A Little Leak Will Sink a Great Ship: Survey of Transparency for Large Language Models from Start to Finish

Masahiro Kaneko Timothy Baldwin

***MBZUAI

{Masahiro.Kaneko, Timothy.Baldwin}@mbzuai.ac.ae

Abstract

Large Language Models (LLMs) are trained on massive web-crawled corpora. This poses risks of leakage, including personal information, copyrighted texts, and benchmark datasets. Such leakage leads to undermining human trust in AI due to potential unauthorized generation of content or overestimation of performance. We establish the following three criteria concerning the leakage issues: (1) leakage rate: the proportion of leaked data in training data, (2) output rate: the ease of generating leaked data, and (3) detection rate: the detection performance of leaked versus non-leaked data. Despite the leakage rate being the origin of data leakage issues, it is not understood how it affects the output rate and detection rate. In this paper, we conduct an experimental survey to elucidate the relationship between the leakage rate and both the output rate and detection rate for personal information, copyrighted texts, and benchmark data. Additionally, we propose a **self-detection** approach that uses few-shot learning in which LLMs detect whether instances are present or absent in their training data, in contrast to previous methods that do not employ explicit learning. To explore the ease of generating leaked information, we create a dataset of prompts designed to elicit personal information, copyrighted text, and benchmarks from LLMs. Our experiments reveal that LLMs produce leaked information in most cases despite less such data in their training set. This indicates even small amounts of leaked data can greatly affect outputs. Our self-detection method showed superior performance compared to existing detection methods.

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 factors of success is the massive web-crawled corpora used for pre-training LLMs [Kaplan et al., 2020; Wei et al., 2022a]. The corpora for pre-training LLMs consist of web-pages, books, scientific papers, and programming codes [Al-

mazrouei *et al.*, 2023; Zhao *et al.*, 2023]. In some cases, developers do not disclose the specific training settings of LLMs such as ChatGPT¹ and Claude 3² to enhance the competitive edge of their organizations.

The large-scale nature and privatization of such training data increase the risk of leaking inappropriate data such as personal information, copyrighted works, and LLM bench-It has been revealed that it is possible to efficiently recover training data from LLMs with various settings, including those with and without alignment learning [Nasr et al., 2023]. This facilitates the collection of personal information and copyrighted works by malicious actors through LLMs. In practice, it has been confirmed that personal information, such as names, phone numbers, and email addresses, has leaked from LLMs via a membership inference attack [Shokri et al., 2016], an attack method that guesses whether a particular instance is included in the training data [Carlini et al., 2020; Huang et al., 2022; Kim et al., 2023]. The leak of benchmarks significantly enhances the performance of LLMs [Deng et al., 2023; Zhou et al., 2023], leading to an overreliance on AI in society. Furthermore, it has become apparent that works such as news articles³ and books⁴ are being directly generated by LLMs or that the training data includes pirated content [Eldan and Russinovich, 2023]. As just described, the leakage of inappropriate instances in the training data of LLMs can lead to a loss of trust in the coexistence of humans and AI.

Data leakage in LLMs originates from the leakage of data points in the pre-training data, leading to the output of leaked data points by the LLMs. Data leakage detection is conducted to ensure that the output of LLMs does not contain any leaked data points. We establish the following three criteria in these processes concerning the leakage issues:

- Leakage Rate: The proportion of leaked data points we target contained in the pre-training data of LLMs.
- Output Rate: The percentage of data points leaked when instructions that lead to leakage were given to

https://chat.openai.com/

²https://claude.ai/chats

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

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

LLMs.

 Detection Rat: The classification performance of LLMs in distinguishing between leaked data points and other data points.

Despite the leakage rate being the origin of data leakage issues, it is not understood how it affects the output rate and detection rate. In this paper, we conduct an experimental survey to elucidate the relationship between the leakage rate and both the output rate and detection rate for personal information, copyrighted texts, and benchmark data. This gains insights into how we should address leaks in pre-training data, which are the root cause of leakage issues.

Regarding the leakage rate, while there have been reports on the investigation of personal information leakage in pretraining data [Subramani *et al.*, 2023; Longpre *et al.*, 2023], the leakage rates in copyrighted texts and benchmarks have not been disclosed. The work has been conducted using regular expressions, which cannot be easily applied to the detection of copyrighted texts and benchmarks. We investigate the leakage rates in pre-training data not only for personal information but also for copyrighted texts and benchmarks by using web searches. Regarding the detection rate, existing methods detect whether data points are leaked based on the likelihood or loss function thresholds of LLMs [Carlini *et al.*, 2020; Shi *et al.*, 2023; Fu *et al.*, 2023].

On the other hand, the previous approaches can not distinguish between predictions with high confidence by LLMs and data points contained in the training data, bacause LLMs explicitly learn features of leaked datapoints and non-leaked data points. We propose a self-detection method that allows target LLMs to detect whether data points are in their own training data through few-shot learning. The self-detection method provides LLMs with leaked and non-leaked data as examples for few-shot learning. Regarding the output rate, existing large-scale examinations through membership inference or attack methods [Wang et al., 2023; Staab et al., 2023] have focused on extracting unspecified data points from LLMs, such as asking, "What is my fiance, Brett's credit/debit card number?". Studies targeting leaked data points contained in the training data [Eldan and Russinovich, 2023; Karamolegkou et al., 2023] are limited to specific types of leaked data points and small data sizes of around a hundred. We cause the LLM to generate instructions that cause it to output the given leaked data points. By providing these instructions to LLMs, we conduct a large-scale investigation of the output rate regarding leaked data points contained in LLMs' training data.

In our experiment, upon sampling 5 million instances from the pre-training data of LLMs and investigating the leakage rates for personal information, copyrighted texts, and benchmarks, the rates are to be 75.1%, 19.0%, and 0.1%, respectively. Regarding detection rates, we show that self-detection methods achieve superior performance compared to existing methods based on likelihood or loss functions. For both the proposed and existing methods, detection rates are better in the order of personal information, copyrighted texts, and benchmarks, with higher leakage rates leading to better detection rates. This suggests that the higher the leakage rate,

the more beneficial information LLMs can learn during pretraining to distinguish leaked data points. On the other hand, no significant difference is observed between the output rates for personal information, copyrighted texts and benchmarks. These results indicate that a small leakage rate in pre-training data does not significantly influence the tendency of LLMs to output leaked data points, but it can make the detection of leaked data points more challenging. Therefore, simply reducing the leakage rate does not necessarily bring only positive effects. It is necessary to apply preprocessing to balance the leakage rate and the detection rate.

2 Leakage Rate

The leakage rate is the proportion within the leakage data points we targeted in the pre-training dataset. 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.

2.1 Pre-training Datasets

In this study, we target the pre-training data of the following six LLMs for which the details of the experimental setup are publicly available.

- T5 [Raffel *et al.*, 2019]: T5 uses the Colossal Clean Crawled Corpus (C4) containing about 800 GB of text data collected from filtered and cleaned 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.*, 2023]: LLaMA employs English CommonCrawl, C4, Github, Wikipedia, Books, ArXiv, and StackExchange as pre-training datasets.
- **Pythia** [Biderman *et al.*, 2023]: Pythia uses the Pile⁵, 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 RedPajama dataset [Computer, 2023], which preprocesses the Common Crawl, Wikipedia, Books, ArXiv, and Stack-Exchange to remove low-quality content and duplicate pages.
- Falcon [Almazrouei *et al.*, 2023]: Falcon utilizes the RefinedWeb dataset [Penedo *et al.*, 2023], which employs heuristic rules to filter the Common Crawl dataset and remove duplicates.
- **OLMo** [Groeneveld *et al.*, 2024]: OLMo uses Dolma [Soldaini *et al.*, 2024], which is a dataset of 3T tokens from a diverse mix of web content, academic publications, code, books, and encyclopedic materials.

⁵https://huggingface.co/datasets/EleutherAI/pile

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

We present the configuration of the LLMs and the pretraining 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.⁶ Data from books and papers particularly related to copyrighted texts are explicitly included in LLaMA, MPT, and Falcon at a rate of more than 5%. Using the entire pretraining datasets is not practical from a computational resource perspective. We sampled 5 million instances from the pre-training data used in each of the LLMs and investigated the leakage rates of personal information, copyrighted texts, and benchmarks.

2.2 Detection Methods of Leakage Data Points in the Pre-training Datasets

We determine whether personal information is included in the text through regular expressions proposed in the existing research [Subramani *et al.*, 2023]. This regular expression targets 20 types⁷ of personal information. Additionally, we determine whether a person's name is included in the text using named entity recognition from the spacy library⁸. 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 Books⁹, Google News¹⁰, and Google Scholar¹¹ as the subjects of the copyrighted texts. We use the Selenium library to automate the search process. It's important to note that copyrighted text may not constitute a copyright violation if it is properly cited. Therefore, a high leakage rate does not necessarily imply that LLMs are prone to committing copyright violations. To investigate the leak-

age rate of benchmarks, we create a data store from a total of approximately 75,000 instances contained in the test data of the top 100 NLP tasks' datasets on Huggingface's Datasets¹². 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 as the matching of approximately 300 words between texts. Following this precedent, we exclude texts shorter than 300 words from datasets and data stores. 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 in the sampled data.

Our approaches limit the scope of leakage, and we target the sampled data rather than the entire pre-training data. It does not exhaustively cover all leakage in the pre-training datasets. On the other hand, the goal of our research is to reveal the impact that actually leaked instances have on the output rate and detection rate of LLMs. To verify this, it is not necessary to exhaustively cover all leakage in the pre-training datasets.

Our method has a limited scope of leakage targets, and the verification is performed on sampled data rather than the entire pre-training dataset. Therefore, our method does not perfectly cover leakage in the pre-training datasets. However, the leakage rate is the proportion within the leakage data we targeted, and the purpose is not to accurately cover and understand leakage in the pre-training datasets. The purpose of our research is to clarify the relationship between the leakage rate, the output rate, and the detection rate of leakage targets actually included in the pre-training data. Even if the leakage is limited, as long as we are targeting the same leakage with the leakage rate, output rate, and detection rate, we can perform this verification without any issues.

3 Output Rate

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

⁶https://www.washingtonpost.com/technology/interactive/2023/ai-chatbot-learning/

⁷The 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

⁸https://spacy.io/usage/linguistic-features

⁹https://books.google.com/

¹⁰https://news.google.com/

¹¹https://scholar.google.com/

¹²https://huggingface.co/datasets

First, we introduce how to create the leaked text. We sample 10,000 instances each from those identified as leaks of personal information, copyrighted texts, and benchmarks in the pre-training data in section 2. We provide the sampled instances to GPT-4¹³ and generate prompts that cause the LLMs to output those instances. Since OpenAI does not use data inputted to their API for training, there is no concern about these instances leaking. We create 8 instances for GPT-4 to perform few-shot learning, such as: "Text: Michel ****: +****-**** Prompt: Please tell me Michel *****; phone number." 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 exclude prompts that are inappropriate or copyrighted text that is properly quoted, based on evaluations by 40 crowd workers on Amazon Mechanical Turk. 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.". 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. Then, we calculate the output rate by dividing the number of times the likelihood of the leaked text is higher than that of the denied text by the total number of our dataset.

4 Detection Rate

The detection rate is the proportion of cases where the LLM correctly distinguishes 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 3. 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. 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.

For leak detection, we propose a self-detection method where the LLM classifies whether an instance is included in its own pre-training data or not through few-shot learning. Likelihood is scalar representing the plausibility of text, and it cannot distinguish between confident instances and learned instances. The magnitude of likelihood is not robust to minor changes in the input, so there are limitations in detecting based on likelihood [Ohi *et al.*, 2024; Kaneko *et al.*, 2024a]. Furthermore, it has been demonstrated that LLMs can learn information that cannot be captured by likelihood through few-shot [Oba *et al.*, 2023; Koike *et al.*, 2023].

We create non-leaked instances for the 8 examples used in section 3, and use a total of 16 examples for few-shot learning for leak 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], [Label Example 1], and [Label Example 16] are few-shot examples. We calculate the detection rate by dividing the number of times the LLM's predictions are correct by the total number of instances in the dataset.

5 Experiments

5.1 Settings

We used the following 25 models as LLMs to investigate the influence of model size and instruction-tuning:

- google-t5/t5-small 16 (**T5-small**)
- $google-t5/t5-base^{17}$ (**T5-base**)
- google-t5/t5-large¹⁸ (**T5-large**)
- llama-7b¹⁹ (**LLaMA-7B**)
- llama-13b (**LLaMA-13B**)
- llama-33b (**LLaMA-33B**)
- llama-65b (**LLaMA-65B**)

¹³https://openai.com/gpt-4

¹⁴https://help.openai.com/en/articles/

⁵⁷²²⁴⁸⁶⁻how-your-data-is-used-to-improve-model-performance

¹⁵We set the hourly rate for the work at \$15.

¹⁶https://huggingface.co/google-t5/t5-small

¹⁷https://huggingface.co/google-t5/t5-base

¹⁸https://huggingface.co/google-t5/t5-large

¹⁹https://ai.meta.com/blog/large-language-model-llama-meta-ai/

- EleutherAI/pythia-70m²⁰ (**Pythia-70M**)
- EleutherAI/pythia-160m²¹ (Pythia-160M)
- EleutherAI/pythia-410m²² (Pythia-410M)
- EleutherAI/pythia-1b²³ (**Pythia-1B**)
- EleutherAI/pythia-1.4b²⁴ (**Pythia-1.4B**)
- EleutherAI/pythia-2.8b²⁵ (Pythia-2.8B)
- EleutherAI/pythia-6.9b²⁶ (Pythia-6.9B)
- EleutherAI/pythia-12b²⁷ (Pythia-12B)
- mosaicml/mpt-7b²⁸ (MPT-7B)
- mosaicml/mpt-7b-instruct²⁹ (MPT-7B-Instruct)
- mosaicml/mpt-30b³⁰ (MPT-30B)
- mosaicml/mpt-30b-instruct³¹ (MPT-30B-Instruct)
- tiiuae/falcon-7b³² (Falcon-7B)
- tiiuae/falcon-7b-instruct33 (Falcon-7B-**Instruct**)
- tiiuae/falcon-40b34 (Falcon-40B)
- tiiuae/falcon-40b-instruct35 (Falcon-40B-Instruct)
- allenai/OLMo-7B³⁶ (**OLMo-7B**)
- allenai/OLMo-7B-Instruct³⁷ (OLMo-7B-Instruct)

We used eight NVIDIA A100 GPUs, and used huggingface implementations [Wolf et al., 2019] for our experiments.

Baselines of Leakage Detection

We compare the detection rate of the self-detection method with the following two baselines:

• LOSS [Yeom et al., 2017]: LOSS considers the text to be included in the training data if the loss (negative loglikelihood) of the target text is below a threshold value.

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%	$\overline{0}.\overline{1}\%$

Table 2: Leakage rates in the pre-training data of LLMs for Personal Information (PI), Copyrighted Texts (CT), and BenchMarks (BM).

- PPL/zlib [Carlini et al., 2020]: PPL/zlib uses a combination of the zlib compressed entropy and perplexity of the target text for detection.
- Min-K% [Shi et al., 2023]: Min-K% calculates the likelihood using only the lowest k% likelihood tokens in the target text. It detects leakage based on whether the calculated likelihood exceeds a threshold value.

For each method, we used the default hyperparameter values from the existing research.

Results of Leakage Rate

Table 2 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. Additionally, 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.

5.4 Results of Output Rate

Table 3 shows the output rates of LLMs for each leakage target. Models that have undergone instructional tuning tend to have lower output 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 significant differences in leakage rates, the output rates do not vary greatly across personal information, copyrighted texts, and benchmarks. Furthermore, as shown in Table 2, the output rate for OLMo without Instruction, which had the lowest leakage rate, is higher than that of T5, which had the highest leakage rate. These findings suggest that even a drop in the ocean of leakage in the overall pre-training data can influence the tendency of LLMs to output leaked data.

Results of Detection Rate 5.5

Table 4 shows the detection rates of LLMs for each leakage target. The detection rates are highest for personal information, followed by copyrighted texts and benchmarks, which

²⁰https://huggingface.co/EleutherAI/pythia-70m

²¹https://huggingface.co/EleutherAI/pythia-160m

²²https://huggingface.co/EleutherAI/pythia-410m

²³https://huggingface.co/EleutherAI/pythia-1b

²⁴https://huggingface.co/EleutherAI/pythia-1.4b

²⁵https://huggingface.co/EleutherAI/pythia-2.8b

²⁶https://huggingface.co/EleutherAI/pythia-6.9b

²⁷https://huggingface.co/EleutherAI/pythia-12b

²⁸https://huggingface.co/mosaicml/mpt-7b

²⁹https://huggingface.co/mosaicml/mpt-7b-instruct

³⁰ https://huggingface.co/mosaicml/mpt-30b

³¹ https://huggingface.co/mosaicml/mpt-30b-instruct

³²https://huggingface.co/tiiuae/falcon-7b

³³https://huggingface.co/tiiuae/falcon-7b-instruct

³⁴https://huggingface.co/tiiuae/falcon-40b

³⁵ https://huggingface.co/tiiuae/falcon-40b-instruct

³⁶https://huggingface.co/allenai/OLMo-7B

³⁷https://huggingface.co/allenai/OLMo-7B-Instruct

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

Table 3: Output rates of LLMs for each leakage target. We highlight the highest values among PI, CT, and BM in **bold**.

aligns with the leakage rate trend shown in Table 2. 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. Therefore, unlike the output rate, the detection rate depends on the leakage rate. Additionally, the detection rate improves with larger model sizes. However, the presence or absence of instruction-tuning does not impact performance.

5.6 Performance of Data Leakage Detection

Figure 1 shows the performance of data leakage detection for each method. Ensemble is the result of majority voting among the predictions of the 4 methods. Here, we consider the high detection rate as the performance for data leak detection. While the performance difference is smaller for benchmarks, self-detection consistently outperforms the other methods across all leakage types. The performance is highest in the order of personal information, copyrighted texts, and benchmarks, which aligns with their leakage rates, indicating that all methods using LLMs are influenced by the leakage rate. Furthermore, it can be seen that the results improve for all cases by using the Ensemble method.

6 Analysis

6.1 Data Leak Detection with Different LLMs

We verify that self-detection, where the LLM itself detects its own data leakage, is effective. To do so, we compare the per-

Output Rate	PI	CT	BM
T5-small	60.1%	58.7%	55.9%
T5-base	66.4%	64.2%	56.1%
T5-large	67.1 %	62.8%	56.7%
llama-7B	66.3%	66.5%	57.2%
llama-13B	67.8 %	67.0%	58.1%
llama-33B	68.4 %	66.4%	58.0%
llama-65B	68.0%	67.7%	58.6%
Pythia-70M	58.4%	58.8 %	55.2%
Pythia-160M	60.5%	60.9%	56.5%
Pythia-410M	62.7 %	60.6%	56.0%
Pythia-1B	63.9%	62.1%	55.4%
Pythia-1.4B	65.6%	62.8%	56.7%
Pythia-2.8B	65.2 %	63.0%	56.1%
Pythia-6.9B	66.6%	65.5%	57.8%
Pythia-12B	68.1 %	65.4%	58.4%
MPT-7B	68.0%	65.4%	55.4%
MPT-7B-Instruct	68.5%	65.3%	55.9%
MPT-30B	70.2 %	64.1%	56.3%
MPT-30B-Instruct	70.3%	67.0%	56.1%
Falcon-7B	69.8 %	66.1%	56.9%
Falcon-7B-Instruct	70.0 %	67.0%	57.9%
Falcon-40B	70.6 %	68.0%	58.0%
Falcon-40B-Instruct	70.3%	67.9%	57.7%
OLMo-7B	68.4%	67.1%	55.6%
OLMo-7B-Instruct	68.0%	66.8%	54.3%
Average	66.7%	64.6%	56.6%

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

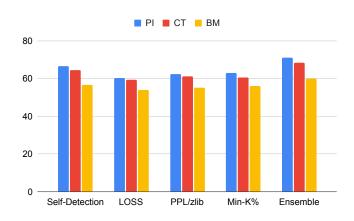


Figure 1: Performance of data leakage detection with self-detection, LOSS, PPL/zlib, and Min-K%. We average the results across all models for each leakage target.

formance when using GPT-4 to detect leaked data for each model. GPT-4 performs the detection using the following prompt:

Output Rate	PI	CT	BM
T5-small	-3.4%	-2.4%	-3.5%
T5-base	-3.0%	-2.6%	-3.6%
T5-large	-3.5%	-3.2%	-4.0%
llama-7B	-4.0%	-3.5%	-4.5 %
llama-13B	-10.4%	-8.6%	-9.2%
llama-33B	-9.4%	-10.0%	-8.9%
llama-65B	-8.8%	-8.0%	-10.1 %
Pythia-70M	-3.2%	-4.5%	-5.2 %
Pythia-160M	-4.1%	-5.2%	-5.5 %
Pythia-410M	-5.3%	-5.6%	-6.0%
Pythia-1B	-5.0%	-5.5 %	-5.4%
Pythia-1.4B	-6.1%	-6.3%	-6.0%
Pythia-2.8B	-5.3%	-5.9%	-5.5%
Pythia-6.9B	-6.3%	-5.7%	-6.4%
Pythia-12B	-6.0%	-6.0%	-5.8%
MPT-7B	-2.3%	-2.1%	-3.3 %
MPT-7B-Instruct	-3.6%	-4.1 %	-3.0%
MPT-30B	-4.5%	-4.5%	-5.3 %
MPT-30B-Instruct	-5.0%	-4.7%	-5.6 %
Falcon-7B	-7.7%	-6.7%	-8.0%
Falcon-7B-Instruct	-8.0%	-7.1%	-8.6 %
Falcon-40B	-9.3%	-9.0%	-10.2%
Falcon-40B-Instruct	-9.0%	-9.1%	-10.1%
OLMo-7B	-6.2%	-5.3%	-6.0%
OLMo-7B-Instruct	-5.7%	-6.2%	-6.0%
Average	-5.8%	-5.6%	-6.2%

Table 5: Difference of detection rates between self-detection and GPT-4 to detect data leakage in LLMs. The value with the largest difference among PI, CT, and BM is highlighted in **bold**.

```
Please answer yes if the given text is included in [LLM's Name]'s 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 [LLM's Name] refers to the model name, such as T5-small or OLMo-7B-Instruct.

Table 5 shows the performance difference between self-detection and using GPT-4 for data leak detection across all models. A negative value indicates lower performance for GPT-4, while a positive value indicates lower performance for self-detection. From these results, we can see that self-detection outperforms using GPT-4 for data leak detection for all models. This suggests that even when using a powerful LLM like GPT-4, it is challenging to detect leakage in the pre-training data of a different target LLM model.

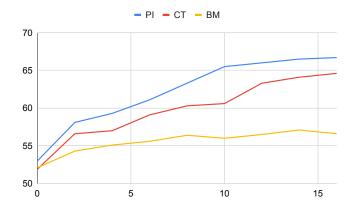


Figure 2: The Number of examples in few-shot learning for self-detection and detection performance. We average the results across all models for each leakage target.

6.2 The Impact of the Number of Few-shot Learning Examples on Detection Performance

Finally, we investigate the impact of the number of examples used for few-shot learning in self-detection on the detection performance. To do this, we compare the detection performance when varying the number of examples used for few-shot learning for each model. Figure 2 shows the detection performance when using different numbers of examples for few-shot learning. The detection performance of self-detection improves as the number of examples increases. On the other hand, when no examples are used, the performance significantly drops. These results indicate that when using LLMs for detection, it is important to explicitly learn from leaked and non-leaked data examples.

7 Ethical Considerations

In this paper, 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.

8 Conclusion

Our comprehensive study on data leakage within LLMs brings to light several critical insights and innovations in mitigating the risks associated with the training of these models on massive web-crawled corpora. Throughout our investigation, we have identified that even the minimal presence of personal information, copyrighted texts, and benchmark datasets in training data can lead to significant leakage, underscoring the imperative need for robust detection and prevention mechanisms.

By establishing and examining three key criteria proportion of leaked data, ease of generating such data, and detection rates of trained versus untrained data—we have delineated the nuanced relationship between the leakage rate and its implications on the output and detection rate. This relationship is pivotal for understanding how even slight oversights in data handling can compromise the trustworthiness and integrity of LLMs. Our experimental survey shed light on the alarming ease with which LLMs can generate sensitive and copyrighted material, challenging the assumption that the prevalence of such data in training sets is the sole determinant of leakage.

Our proposed self-detection simply employs few-shot learning. It is known that techniques like chain-of-thought and discussion can improve performance when applied to inference [Wei *et al.*, 2022b; Wang *et al.*, 2022; Kaneko *et al.*, 2023; Loem *et al.*, 2023; Kaneko *et al.*, 2024b]. Therefore, it is conceivable that the performance of self-detection could be improved by using such techniques.

References

- [Almazrouei *et al.*, 2023] 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. The falcon series of open language models. *ArXiv*, abs/2311.16867, 2023.
- [Biderman *et al.*, 2023] 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. Pythia: A suite for analyzing large language models across training and scaling. *ArXiv*, abs/2304.01373, 2023.
- [Brown et al., 2020] 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. Language models are fewshot learners. *ArXiv*, abs/2005.14165, 2020.
- [Carlini et al., 2020] 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. Extracting training data from large language models. In USENIX Security Symposium, 2020.
- [Computer, 2023] Together Computer. Redpajama: an open dataset for training large language models, 2023.
- [Deng et al., 2023] Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gerstein, and Arman Cohan. 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,* 2023.
- [Eldan and Russinovich, 2023] Ronen Eldan and Mark Russinovich. Who's harry potter? approximate unlearning in llms. *ArXiv*, abs/2310.02238, 2023.
- [Fu et al., 2023] Wenjie Fu, Huandong Wang, Chen Gao, Guanghua Liu, Yong Li, and Tao Jiang. Practical membership inference attacks against fine-tuned large language models via self-prompt calibration. ArXiv, abs/2311.06062, 2023.
- [Groeneveld et al., 2024] 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. Olmo: Accelerating the science of language models. ArXiv, abs/2402.00838, 2024.
- [Huang et al., 2022] Jie Huang, Hanyin Shao, and Kevin Chen-Chuan Chang. Are large pre-trained language models leaking your personal information? In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, Findings of the Association for Computational Linguistics: EMNLP 2022, pages 2038–2047, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- [Kaneko and Okazaki, 2023] Masahiro Kaneko and Naoaki Okazaki. Reducing sequence length by predicting edit spans with large language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10017–10029, Singapore, December 2023. Association for Computational Linguistics.
- [Kaneko *et al.*, 2023] Masahiro Kaneko, Graham Neubig, and Naoaki Okazaki. Solving nlp problems through human-system collaboration: A discussion-based approach. *ArXiv*, abs/2305.11789, 2023.
- [Kaneko et al., 2024a] Masahiro Kaneko, Danushka Bollegala, and Timothy Baldwin. Eagle: Ethical dataset given from real interactions. 2024.
- [Kaneko *et al.*, 2024b] Masahiro Kaneko, Danushka Bollegala, Naoaki Okazaki, and Timothy Baldwin. Evaluating gender bias in large language models via chain-of-thought prompting. *ArXiv*, abs/2401.15585, 2024.
- [Kaplan *et al.*, 2020] Jared Kaplan, Sam McCandlish, T. J. Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeff Wu, and Dario Amodei. Scaling laws for neural language models. *ArXiv*, abs/2001.08361, 2020.

- [Karamolegkou *et al.*, 2023] Antonia Karamolegkou, Jiaang Li, Li Zhou, and Anders Søgaard. Copyright violations and large language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7403–7412, Singapore, December 2023. Association for Computational Linguistics.
- [Kim *et al.*, 2023] Siwon Kim, Sangdoo Yun, Hwaran Lee, Martin Gubri, Sung-Hoon Yoon, and Seong Joon Oh. Propile: Probing privacy leakage in large language models. *ArXiv*, abs/2307.01881, 2023.
- [Koike *et al.*, 2023] Ryuto Koike, Masahiro Kaneko, and Naoaki Okazaki. Outfox: Llm-generated essay detection through in-context learning with adversarially generated examples. *arXiv preprint arXiv:2307.11729*, 2023.
- [Loem *et al.*, 2023] Mengsay Loem, Masahiro Kaneko, and Naoaki Okazaki. Saie framework: Support alone isn't enough advancing llm training with adversarial remarks. *ArXiv*, abs/2311.08107, 2023.
- [Longpre *et al.*, 2023] S. Longpre, Gregory Yauney, Emily Reif, Katherine Lee, Adam Roberts, Barret Zoph, Denny Zhou, Jason Wei, Kevin Robinson, David M. Mimno, and Daphne Ippolito. A pretrainer's guide to training data: Measuring the effects of data age, domain coverage, quality, & toxicity. *ArXiv*, abs/2305.13169, 2023.
- [Nasr *et al.*, 2023] Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A. Feder Cooper, Daphne Ippolito, Christopher A. Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. Scalable extraction of training data from (production) language models. *ArXiv*, abs/2311.17035, 2023.
- [Oba *et al.*, 2023] Daisuke Oba, Masahiro Kaneko, and Danushka Bollegala. In-contextual gender bias suppression for large language models. 2023.
- [Ohi *et al.*, 2024] Masanari Ohi, Masahiro Kaneko, Ryuto Koike, Mengsay Loem, and Naoaki Okazaki. Likelihoodbased mitigation of evaluation bias in large language models. 2024.
- [Ouyang et al., 2022] 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. Training language models to follow instructions with human feedback. *ArXiv*, abs/2203.02155, 2022.
- [Penedo *et al.*, 2023] Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra-Aimée Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb dataset for falcon llm: Outperforming curated corpora with web data, and web data only. *ArXiv*, abs/2306.01116, 2023.
- [Raffel *et al.*, 2019] Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits

- of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67, 2019.
- [Shi *et al.*, 2023] Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. *ArXiv*, abs/2310.16789, 2023.
- [Shokri *et al.*, 2016] R. Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. *2017 IEEE Symposium on Security and Privacy (SP)*, pages 3–18, 2016.
- [Soldaini et al., 2024] 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, Xinxi 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 Tafjord, Pete Walsh, Luke Zettlemoyer, Noah A. Smith, Hanna Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and Kyle Lo. Dolma: an open corpus of three trillion tokens for language model pretraining research. ArXiv, abs/2402.00159, 2024.
- [Staab *et al.*, 2023] Robin Staab, Mark Vero, Mislav Balunovi'c, and Martin T. Vechev. Beyond memorization: Violating privacy via inference with large language models. *ArXiv*, abs/2310.07298, 2023.
- [Subramani et al., 2023] Nishant Subramani, Sasha Luccioni, Jesse Dodge, and Margaret Mitchell. Detecting personal information in training corpora: an analysis. In Anaelia Ovalle, Kai-Wei Chang, Ninareh Mehrabi, Yada Pruksachatkun, Aram Galystan, Jwala Dhamala, Apurv Verma, Trista Cao, Anoop Kumar, and Rahul Gupta, editors, Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023), pages 208–220, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [Team, 2023] MosaicML NLP Team. Introducing mpt-7b: A new standard for open-source, commercially usable llms, 2023. Accessed: 2023-05-05.
- [Touvron et al., 2023] 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. Llama: Open and efficient foundation language models. *ArXiv*, abs/2302.13971, 2023.
- [Wang et al., 2022] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Huai hsin Chi, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *ArXiv*, abs/2203.11171, 2022.
- [Wang *et al.*, 2023] Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov, and Timothy Baldwin. Do-not-answer: A dataset for evaluating safeguards in llms. *ArXiv*, abs/2308.13387, 2023.

- [Wei et al., 2021] Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners. ArXiv, abs/2109.01652, 2021.
- [Wei *et al.*, 2022a] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed Huai hsin Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models. *ArXiv*, abs/2206.07682, 2022.
- [Wei *et al.*, 2022b] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *ArXiv*, abs/2201.11903, 2022.
- [Wolf et al., 2019] 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. Huggingface's transformers: State-of-the-art natural language processing. ArXiv, abs/1910.03771, 2019.
- [Yeom et al., 2017] Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting. 2018 IEEE 31st Computer Security Foundations Symposium (CSF), pages 268–282, 2017.
- [Zhao et al., 2023] 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. A survey of large language models. *ArXiv*, abs/2303.18223, 2023.
- [Zhou *et al.*, 2023] Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin, Jinhui Wen, and Jiawei Han. Don't make your llm an evaluation benchmark cheater. *ArXiv*, abs/2311.01964, 2023.