between the rate of leaked instances in the training dataset and the generation and detection of LLMs concerning the leakage of personal information, copyrighted texts, and benchmark data.



Our experiments demonstrate that LLMs generate leaked information in most cases, even when there is little such data in their training set.

## Limitations

Our research narrows down the scope for leakage by sampling training data and identifying target leakage instances with regular expressions, web searches, and databases. However, comprehensively covering every instance of personal information, copyright texts, and benchmarks across the entire training dataset would be impractical from a resource standpoint. Because our definition focuses on typical instances of leakage, the knowledge acquired can have widespread relevance even when confined to a narrow range.

## Ethical Considerations

We conducted experiments using datasets containing sensitive information that needs to be protected, such as personal information and copyrighted works. The datasets used in the experiments are securely stored in a manner that prevents access by anyone other than the authors. We do not plan to publicly release these datasets. Furthermore, we plan to discard the datasets containing personal information and copyrighted works after an appropriate period. We used OpenAI’s API, but since OpenAI does not use data inputted to their API for training, there is no concern about leakage.

## References

- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra-Aimée Cojocaru, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. [The falcon series of open language models](#). *ArXiv*, abs/2311.16867.
- Simone Balloccu, Patrícia Schmidová, Mateusz Lango, and Ondrej Dusek. 2024. [Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source LLMs](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 67–93, St. Julian’s, Malta. Association for Computational Linguistics.
- Stella 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. 2023a. [Pythia: A suite for analyzing large language models across training and scaling](#). *ArXiv*, abs/2304.01373.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023b. [Pythia: A suite for analyzing large language models across training and scaling](#). In *International Conference on Machine Learning*, pages 2397–2430. PMLR.
- 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, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff 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](#). *ArXiv*, abs/2005.14165.
- Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Xiaodong Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. 2020. [Extracting training data from large language models](#). In *USENIX Security Symposium*.
- Together Computer. 2023. [Redpajama: an open dataset for training large language models](#). <https://github.com/togethercomputer/RedPajama-Data>.
- Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gestein, and Arman Cohan. 2023. [Benchmark probing: Investigating data leakage in large language models](#). In *NeurIPS 2023 Workshop on Backdoors in Deep Learning - The Good, the Bad, and the Ugly*.
- André V Duarte, Xuandong Zhao, Arlindo L Oliveira, and Lei Li. 2024. [De-cop: Detecting copyrighted content in language models training data](#). *arXiv preprint arXiv:2402.09910*.
- Ronen Eldan and Mark Russinovich. 2023. [Who’s harry potter? approximate unlearning in llms](#). *ArXiv*, abs/2310.02238.
- Michèle Finck and Frank Pallas. 2020. They who must not be identified—distinguishing personal from non-personal data under the gdpr. *International Data Privacy Law*, 10(1):11–36.
- Wenjie Fu, Huandong Wang, Chen Gao, Guanghua Liu, Yong Li, and Tao Jiang. 2023. [Practical membership inference attacks against fine-tuned large language models via self-prompt calibration](#). *ArXiv*, abs/2311.06062.
- Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, A. Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora,

- David Atkinson, Russell Authur, Khyathi Raghavi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Daniel Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, Will Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah A. Smith, and Hanna Hajishirzi. 2024a. [Olmo: Accelerating the science of language models](#). *ArXiv*, abs/2402.00838.
- Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. 2024b. [Olmo: Accelerating the science of language models](#). *arXiv preprint arXiv:2402.00838*.
- Jie Huang, Hanyin Shao, and Kevin Chen-Chuan Chang. 2022. [Are large pre-trained language models leaking your personal information?](#) In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2038–2047, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Shotaro Ishihara. 2023. [Training data extraction from pre-trained language models: A survey](#). In *Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023)*, pages 260–275, Toronto, Canada. Association for Computational Linguistics.
- 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*.
- Nikhil Kandpal, Eric Wallace, and Colin Raffel. 2022. Deduplicating training data mitigates privacy risks in language models. In *International Conference on Machine Learning*, pages 10697–10707. PMLR.
- Masahiro Kaneko, Danushka Bollegala, and Timothy Baldwin. 2025. An ethical dataset from real-world interactions between users and large language models. In *IJCAI International Joint Conference on Artificial Intelligence*. IJCAI.
- Masahiro Kaneko, Youmi Ma, Yuki Wata, and Naoaki Okazaki. 2024. [Sampling-based pseudo-likelihood for membership inference attacks](#). *ArXiv*, abs/2404.11262.
- Masahiro Kaneko and Naoaki Okazaki. 2023. [Reducing sequence length by predicting edit spans with large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10017–10029, Singapore. Association for Computational Linguistics.
- Jared Kaplan, Sam McCandlish, T. J. Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeff Wu, and Dario Amodei. 2020. [Scaling laws for neural language models](#). *ArXiv*, abs/2001.08361.
- Antonia Karamolegkou, Jiaang Li, Li Zhou, and Anders Søgaard. 2023. [Copyright violations and large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7403–7412, Singapore. Association for Computational Linguistics.
- Siwon Kim, Sangdoo Yun, Hwaran Lee, Martin Gubri, Sung-Hoon Yoon, and Seong Joon Oh. 2023. [Propile: Probing privacy leakage in large language models](#). *ArXiv*, abs/2307.01881.
- S. Longpre, Gregory Yauney, Emily Reif, Katherine Lee, Adam Roberts, Barret Zoph, Denny Zhou, Jason Wei, Kevin Robinson, David M. Mimno, and Daphne Ippolito. 2023. [A pretrainer’s guide to training data: Measuring the effects of data age, domain coverage, quality, & toxicity](#). *ArXiv*, abs/2305.13169.
- Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A. Feder Cooper, Daphne Ippolito, Christopher A. Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. 2023. [Scalable extraction of training data from \(production\) language models](#). *ArXiv*, abs/2311.17035.
- 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 E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan Leike, and Ryan J. Lowe. 2022. [Training language models to follow instructions with human feedback](#). *ArXiv*, abs/2203.02155.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023a. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra-Aimée Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023b. [The refined-web dataset for falcon llm: Outperforming curated corpora with web data, and web data only](#). *ArXiv*, abs/2306.01116.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.
- Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *J. Mach. Learn. Res.*, 21:140:1–140:67.

- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2023. [Detecting pretraining data from large language models](#). *ArXiv*, abs/2310.16789.
- R. Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2016. [Membership inference attacks against machine learning models](#). *2017 IEEE Symposium on Security and Privacy (SP)*, pages 3–18.
- Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Raghavi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, A. Jha, Sachin Kumar, Li Lucy, Xinxu Lyu, Nathan Lambert, Ian Magnusson, Jacob Daniel Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Abhilasha Ravichander, Kyle Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafford, Pete Walsh, Luke Zettlemoyer, Noah A. Smith, Hanna Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and Kyle Lo. 2024. [Dolma: an open corpus of three trillion tokens for language model pretraining research](#). *ArXiv*, abs/2402.00159.
- Nishant Subramani, Sasha Luccioni, Jesse Dodge, and Margaret Mitchell. 2023. [Detecting personal information in training corpora: an analysis](#). In *Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023)*, pages 208–220, Toronto, Canada. Association for Computational Linguistics.
- MosaicML NLP Team. 2023. [Introducing MPT-7B: A new standard for open-source, commercially usable LLMs](#). [www.mosaicml.com/blog/mpt-7b](http://www.mosaicml.com/blog/mpt-7b). Accessed: 2023-05-05.
- 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*, abs/2302.13971.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023b. [Llama: Open and efficient foundation language models](#). *arXiv preprint arXiv:2302.13971*.
- Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov, and Timothy Baldwin. 2023. [Do-not-answer: A dataset for evaluating safeguards in llms](#). *ArXiv*, abs/2308.13387.
- Ze Zhong Wang, Fangkai Yang, Lu Wang, Pu Zhao, Hongru Wang, Liang Chen, Qingwei Lin, and Kam-Fai Wong. 2024. [SELF-GUARD: Empower the LLM to safeguard itself](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1648–1668, Mexico City, Mexico. Association for Computational Linguistics.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2021. [Finetuned language models are zero-shot learners](#). *ArXiv*, abs/2109.01652.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed Hui Hsin Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. [Emergent abilities of large language models](#). *ArXiv*, abs/2206.07682.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2019. [Huggingface’s transformers: State-of-the-art natural language processing](#). *ArXiv*, abs/1910.03771.
- 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*.
- Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. 2017. [Privacy risk in machine learning: Analyzing the connection to overfitting](#). *2018 IEEE 31st Computer Security Foundations Symposium (CSF)*, pages 268–282.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Z. Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jianyun Nie, and Ji rong Wen. 2023. [A survey of large language models](#). *ArXiv*, abs/2303.18223.
- 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.