

---

# Unlocking Post-hoc Dataset Inference with Synthetic Data

---

Bihe Zhao<sup>1</sup> Pratyush Maini<sup>2,3</sup> Franziska Boenisch<sup>1</sup> Adam Dziedziec<sup>1</sup>

## Abstract

The remarkable capabilities of large language models stem from massive internet-scraped training datasets, often obtained without respecting data owners’ intellectual property rights. Dataset Inference (DI) enables data owners to verify unauthorized data use by identifying whether a suspect dataset was used for training. However, current DI methods require private held-out data with a distribution that closely matches the compromised dataset. Such held-out data are rarely available in practice, severely limiting the applicability of DI. In this work, we address this challenge by synthetically generating the required held-out set through two key contributions: (1) creating high-quality, diverse synthetic data via a data generator trained on a carefully designed suffix-based completion task, and (2) bridging likelihood gaps between real and synthetic data, which is realized through post-hoc calibration. Extensive experiments on diverse text datasets show that using our generated data as a held-out set enables DI to detect the original training sets with high confidence, while maintaining a low false positive rate. This result empowers copyright owners to make legitimate claims on data usage and demonstrates our method’s reliability for real-world litigations.

## 1. Introduction

Large language models (LLMs) have recently achieved remarkable success in a broad range of tasks, fueled by the availability of massive high-quality text corpora often scraped from the internet (Weber et al., 2024; Penedo et al., 2024). While this practice has enabled LLMs to generate high-quality text and to excel on benchmarks, it

also raises serious concerns related to intellectual property rights (Reuters, 2023; Gry, 2023; Sil, 2023), data privacy, and transparency (Rahman & Santacana, 2023; Wu et al., 2023). The reliance on potentially unauthorized data creates an urgent need for methods that allow independent authors to verify whether a given dataset has been used to train an LLM without the explicit consent of the model provider.

A promising approach to addressing these concerns is *dataset inference* (DI) (Maini et al., 2021; Dziedziec et al., 2022; Maini et al., 2024; Dubiński et al., 2024), which aims to determine whether a suspect dataset has contributed to a model’s training. This puts power in the hands of data owners to monitor and exercise their intellectual property rights. Despite its potential, DI currently faces a critical bottleneck: it requires a held-out set—a dataset known to be absent from training—that shares the same distribution as the suspect dataset (Zhang et al., 2024a). In practice, however, such an in-distribution held-out set is rarely available. Data creators do not typically reserve a dedicated held-out set for auditing purposes, and any disclosed held-out data could itself be repurposed for future training. Moreover, even when a dataset owner can provide held-out data, any slight distributional shifts from the original suspect data can undermine DI by inflating false positives (Das et al., 2024; Duan et al., 2024; Meeus et al., 2024; Maini & Suri, 2024).

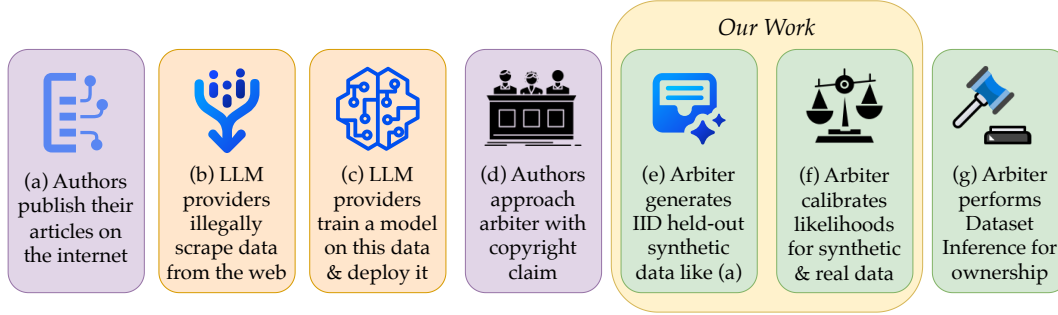
To illustrate the brittleness of using seemingly IID (Independent and Identically Distributed) held-out data, we demonstrate in Section 3 that even in a simple scenario—where an LLM is fine-tuned on blog posts from a *single* author—there exists a distributional shift between training data (members) and randomly held-out blog posts from the same author. This highlights how even subtle variations in held-out data can undermine DI. Malicious actors may exploit this vulnerability by strategically introducing *shifted* held-out data, falsely accusing model owners of copyright infringement and further reducing the reliability of DI methods.

In this work, we address these challenges by proposing to *synthetically generate* held-out data for DI, bypassing the need for in-distribution held-out data. This vision, however, is non-trivial to achieve. First, the generated texts must be realistic, high-quality, and sufficiently diverse to approximate the distribution of the original data. Second, the generation process itself may introduce a distribution shift between nat-

---

<sup>1</sup>CISPA Helmholtz Center for Information Security <sup>2</sup>Carnegie Mellon University <sup>3</sup>DatologyAI. Correspondence to: Bihe Zhao <bihe.zhao@cispa.de>, Pratyush Maini <pratyush-maini@cmu.edu>, Franziska Boenisch <boenisch@cispa.de>, Adam Dziedziec <adam.dziedziec@cispa.de>.

Published at Data in Generative Models Workshop: The Bad, the Ugly, and the Greats (DIG-BUGS) at ICML 2025, Vancouver, Canada. Copyright 2025 by the author(s).



**Figure 1. Dataset Inference Procedure with Synthetic Held-Out Data.** This figure presents a high-level overview of how the proposed dataset inference (DI) process will take place in real-world use cases. While original setting of DI requires a held-out dataset that is IID to the suspect set ((a-d),(g)), we highlight the contributions of **our work**: (e) The arbiter generates IID synthetic held-out data that mimics the author’s original data. (f) The arbiter calibrates likelihoods between real and synthetic data to ensure fair comparison, enabling them to reliably perform dataset inference.

ural and synthetic held-out data. Such a shift complicates DI: if a difference is observed between the suspect and held-out sets, it becomes unclear whether this difference arises from a genuine membership signal (*i.e.*, the target model behaves differently on the suspect data because it has seen it during training) or merely from the distribution shift (*i.e.*, the model behaves differently on suspect data because it is natural data). Recent studies have extensively highlighted this issue in the context of Membership Inference Attacks (MIAs) (Shokri et al., 2017), where distribution shifts lead to misleading evaluation results (Das et al., 2024; Zhang et al., 2024a; Maini et al., 2024; Dubiński et al., 2024).

To this end, we first train a carefully designed text generator on the suspect dataset itself, on a suffix completion task (Section 4.1). This approach produces high-quality datasets with only a small distributional shift from the suspect texts. However, even small shifts in distribution can undermine DI’s reliability. To address this, we introduce a *post-hoc* calibration step (Section 4.2) to ensure that the generated held-out set can serve as a reliable reference for DI. Specifically, we disentangle the effects of distributional shifts from the actual membership signal—a critical factor in DI. To achieve this, we propose a dual-classifier approach: (1) A *text-only classifier*, trained to distinguish natural (original) from generated data. (2) A *membership-aware classifier*, which incorporates both the textual features and DI’s standard membership indicators. The key insight is that any performance advantage of the membership-aware classifier over the text-only classifier must arise from the presence of membership signals rather than distributional artifacts. This difference serves as our DI signal for inferring whether the suspect dataset was used in the target model’s training. This calibration strategy enhances DI’s robustness, reducing false positives while maintaining high detection accuracy.

We demonstrate the effectiveness of our approach on diverse textual datasets, ranging from single-author datasets to large-scale, multi-author collections such as Wikipedia.

Our results show that using *synthetic* held-out data, combined with calibration, enables DI to detect unauthorized training data use with high confidence while keeping false positives low. This expands the practical applicability of DI and provides a pathway for data owners to safeguard their intellectual property in an era of LLMs.

## 2. Background and Related Work

### 2.1. Membership Inference

MIAs focus on deciding if a single data point was included in a given model’s training dataset and often serve as features extractors for DI. In the LLM domain, MIAs exploit different signals to distinguish between members (training data points) and non-members (data points not used during training). For instance, LOSS exploits the perplexity or loss function of the target model (Yeom et al., 2018). Shi et al. (2024) find that the rare words in a sequence can leak more privacy information, and select K% tokens with the smallest probabilities for evaluation. Min-K%++ further improves upon the Min-K% approach by introducing two calibration factors (Zhang et al., 2024b). Zlib ratio (Carlini et al., 2021) uses the compression rate of z-library to normalize the perplexity of the target model. Neighborhood-based methods compare a suspect sequence with its neighboring texts, which can be produced by synonym substitution (Mattern et al., 2023) or paraphrasing (Duarte et al., 2024). Moreover, reference-based methods compare the output signals on a suspect sample between the target model and a reference model (Fu et al., 2024). Yet, many recent works have shown that the evaluation of MIAs suffers from a falsified experimental setup, where a distributional shift exists between the member and non-member sets (Zhang et al., 2024a; Maini et al., 2024; Das et al., 2024). Duan et al. (2024) show that most MIAs only perform slightly better than random guessing if evaluated correctly on non-biased benchmarks. Recently, Kazmi et al. (2024) proposed how to de-bias MIAs from this distribution shift—which we

use as a foundation for our DI calibration.

## 2.2. Dataset Inference

To strengthen the signal from training data further beyond MIAs, Maini et al. (2021) introduced DI. DI aggregates the membership signal over multiple data points, often referred to as *suspect set*, to decide whether a given model was trained on this data. More formally, given a target model  $f$ , DI aims to detect whether  $f$  was trained on the suspect dataset  $\mathcal{D}_{\text{sus}}$ . Therefore, it needs an additional held-out dataset  $\mathcal{D}_{\text{val}}$  from the same distribution as  $\mathcal{D}_{\text{sus}}$ . Given both sets, DI extracts membership features from the data points in  $\mathcal{D}_{\text{sus}}$  and  $\mathcal{D}_{\text{val}}$ , aggregates all features per given sample, and then scores these aggregate features through a scoring model. The scores should be lower for members than for non-members. Then, DI performs statistical hypothesis testing on the scores of  $\mathcal{D}_{\text{sus}}$  and  $\mathcal{D}_{\text{val}}$ . The null hypothesis is that the average scores for  $\mathcal{D}_{\text{sus}}$  are higher than for  $\mathcal{D}_{\text{val}}$ . If the statistical test manages to reject this null hypothesis, this is a confident indicator that the data points from  $\mathcal{D}_{\text{sus}}$  are indeed members of model  $f$ ’s training data.

How to extract the best membership features from the data points varies based on the learning paradigm. For example, the original DI for supervised models (Maini et al., 2021) designs a random walk strategy to estimate the distance between data points and the decision boundary of a supervised model. This is based on the intuition that member data points are further to the decision boundaries than non-member data points. For self-supervised models, Dziedzic et al. (2022) use Gaussian Mixture Model to estimate the representational differences between the training dataset (members) and the test data. Recent work for DI on LLMs (Maini et al., 2024) relies on existing LLM MIAs to extract membership features and uses a linear model to weight the respective features. We follow this approach in our evaluations.

## 3. Failure Cases of DI

In this section, we dive deeper into the difficulties that arise from DI’s assumption on the availability of an additional in-distribution held-out dataset. More precisely, we show that this assumption is extremely hard to meet in practice, even in the simplest setups—only the articles for a single author are used for DI.

### 3.1. DI on a Single Author’s Data

We consider a practical application of DI in copyright protection as detailed in Figure 1. In this scenario, an author has some published texts on the internet of which they believe that they were illegitimately used by an LLM provider to train their model. The author provides this published works to an arbiter, as a suspect set and some non-published

Table 1. The distributional shift (GPT2 AUC) and DI p-value between a suspect set that consists of *non-members* and held-out blog posts. Here, p-value  $< 0.05$  indicates DI incorrectly suggests that the suspect set is a member set.

Sequences per Blog	5	10	15	20	25
GPT2 AUC (%)	52.0	55.2	53.2	58.2	58.6
DI p-value	0.002	$<0.001$	$<0.001$	$<0.001$	$<0.001$
True Membership	×	×	×	×	×
Inferred Membership	✓	✓	✓	✓	✓

blog-posts as held-out set from the same distribution, *i.e.*, with the same style, topics, etc. Then, the arbiter performs DI to resolve the copyright claims.

To evaluate this setup in practice, we collect blog posts of a public blogger. The blogs are split into member, non-member, and held-out sets. To avoid any potential temporal or topic distributional shifts, we randomly shuffle all the collected blogs before splitting. In lack of the computational capacities to train an LLM from scratch, we finetune a Pythia model (Biderman et al., 2023) on the member set. The Pythia model is trained on the Pile dataset (Gao et al., 2020), so we only used blogs after the release date of the Pile to ensure that none of the blogs is part of the pre-training data. Also, we only finetune the target model on the member set for one epoch. Finally, we run DI. More detailed experiment configurations can be found in Section 5.1.

### 3.2. Metrics of Distributional Gap

Before analyzing the results, we introduce the metrics we use to quantify the distributional shift between the suspect and held-out sets. Following the approach of Blind Baselines (Das et al., 2024), we formulate the measurement of the distribution gap between two text datasets as a classification problem. In particular, the suspect set  $\mathcal{D}_{\text{sus}}$  is randomly split into a classifier training split  $\mathcal{D}_{\text{sus}}^{\text{train}}$  and a test split  $\mathcal{D}_{\text{sus}}^{\text{test}}$ . The held-out set  $\mathcal{D}_{\text{val}}$  is also split into  $\mathcal{D}_{\text{val}}^{\text{train}}$  and  $\mathcal{D}_{\text{val}}^{\text{test}}$  in the same vein. Then, a classifier  $g$  is optimized to distinguish the training splits  $\mathcal{D}_{\text{sus}}^{\text{train}}$  and  $\mathcal{D}_{\text{val}}^{\text{train}}$ . Finally, we calculate the area under the curve (AUC) score of the classifier on the test splits  $\mathcal{D}_{\text{sus}}^{\text{test}}$  and  $\mathcal{D}_{\text{val}}^{\text{test}}$ , which is used to measure the distributional gap between  $\mathcal{D}_{\text{sus}}$  and  $\mathcal{D}_{\text{val}}$ .

The design of the classifier decides how the texts are vectorized and if the discrepancies between texts can be sufficiently captured. Das et al. (2024) apply a bag-of-words (BoW) classifier, which can only detect the differences in terms of word frequency. Instead, we build a GPT2-based classifier with two transformer blocks to also find the differences in grammar, content, styles, etc. between two text distributions. We train the classifier from scratch to avoid the impact of any pre-training data. Using only two transformer blocks of the GPT2 architecture avoids overfitting.

Table 2. Distributional shifts between the suspect set and generated held-out set measured by Bag of Word (BoW) classifier vs GPT2.

Generation Method	BoW AUC (%)	GPT2 AUC (%)
ICL Paraphrasing	76.2	99.0
Preference Optimization	50.2	58.9
Suffix Completion	<b>50.0</b>	<b>52.2</b>

### 3.3. False Positive of DI

The AUC scores of the GPT2-based classifier in Table 1 show that there is a non-negligible distributional shift between the non-member and the held-out sets. The intuition behind this observation is that each blog has different content and topics, which brings different words across the non-member and held-out documents. The gap is enlarged when we sample more sequences from each blog post. This small distributional shift between texts leads to very low p-value in the t-test, causing significant false positive rates during DI. This means that the DI falsely accuses the LLM provider of violating the copy right of an author. What is more is that this shortcoming of DI can be maliciously exploited: authors could deliberately provide held-out data from a different distribution than their suspect data to mislead DI and *illegitimately* accuse the LLM provider. Please also refer to Appendix M for a visualized demonstration of such shifts during DI. As a solution to this problem, in the next section, we propose our approach on generating an adequate in-distribution held-out dataset synthetically.

## 4. Synthesizing Held-out Data

Our approach consists of two subsequent steps. First, we generate high-quality held-out data, then, we perform a calibration to account for the distribution shift that such generation can introduce.

### 4.1. Held-out Data Generation

We explore three approaches that leverage LLMs for generating held-out data based on provided suspect data with minimal distribution shift.

**Baselines.** We adapt two existing generative methods as the baseline for the held-out data generation: 1) in-context learning (ICL) and 2) preference optimization. Please refer to Appendix A and Appendix B for more detailed explanations of the two approaches.

**Suffix Completion.** The failure of the above methods demonstrates the difficulty of producing high-quality held-out data with a small enough distributional gap to the suspect data. To solve this problem, we design a generator training scheme that enables the generator to derive a suspect set from the author’s provided documents, together with a held-out set from the same distribution as this suspect set. As

shown in Figure 2, we ① first segment the provided documents into multiple short sequences. ② All the sequences are shuffled and randomly split into a generator training split and a generator inference split. Then, ③ a low-rank adaptation (LoRA) generator is finetuned on the training split with the cross-entropy loss for next-token prediction. Finally, ④ we segment each sequence in the generator inference split into two parts, and the generator predicts a synthetic suffix based on the prefix. Here, the original suffixes are used as the suspect set, and the synthetic suffixes as the held-out set. Note that, the training and inference sets are split on the shuffled text sequences rather than on the documents. This is to ensure that the text snippets from the generator training and inference splits are from the same distribution, such that the generator can achieve better generalization from the training to the inference set. Furthermore, we design a suffix completion task for generator inference. In this task, both the original suffix and the synthetic suffix share a common prefix. This approach ensures that the synthetic text maintains the same position within a sentence as its original counterpart, making the two suffixes directly comparable. Another important insight is that the generator can produce suffixes of higher quality when the sequence length is relatively short. Therefore, we limit the length of the sequences to no longer than 64 tokens for a smaller distributional gap. The results in Table 2 show that our method achieves a significantly small distributional shift, and even GPT2-based classifier can only achieve an AUC as low as 52.2%. For examples of our generative approach, please refer to Appendix E.

### 4.2. Post-hoc Calibration

Since the generation itself can introduce a distributional shift (natural vs generated) data, DI might yield false positives. This is because it would detect differences between suspect and held-out data also when they only differ in terms of distribution but not necessarily in membership. Therefore, we need to identify and mitigate this distribution shift.

To do so, we rely on an important observation: the generation shift between natural and synthesized data occurs in the textual space, while the shift caused by the potential membership of the suspect set exists in the target LLM’s output space. This allows us to disentangle the two signals. By relying on our GPT-based **text-classifier** from Section 3.2, we can quantify the textual distribution shift caused by the generation. We denote this classifier by  $c_{\text{text}}(x)$ , where  $x$  is the text input for which the classifier should decide if it is original or generated data. Inspired by Kazmi et al. (2024), we also define a second **MIA-classifier** with input signals from both the texts and the outputs of the target model, such that we can quantify the combined effects of generation and the membership signal. Concretely, we train a combined classifier  $c_{\text{comb}}(x, \text{MIA}(f(x)))$  with inputs from both text  $x$  and the MIA signal  $\text{MIA}(x)$  based on the outputs of  $f$ .



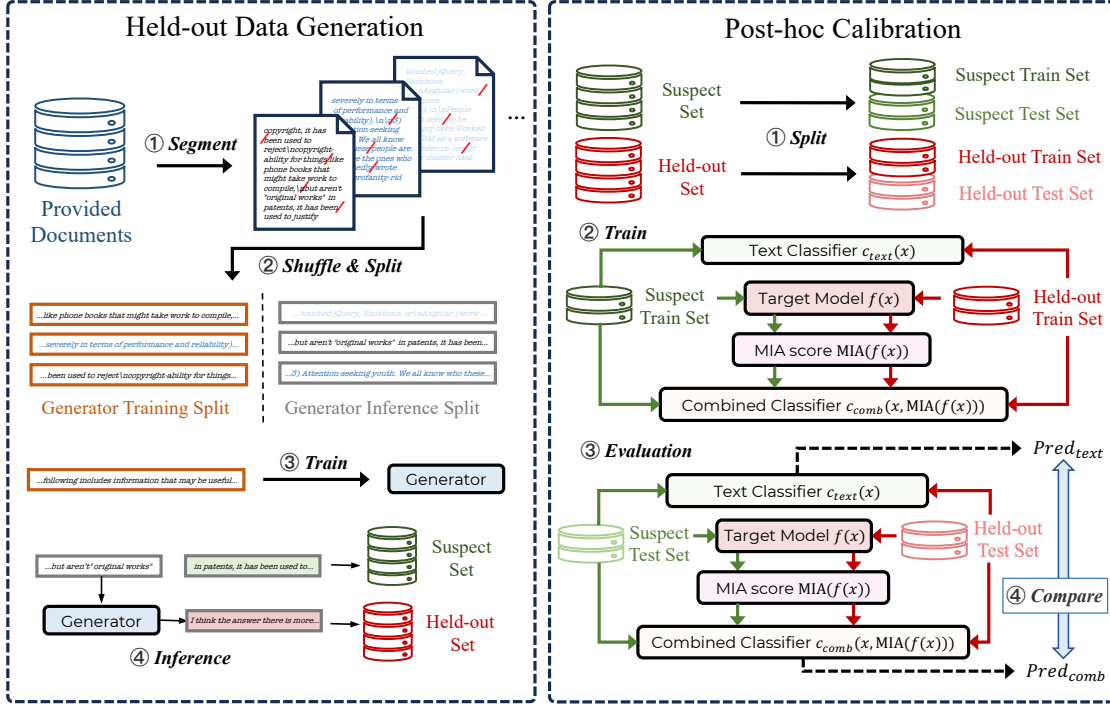


Figure 2. An illustration of our proposed **Held-out Data Generation** (Left Panel) and **Post-hoc Calibration** (Right Panel).

Here,  $MIA(x)$  can also be a vector by concatenating multiple MIA scores. We split both the suspect set and held-out sets into training and test splits. The two classifiers are optimized on the suspect train split  $\mathcal{D}_{sus}^{train}$  and the held-out train split  $\mathcal{D}_{val}^{train}$ , and evaluated on the suspect test split  $\mathcal{D}_{sus}^{test}$  and the held-out test split  $\mathcal{D}_{val}^{test}$ . By comparing the distributional shifts quantified by the MIA classifier and the shifts identified by the text classifier, we can separate the membership signals from the distribution gap caused by generation.

We design a hypothesis test to statistically verify if the combined classifier quantifies a larger distributional shift between the suspect and held-out data than the text classifier, namely the *difference comparison t-test*. First, we sample a suspect data point  $x_{sus} \in \mathcal{D}_{sus}$  and its generated counterpart  $x_{val} \in \mathcal{D}_{val}$ . In every such original/held-out pair, we quantify the shift caused by *generation* with the text classifier as  $c_{text}(x_{val}) - c_{text}(x_{sus})$ . We also quantify the combined effects caused by *generation and membership* with the combined classifier as  $c_{comb}(x_{val}) - c_{comb}(x_{sus})$ . If the membership signal is present, the combined effects will be stronger than the generation effect, and the predicted probability will be slightly more accurate for the combined classifier. To this end, we design a t-test to statistically verify if this prediction difference is significantly larger for  $c_{comb}$  than for  $c_{text}$ . The null hypothesis of the t-test is formalized as follows:

$$\mathcal{H}_0 : \mathbb{E}_{x_{val} \in \mathcal{D}_{val}, x_{sus} \in \mathcal{D}_{sus}} [c_{comb}(x_{val}) - c_{comb}(x_{sus})] \leq \mathbb{E}_{x_{val} \in \mathcal{D}_{val}, x_{sus} \in \mathcal{D}_{sus}} [c_{text}(x_{val}) - c_{text}(x_{sus})]. \quad (1)$$

Table 3. Results for single author blog posts. Here, p-value < 0.05 indicates the suspect set is member set.

True Membership	AUC <sub>Text</sub> (%)	AUC <sub>Comb</sub> (%)	P-value	Inferred Membership
✓	53.8	55.6	0.01	✓
✗	53.8	53.9	0.13	✗

Here, the groundtruth label  $x_{val}$  is defined as 1 and  $x_{sus}$  as 0. The difference comparison t-test is performed multiple times with different random seeds, and the p-values are aggregated with Sidac correction (Šidák, 1967).

## 5. Experimental Evaluation

We start by introducing our experimental setup, further detailed in Appendix D. Then, we present the results of DI executed based on our generated held-out data. We also perform ablation studies to investigate the contribution of each component in our proposed method. Finally, we analyze the impact of t-test sample size and the classifier architecture.

### 5.1. Experimental Setup

**Single author data.** We collect 1400 blog posts from a single author. All figures, tables, videos, and hyperlinks are removed during pre-processing and only plain text is used for evaluation. We sample 450 posts as member data and finetune a Pythia 410M deduplicated model as target model.

The other posts are held out as non-member and held-out sets for the evaluation.

**More Complicated Dataset and Model.** We also evaluate our method on the Pile dataset (Gao et al., 2020), which is much more complicated and has subsets of diverse types of texts. We use the de-duplicated version of Pythia 410M model as the target model. The training split of the Pile dataset is used as member data, and the held-out and test split is used as non-member data. Here, we only evaluate Pile subsets that are free from copyright issues. Please also refer to Appendix C for detailed configuration on the Pile.

**Implementation Details** We finetune a Llama 3 8B model (Dubey et al., 2024) with LoRA as the generator. For both types of datasets, we split 2,000 sequences as the generator inference set, and the others as the generator training split. Both text classifier and combined classifier are trained on 1,000 synthetic held-out data and 1,000 suspect data for each dataset. Our proposed t-test is also conducted on 1,000 synthetic held-out data and 1,000 suspect data. More implementation details can be found in Appendix D. We also provide an analysis of hyperparameter sensitivity in Appendix I.

## 5.2. Results for Single Author Dataset

The experimental results on the single author dataset are presented in Table 3. On the member set, the combined classifier  $c_{\text{comb}}$  outperforms the text classifier  $c_{\text{text}}$ , by a large margin of 1.8% AUC score. Moreover, the observed p-value of 0.01 strongly supports the alternative hypothesis, indicating that the superior performance of  $c_{\text{comb}}$  over  $c_{\text{text}}$  is statistically significant. This enables our method to correctly identify that the target set is part of the training set. For the non-member set,  $c_{\text{comb}}$  and  $c_{\text{text}}$  achieve comparable AUC scores, with a p-value of 0.13 that significantly exceeds the threshold of 0.05. This result confirms the ability of our approach to correctly identify non-member texts as such, thus avoiding the false positives that occur with the original LLM DI approach. Here, we finetune the target model on the single author dataset with LoRA for one epoch. We also present the results with other fine-tuning setups in Appendix F.

## 5.3. Results for Pile Datasets

The results of different Pile subsets are shown in Table 4. We observe that DI correctly predicts the membership of datasets from diverse domains and styles, including plain text, academic writing, and code using our method for generating the held-out data. The results also show that our generation method generalizes well to documents with different lengths, ranging from 1 KB (Wikipedia) to 70 KB (PhilPapers). Moreover, our proposed method generalizes well to texts from different domains and languages, e.g., medical (PubMed Central), legal (FreeLaw), technical (ArXiv), and multilingual (EuroParl) domains. Notably, the p-values for

Table 4. Results for different Pile subsets. *True* represents the true membership while *Inferred* denotes the inferred membership. Our generation is successful if these two align.

Subset	True	AUC <sub>Text</sub> (%)	AUC <sub>Comb</sub> (%)	P-value	Inferred
Pile-CC	✓	53.1	60.3	<0.001	✓
	✗	52.5	48.3	1.0	✗
Wikipedia	✓	51.7	58.6	<0.001	✓
	✗	52.2	52.0	0.43	✗
ArXiv	✓	53.9	57.3	<0.001	✓
	✗	53.1	44.7	1.0	✗
NIH Exporter	✓	51.4	54.1	0.005	✓
	✗	53.3	51.6	1.0	✗
FreeLaw	✓	55.6	56.7	0.003	✓
	✗	51.6	51.6	0.84	✗
Ubuntu IRC	✓	52.7	54.5	0.002	✓
	✗	52.5	54.2	0.12	✗
PubMed Central	✓	54.1	54.4	0.004	✓
	✗	52.4	49.5	0.66	✗
Github	✓	54.0	59.3	<0.001	✓
	✗	53.3	51.9	0.97	✗
EuroParl	✓	51.3	65.0	<0.001	✓
	✗	51.9	47.7	1.0	✗
PhilPapers	✓	58.5	57.0	<0.001	✓
	✗	58.1	55.3	0.13	✗
HackerNews	✓	56.4	57.5	<0.001	✓
	✗	57.1	56.3	0.14	✗
Enron Emails	✓	56.9	58.2	0.001	✓
	✗	58.4	53.8	0.99	✗
StackExchange	✓	54.0	60.0	<0.001	✓
	✗	52.7	50.8	1.0	✗
PubMed Abstract	✓	54.9	59.9	<0.001	✓
	✗	54.7	53.0	0.66	✗
USPTO Backgrounds	✓	56.7	58.1	<0.001	✓
	✗	55.8	55.7	0.13	✗

our difference comparison t-test are significantly lower than 0.05 on all the evaluated member sets, and higher than 0.1 on all the non-member sets.

## 6. Conclusions

We propose how to *synthetically generate* an in-distribution held-out dataset to enable the real-world application of DI. Therefore, we solve two critical challenges, namely (1) creating high-quality, diverse synthetic data that accurately reflects the original distribution and (2) bridging likelihood gaps between real and synthetic data. Our solution relies on designing a data generator training scheme based on a suffix-based completion task and post-hoc calibration to align the likelihood gaps between real and synthetic data. Through extensive experimental evaluation, we highlight that our method enables a robust DI and correctly identifies training data while achieving a low false positive rate. This shows our method’s reliability to support copyright owners to make legitimate claims on data usage for real-world litigations.

## References

- The times sues openai and microsoft over a.i. use of copyrighted work <https://www.nytimes.com/2023/12/27/business/media/new-york-times-open-ai-microsoft-lawsuit.html>. 2023. URL <https://www.nytimes.com/2023/12/27/business/media/new-york-times-open-ai-microsoft-lawsuit.html>.
- Sarah silverman and authors sue openai and meta over copyright infringement. 2023. URL <https://www.nytimes.com/2023/07/10/arts/sarah-silverman-lawsuit-openai-meta.html>.
- Balloccu, S., Schmidová, P., Lango, M., and Dušek, O. 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)*, pp. 67–93, 2024.
- Biderman, S., Schoelkopf, H., Anthony, Q. G., Bradley, H., O’Brien, K., Hallahan, E., Khan, M. A., Purohit, S., Prashanth, U. S., Raff, E., et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T., Song, D., Erlingsson, U., et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pp. 2633–2650, 2021.
- Das, D., Zhang, J., and Tramèr, F. Blind baselines beat membership inference attacks for foundation models. *arXiv preprint arXiv:2406.16201*, 2024.
- Duan, M., Suri, A., Mireshghallah, N., Min, S., Shi, W., Zettlemoyer, L., Tsvetkov, Y., Choi, Y., Evans, D., and Hajishirzi, H. Do membership inference attacks work on large language models? *arXiv preprint arXiv:2402.07841*, 2024.
- Duarte, A. V., Zhao, X., Oliveira, A. L., and Li, L. Decop: Detecting copyrighted content in language models training data. *arXiv preprint arXiv:2402.09910*, 2024.
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Dubiński, J., Kowalczyk, A., Boenisch, F., and Dziedzic, A. Cdi: Copyrighted data identification in diffusion models. *arXiv preprint arXiv:2411.12858*, 2024.
- Dziedzic, A., Duan, H., Kaleem, M. A., Dhawan, N., Guan, J., Cattani, Y., Boenisch, F., and Papernot, N. Dataset inference for self-supervised models. *Advances in Neural Information Processing Systems*, 35:12058–12070, 2022.
- Fu, W., Wang, H., Gao, C., Liu, G., Li, Y., and Jiang, T. Membership inference attacks against fine-tuned large language models via self-prompt calibration. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Gao, L., Biderman, S., Black, S., Golding, L., Hoppe, T., Foster, C., Phang, J., He, H., Thite, A., Nabeshima, N., et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- Golchin, S. and Surdeanu, M. Time travel in llms: Tracing data contamination in large language models. *arXiv preprint arXiv:2308.08493*, 2023.
- Groeneveld, D., Beltagy, I., Walsh, E., Bhagia, A., Kinney, R., Tafjord, O., Jha, A., Ivison, H., Magnusson, I., Wang, Y., et al. Olmo: Accelerating the science of language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15789–15809, 2024.
- Kazmi, M., Lautraite, H., Akbari, A., Tang, Q., Soroco, M., Wang, T., Gambs, S., and Lécuyer, M. PANORAMIA: Privacy auditing of machine learning models without retraining. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Li, C. and Flanigan, J. Task contamination: Language models may not be few-shot anymore. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 18471–18480, 2024.
- Magnusson, I., Bhagia, A., Hofmann, V., Soldaini, L., Jha, A. H., Tafjord, O., Schwenk, D., Walsh, E., Elazar, Y., Lo, K., et al. Paloma: A benchmark for evaluating language model fit. *Advances in Neural Information Processing Systems*, 37:64338–64376, 2024.
- Maini, P. and Suri, A. Reassessing emnlp 2024’s best paper: Does divergence-based calibration for membership inference attacks hold up? 2024. URL <https://www.anishumansuri.com/blog/2024/calibrated-mia/>. Accessed January 29, 2025.
- Maini, P., Yaghini, M., and Papernot, N. Dataset inference: Ownership resolution in machine learning. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.

- Maini, P., Jia, H., Papernot, N., and Dziedzic, A. LLM dataset inference: Did you train on my dataset? In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Mattern, J., Mireshghallah, F., Jin, Z., Schoelkopf, B., Sachan, M., and Berg-Kirkpatrick, T. Membership inference attacks against language models via neighbourhood comparison. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 11330–11343, 2023.
- Meeus, M., Shilov, I., Jain, S., Faysse, M., Rei, M., and de Montjoye, Y.-A. Sok: Membership inference attacks on llms are rushing nowhere (and how to fix it). *arXiv preprint arXiv:2406.17975*, 2024.
- Meng, Y., Xia, M., and Chen, D. SimPO: Simple preference optimization with a reference-free reward. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Oren, Y., Meister, N., Chatterji, N. S., Ladhak, F., and Hashimoto, T. Proving test set contamination in black-box language models. In *The Twelfth International Conference on Learning Representations*, 2024.
- Penedo, G., Kydlíček, H., Lozhkov, A., Mitchell, M., Raffel, C., Von Werra, L., Wolf, T., et al. The fineweb datasets: Decanting the web for the finest text data at scale. *arXiv preprint arXiv:2406.17557*, 2024.
- Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Ermon, S., and Finn, C. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- Rahman, N. and Santacana, E. Beyond fair use: Legal risk evaluation for training llms on copyrighted text. 2023. URL <https://genlaw.org/CameraReady/57.pdf>.
- Reuters. Getty images lawsuit says stability ai misused photos to train AI, 2023. URL <https://www.reuters.com/legal/getty-images-lawsuit-says-stability-ai-misused-photos-train-ai-2023-02-06/>.
- Roberts, M., Thakur, H., Herlihy, C., White, C., and Dooley, S. To the cutoff... and beyond? a longitudinal perspective on llm data contamination. In *The Twelfth International Conference on Learning Representations*, 2024.
- Shi, W., Ajith, A., Xia, M., Huang, Y., Liu, D., Blevins, T., Chen, D., and Zettlemoyer, L. Detecting pretraining data from large language models. In *The Twelfth International Conference on Learning Representations*, 2024.
- Shokri, R., Stronati, M., Song, C., and Shmatikov, V. Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pp. 3–18. IEEE, 2017.
- Šidák, Z. Rectangular confidence regions for the means of multivariate normal distributions. *Journal of the American statistical association*, 62(318):626–633, 1967.
- Soldaini, L., Kinney, R., Bhagia, A., Schwenk, D., Atkinson, D., Authur, R., Bogin, B., Chandu, K., Dumas, J., Elazar, Y., et al. Dolma: an open corpus of three trillion tokens for language model pretraining research. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15725–15788, 2024.
- Weber, M., Fu, D. Y., Anthony, Q., Oren, Y., Adams, S., Alexandrov, A., Lyu, X., Nguyen, H., Yao, X., Adams, V., Athiwaratkun, B., Chalamala, R., Chen, K., Ryabinin, M., Dao, T., Liang, P., Ré, C., Rish, I., and Zhang, C. Redpajama: an open dataset for training large language models. *NeurIPS Datasets and Benchmarks Track*, 2024.
- Wu, X., Duan, R., and Ni, J. Unveiling security, privacy, and ethical concerns of chatgpt. *Journal of Information and Intelligence*, 2023.
- Xu, H., Sharaf, A., Chen, Y., Tan, W., Shen, L., Durme, B. V., Murray, K., and Kim, Y. J. Contrastive preference optimization: Pushing the boundaries of LLM performance in machine translation. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=5liwkioZpn>.
- Yeom, S., Giacomelli, I., Fredrikson, M., and Jha, S. Privacy risk in machine learning: Analyzing the connection to overfitting. In *2018 IEEE 31st computer security foundations symposium (CSF)*, pp. 268–282. IEEE, 2018.
- Zhang, J., Das, D., Kamath, G., and Tramèr, F. Membership inference attacks cannot prove that a model was trained on your data. *arXiv preprint arXiv:2409.19798*, 2024a.
- Zhang, J., Sun, J., Yeats, E., Ouyang, Y., Kuo, M., Zhang, J., Yang, H. F., and Li, H. Min-k%++: Improved baseline for detecting pre-training data from large language models. *arXiv preprint arXiv:2404.02936*, 2024b.



## A. Details of ICL Generation

We experiment with GPT-4-Turbo and prompt it using in-context-learning learning to paraphrase the suspect data. Each prompt includes a few data points as demonstrations (shots) and requests the model to produce paraphrases for the suspect set. We test two types of templates to prompt GPT-4-Turbo model with in-context-learning (ICL). One is to paraphrase the given sample based on the examples, and the other is to complete the text with part of the sentence.

### A.1. Text Completion Prompt

In the text completion prompt, we give ten examples from the author provided documents and truncate a target sample. Then we prompt the GPT-4-Turbo model to complete the target sample, which is shown at the end of the prompt. Here is an example of this kind of queries:

*Input:*

I have some text samples. Please help me complete the last sample based on my example samples.  
 A device may be provided to sense an environment's present temperature. For example, [...] a thermistor that is robust and reliable and that can be fabricated at a reasonable cost.  
 [...]  
 This invention relates to a method of synthesizing musical tones. A musical tone or natural tone produced by an acoustic or natural musical instrument can [...] With the music synthesizer it is considerably difficult to obtain dynamic audio spectra.pneumoniae.  
 Radio communications receivers commonly receive and decode messages which contain information. A part of this information may be [...] a single radio will transmit messages to the remainder of the group. These radio

*Output:*

messages are selectively received by the grouped radios which decode the identity information within the [...] and secure group communications in diverse and often challenging operational environments.

### A.2. Paraphrasing Prompt

In the paraphrasing prompt, we give ten examples from the author provided documents, and prompt the GPT-4-Turbo model to paraphrase the last. Here is an example of this kind of queries:

*Input:*

Paraphrase the last sample according to the other samples.  
 1. Field of the Invention This invention relates generally to fluid injectors. More particularly, this invention relates to a proportional needle control fuel injector having a [...] The needle valve is moveable between a first position closing the injection fuel outlet port and a second position to open to injection fuel outlet.  
 [...]  
 The present invention relates to the structure of a casing of a miniature portable [...] is forcibly urged away from the casing by a tool, it is likely that the casing is scratched or the base plate is deformed.  
 This invention relates to fuel cells and, in particular, to improvements in the catalyst used in internal [...] fuel cell in which the catalyst assembly has minimum displacement and permits vertical and horizontal orientations of the fuel cell.

*Output:*

This invention pertains to fuel cells, specifically addressing enhancements in catalysts utilized [...] and efficiency of internal reforming fuel cells by mitigating previous limitations and improving overall system functionality.

### A.3. Distribution Shift Analysis

The distributional shifts are quantified in BoW AUC and GPT2 AUC as shown in Table A1. As observed, both types of prompt templates produces synthetic texts with large distributional shifts to the suspect sets. Notably, the GPT2-based classifier can achieve as much as an AUC of 99.2%. The reason is that there are many words (such as "remarkable" and "moreover") that appear much more frequently in the synthetic text than in the human-written text.

Table A1. Distributional shifts between the suspect set and GPT-4-Turbo generated validation set.

Template Type	BoW AUC(%)	GPT2 AUC(%)
ICL Text Completion	79.2	99.2
ICL Paraphrasing	76.2	99.0

## B. Details of Preference Optimization Generation

Preference optimization methods focus on optimizing a pre-trained LLM based on human preference (Rafailov et al., 2024; Xu et al., 2024). Particularly, LLMs iteratively produce random generations, then human annotators are requested to label the generations as chosen or rejected, and the LLMs are further optimized according to this human feedback. We note that, we can leverage preference optimization approaches to make our generator model prefer the human-written texts over synthetic data, thus producing texts with a more similar distribution to natural texts. Here, we instantiate the preference optimization scheme with a state-of-the-art method, the simple preference optimization (SimPO) (Meng et al., 2024). During each training iteration, the human-written suspect data are always labeled as chosen and the generations from the last iteration are marked as rejected. As noted in Section 4.1, this approach improves significantly upon prompted paraphrasing, but still causes a large distributional shift between the suspect set and the generated held-out set.

Table A2. Segmentation configurations for different Pile subsets.

Subset	Number of Test Set in Pile	Chosen Split Size	Max. Snippets per Document	Number of Tokens per Snippet
Pile-CC	>4000	4000	20	32
StackExchange	>4000	4000	5	64
PubMed Abstracts	>4000	4000	20	64
Wikipedia (en)	>4000	4000	5	32
USPTO Backgrounds	>4000	4000	20	64
PubMed Central	>4000	4000	10	32
FreeLaw	>4000	4000	5	32
ArXiv	>4000	200	100	32
NIH ExPorter	>3000	3000	10	32
HackerNews	>3000	3000	10	64
Github	>1000	1000	30	32
Enron Emails	1957	1957	30	64
EuroParl	290	290	200	32
PhilPapers	132	132	500	64
Ubuntu IRC	43	43	500	32

## C. Pile Dataset Segmentation

We present the details for the configurations of Pile subset in Table A2. We note that, it is claimed that the following Pile subsets may have copyright issues and cannot be included for evaluation: Books3, OpenWebText2, Gutenberg (PG-19), OpenSubtitles, BookCorpus2, and YoutubeSubtitles. For most subset there are documents that are much longer than the other documents, which causes that too many snippets are sample from these documents if all snippets are used. Therefore, we set a maximum snippet for each document on each subset according to the median lengths of the documents. Also, we note that our approach can achieve good performance on most subsets with only 32 tokens. For certain subsets, we use

a token length of 64 for a stronger membership signal. The average token number for the Pile subset is 45. In practical applications, we suggest choosing the token numbers when the text-classifier has the minimal AUC to avoid distributional shift as much as possible. As a more conservative approach, this avoids making false accusations of copyright violation.

## D. Implementation Details

### D.1. Generator

The LoRA rank for the generator is 32. The generator is trained for 100 epochs, and the learning rate is set to  $2 \times 10^{-4}$ . We set a warm-up ratio of 0.03, and a linear scheduler is used to dynamically adjust the learning rate.

## D.2. Text and Combined Classifiers

For both the text and the combined classifier, we leverage the basic architecture of the GPT2 classifier with an extra linear layer. Specifically, the classifier has only two layers, with an embedding dimension of 1600 and an attention head number of 25. The GPT2-based classifier is optimized for 20 epochs, and the linear layer is further optimized for 200 epochs.

### E. Examples of Synthetic Texts

In this section, we provide some examples of the synthetic texts on the Pile dataset. Here, prefix denotes the first half of the generated text, real suffix refers to the original suffix of the natural text, and generated suffix refers to the synthetic completion based on the prefix. We observe that, the generated suffixes are reasonable continuation of the prefixes. The generated suffixes also align with the style of each dataset and do not overfit to the content of the real suffixes.

### E.1. Pile-CC

*Prefix:*

are excited about and also what we hoped to see from this years E3!

*Real suffix:*

From the surprising new Spider-Man PS4 game to the bizarre We Happy Few and

*Generated suffix:*

Let us know your thoughts on this monologue as we are preparing for our next

### E.2. StackExchange

*Prefix:*

```
var FKEntityListWithCastCopy = new debiteur().GetType().GetProperty("\"
```

*Real suffix:*

```
schakeling\").GetValue(dbEntry) as List<FKEntity>;//Just
```

*Generated suffix:*

```
FKEntityList\").GetValue(instance, null);\n\nforeach (var t in FKEntity
```

### E.3. PubMed Abstracts

*Prefix:*

were calculated using the Kaplan-Meier method. Of the 117 patients in

*Real suffix:*

whom data were analyzed, 103 had follow-up MR or CT images and 14 patients were

*Generated suffix:*

the study (76 with UC and 41 with DC), 45 patients required proctocolic resection

#### E.4. Wikipedia (en)

*Prefix:*

Em is going away for a while. While it's not up to the standard

*Real suffix:*

of "Mockingbird," it is more fully realized than the two other new

*Generated suffix:*

of their three previous albums, cattle call is still an enjoyable romp,

#### E.5. USPTO Backgrounds

*Prefix:*

1. Field of the Invention\nThis invention relates to a storage device for athletic equipment and, in particular, to a portable storage device for transporting and retaining

*Real suffix:*

elongate items of athletic equipment such as hockey sticks and related athletic equipment.\n2. Discussion of Related Art\nNumerous team athletic activities require individual players on the

*Generated suffix:*

multiple pairs of basketballs.\n2. Description of the Related Art\nDuring the summer and other periods when there is an extended break from an athletic school or program

#### E.6. PubMed Central

*Prefix:*

example, both cycles apply Lewis acidic metal centers to bind the monomers (ep

*Real suffix:*

oxide or lactone), and both invoke labile metal alkoxide intermediates as

*Generated suffix:*

oxides or cyclic carbonates), but the axes of the metallacycle in

#### E.7. FreeLaw

*Prefix:*

Court, 638 P.2d 65 (Colo.1981



*Real suffix:*

Here, the juvenile court denied the GAL's motions because it did not want

*Generated suffix:*

), cert. denied, 454 U.S. 1146, 102

### **E.8. Arxiv**

*Prefix:*

up and vice versa. In contrast, fundamentalists expect the price to track its

*Real suffix:*

fundamental value. Orders from this type of agent may be written as\n\n\$\$D

*Generated suffix:*

underlying fundamentals up and down, but given sufficient acceleration the price  
might \u201crun away

### **E.9. NIH ExPorter**

*Prefix:*

attachment and growth, respectively. Together with an industrial sponsor, Vaxiron  
,

*Real suffix:*

Inc., we will develop quality control tools and metrics for assessing vaccine  
antigen formulations,

*Generated suffix:*

the applicant has carried out clinical trials of different vaccine candidates  
based on different viruses for

### **E.10. Github**

*Prefix:*

.string \"reach only by using a BIKE technique.\$\" \n\nRoute110\_Text\_

*Real suffix:*

16EEF6:: @ 816EEF6\n\t.string \"Which

*Generated suffix:*

16F381:: @ 816F381\n\t.string \"ROUTE {ROAD

### **E.11. Enron Emails**

*Prefix:*

Lay. He went on to say that Kenneth was Dewayne Re

*Real suffix:*

es' cousin and started telling about all of your fine attributes and what a

*Generated suffix:*

ams' direct \nreport and that it would be extremely difficult for Kenneth to get

### **E.12. EuroParl**

*Prefix:*

het mondeling amendement op schrift heeft gekregen.\nIk st

*Real suffix:*

el voor om niet te spreken over \"de Raad en de lidstat

*Generated suffix:*

akk voor de uitnodiging om tijdens uw volgende bij

### **E.13. PhilPapers**

*Prefix:*

distribute well among [the gods who fought with him] their titles and privileges

*Real suffix:*

" (885, cf. 66\u201367 and 74); to swallow

*Generated suffix:*

(17.1). Orderly distribution of praise for the victory is re

### **E.14. Ubuntu IRC**

*Prefix:*

about setting up reoccurring status meetings?\n<dfarning> should we start

*Real suffix:*

holding those or is it too soon?\n<dfarning> Luke will be joining

*Generated suffix:*

with a status meeting or a design meeting?\n<manusheel> dfarning

### **E.15. HackerNews**

*Prefix:*

Angular (work just uses Dojo).\n\nPeople don't seem to

*Real suffix:*

be hungry here.\n\n-----\nlewispollard\nWorked for IBM as a software engineer on  
one of

*Generated suffix:*

care that it's adding yet another ~20KB per page. We're\nsaying no

## F. Other Finetuning Configurations for Single Author Dataset

We evaluate our proposed approach on the single author dataset under different finetuning settings in Table A3. The results show that, the membership signal is stronger with more iterations or larger parameter size, and therefore easier to detect.

Table A3. Results for different fine-tuning methods. *True* represents the true membership while *Inferred* denotes the inferred membership. Our generation is successful if these two align.

Fine-tuning Method	True	AUC <sub>Text</sub> (%)	AUC <sub>Comb</sub> (%)	P-value	Inferred
LoRA (1 epoch)	✓	53.8	55.6	0.01	✓
	✗	53.8	53.9	0.13	✗
LoRA (10 epochs)	✓	53.7	56.2	0.005	✓
	✗	53.6	53.5	0.14	✗
Full- finetuning	✓	53.7	56.8	0.008	✓
	✗	53.8	53.7	0.21	✗

## G. Results on the OLMo Model

We conduct the experiments to analyze the performance with OLMo-7B (Groeneveld et al., 2024). The OLMo-7B model is trained on the Dolma V.1.7 dataset (Soldaini et al., 2024), which serves as our member set. Following Duan et al. (2024), we employ Paloma (Magnusson et al., 2024) as the non-member set. The results in Table A4 demonstrate that our method successfully detects both member and non-member sets for Wikipedia and Common Crawl subsets when using the OLMo-7B model as the target model.

Table A4. Results for OLMo-7B on different data subsets. *True* represents the true membership while *Inferred* denotes the inferred membership. Our generation is successful if these two align.

Subset	True	AUC <sub>Text</sub> (%)	AUC <sub>Comb</sub> (%)	P-value	Inferred
Wikipedia	✓	52.9	55.4	0.009	✓
	✗	52.1	50.6	1.0	✗
Common Crawl	✓	53.5	55.7	0.01	✓
	✗	54.2	53.8	0.68	✗

## H. Ablation Studies on Single Author Dataset

Besides the ablation studies on the Pile presented in Appendix H.1, we also perform the ablation studies on the single author dataset. The results in Table A5 follow a similar trend to the Pile, showing the importance of each component in our framework.

Table A5. Results for different configurations. *True Membership* represents the true membership while *Inferred Membership* denotes the inferred membership. Our generation is successful if these two align.

Configuration	True Membership	P-value	Inferred Membership
w/o Suffix Completion (ICL Paraphrasing)	✓	1.0	✗
	✗	1.0	✗
w/o Post-hoc Calibration (Original T-test in DI)	✓	<0.001	✓
	✗	<0.001	✓
Ours	✓	0.01	✓
	✗	0.13	✗

### H.1. Ablation on Post-hoc Dataset Inference

We conduct ablation studies to separately analyze the contribution of the three components in our held-out data generation: suffix completion and post-hoc calibration.

**Suffix Completion.** As presented in Table 2, our proposed sequence completion scheme can synthesize held-out texts with a distribution much more similar to the suspect texts when compared with the baseline methods. In addition to the AUC results, we also show that the baseline generation methods cannot produce reliable held-out sets even when combined with our post-hoc calibration and weight constraint in Table A6. In particular, we replace our generation scheme with three baselines, including ICL paraphrasing, ICL text completion, and preference optimization. The p-values are presented as Setting 1-3. We also remove two key designs in our generation method, 1) Segment and Shuffle, and 2) Suffix Comparison, as shown in Setting 4-5. In all the above settings, the p-values for both member and non-member sets are 1.0, which indicates that the  $c_{\text{text}}$  has better or similar performance when compared with  $c_{\text{comb}}$ . The reason behind the observation is that the distributional shift caused by the generation is much larger than the shift induced by the membership signal, such that  $c_{\text{comb}}$  does not outperform  $c_{\text{text}}$  even with extra membership inputs on the member set. Consequently, the DI predicts both sets as non-member and suffers from false negatives.

**Post-hoc Calibration.** We replace our calibration method with the original DI without calibration, as shown in Setting 6 of Table A6. Specifically, only a linear classifier is optimized to aggregate different MIA metrics and output the final prediction score. Furthermore, the t-test is conducted directly between the predictions on the target set and the ones on the held-out set. We observe that the p-values under this condition are extremely low for both member and non-member sets, and DI has false positive in this case. This observation aligns with results in Section 3, where we show that even a small distributional shift causes a significantly small p-value in the original DI. Therefore, our post-hoc calibration approach is crucial to evaluating the distributional shift caused only by membership signals.

Table A6. Ablation studies of our approach. **Setting 1-3:** replacing our generation method with baselines. **Setting 4-5:** removing key designs from our generation method. **Setting 6:** without post-hoc calibration. **Setting 7:** our complete method.

Setting	Configuration	True Membership	P-value	Inferred Membership
1	w/o Suffix Completion (ICL Paraphrasing)	✓	1.0	✗
		✗	1.0	✗
2	w/o Suffix Completion (ICL Text Completion)	✓	1.0	✗
		✗	1.0	✗
3	w/o Suffix Completion (Preference Optimization)	✓	1.0	✗
		✗	1.0	✗
4	w/o Segment and Shuffle	✓	1.0	✗
		✗	1.0	✗
5	w/o Suffix Comparison	✓	1.0	✗
		✗	1.0	✗
6	w/o Post-hoc Calibration (Original T-test in DI)	✓	<0.001	✓
		✗	<0.001	✓
7	Ours	✓	<0.001	✓
		✗	1.0	✗

## I. Analysis of Hyperparameter Sensitivity

We conducted a comprehensive analysis of hyperparameter sensitivity, focusing on two key parameters: the number of epochs and the number of t-test samples. The number of epochs represents the training epochs for our linear model that aggregates MIA scores. The number of t-test samples indicates the total sample size used in our statistical analysis, including both the suspect and synthetic held-out sets. Our experimental results in Table A7 demonstrate that our proposed method exhibits robust performance across a wide range of values for both hyperparameters, indicating low sensitivity to these configuration choices.



Table A7. Performance of our method across different numbers of epochs and T-test samples.

Hyperparameter	Value	True membership	P-value	Inferred membership
Number of Epochs	100	✓	<0.001	✓
		✗	1.0	✗
	200	✓	<0.001	✓
		✗	1.0	✗
Number of T-test Samples	500	✓	<0.001	✓
		✗	1.0	✗
	1000	✓	0.003	✓
		✗	1.0	✗
Number of T-test Samples	1000	✓	<0.001	✓
		✗	1.0	✗
	2000	✓	<0.001	✓
		✗	1.0	✗
Number of T-test Samples	3000	✓	<0.001	✓
		✗	0.41	✗
	4000	✓	<0.001	✓
		✗	0.25	✗

Table A8. The AUC of different classifier architectures.

Architecture	AUC <sub>Text</sub>	Training Time
all-MiniLM-L6-v2	50.8	<b>0.3</b>
BERT	51.2	2.0
Llama3-8B, Pre-trained+LoRA	53.2	65.1
GPT2, Pre-trained+LoRA	53.0	26.2
GPT2, Pre-trained+Full Finetuned	52.3	36.8
GPT2, 2 Layers+Initialized	<b>53.3</b>	0.5

## J. Other Related Works about Test Set Contamination Detection

Test set contamination is a newly identified risk, where the public test benchmarks are involved during LLM training (Balloccu et al., 2024). For example, Roberts et al. (2024) observe that LLMs are better at generating code with more appearances on GitHub, revealing that LLMs can be contaminated with open-source GitHub data and are overestimated on coding tasks. Similarly, Li & Flanigan (2024) demonstrate that some LLMs have a better performance on few-shot benchmarks constructed before the model training, which indicates test set contamination for LLMs. To detect test set contamination, Golchin & Surdeanu (2023) design prompts that guide LLM to reproduce exact or near-exact test set instances, such that the model encloses the contaminated samples memorized during the pre-training phase. Oren et al. (2024) compare the target model predictions between a test set and all of its permutations. However, this method is based on the assumption that the test set is involved in the training set in its exact order, which could be interrupted by a random shuffle before training. Test set contamination can also be a potential application of our method, as the proposed approach can perform training data detection on complex datasets composed by different authors.

## K. Analysis of Sample Size

We also set out to analyze how the sample size in our proposed t-test affects the statistical confidence of DI with our generated held-out data. Here, the sample size is the total number of the suspect and held-out set, which is also the number of queries made to the target model. The two sets are of the same size, as they are produced in a pairwise manner. We observe from Figure A1 that, as the number of samples increases, DI exhibits improved detection capability of training data. Notably, with fewer than 1,000 samples, DI achieves statistical significance ( $p < 0.05$ ) across most of the evaluated datasets. When increasing the sample size to 2k queries, the method demonstrates even stronger statistical significance ( $p < 0.01$ ) consistently across all datasets.

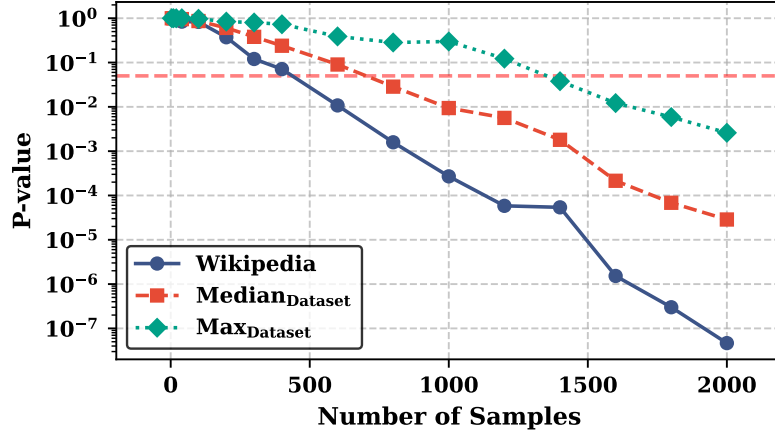


Figure A1. The p-values of member sets with change in sample size. Median<sub>Dataset</sub> denotes the median p-value of different datasets, and Mean<sub>Dataset</sub> is the maximum p-value of all subsets. Number of samples refers to the total size of both suspect and validation sets.

## L. Choice of Classifier

We explore different text classifier backbones and chose the simple 2-layer GPT2-based classifier as our text classifier. Considering the limited number of tokens provided by the author in the DI scenario, stronger text classifiers, such as Llama and full GPT2 architecture can be easily overfitted, especially for SoTA LLM-based text classifiers. We present the results for different architectures with different parameter sizes in Table A8. The results show that the simple GPT2-based classifier with 2 layers and trained from scratch can achieve the best AUC in our experimental settings. Additionally, this simple classifier has a significantly shorter training time, making the method more practical when faced with more queries. In real-world applications, an arbitrator can select the most suitable text classifier based on their specific conditions regarding data size, data type, and computation resources.

## M. Visualization of DI on Single Author Data

In Section 3, we show that there is a distributional shift between the non-member data and held-out data, even for texts composed by a single author. Here, we show this distributional shift in texts also lead to a shift in the MIA score. As presented in Figure A2, the distributional shift in perplexities exists not only between member and held-out sets, but also between *non-member and held-out sets*. This shows that the inherent distributional shift among documents is entangled with the shift caused by membership signals in the MIA score, and makes DI fail to determine membership by simply detecting any distributional shift in the MIA score.

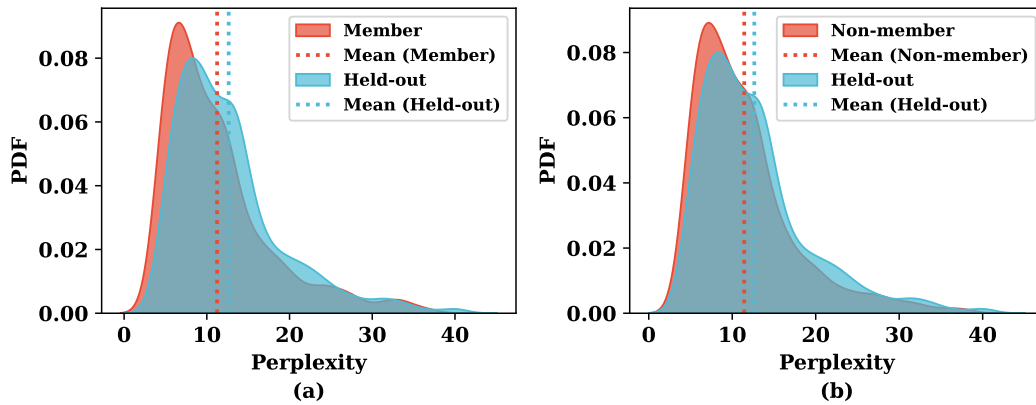


Figure A2. Probability distribution function of target model perplexities on different sets. We show the comparison between (a) the member and held-out, and (b) the non-member and held-out sets.

## N. Algorithm of Our Work

We present the detailed algorithms for our held-out data generation in Algorithm 1, and post-hoc calibration in Algorithm 2.

---

### Algorithm 1 Held-out Data Generation

---

**Require:** Documents  $Doc = \{Doc_1, \dots, Doc_m\}$

**Require:** Hyperparameters: Document number  $m$ , Maximum sequence in each document  $MaxSeq$

**Ensure:** Suspect set  $\mathcal{D}_{\text{sus}}$  and held-out set  $\mathcal{D}_{\text{val}}$  are nearly IID

- 1: Initialize:  $Seq, \mathcal{D}_{\text{sus}}, \mathcal{D}_{\text{val}} = \{\}, \{\}, \{\}$
  - 2: **for** each document  $doc_i \in Doc$  **do**
  - 3:   Segment  $doc_i$  into multiple sequences  $\{seq_i^1, \dots, seq_i^{m_i}\}$
  - 4:   **if**  $m_i < MaxSeq$  **then**
  - 5:      $Seq_i = \{seq_i^1, \dots, seq_i^{m_i}\}$
  - 6:   **else**
  - 7:      $Seq_i =$  randomly sampled  $MaxSeq$  sequences from  $\{seq_i^1, \dots, seq_i^{m_i}\}$
  - 8:   **end if**
  - 9:    $Seq = Seq \cup Seq_i$
  - 10: **end for**
  - 11: Randomly split  $Seq$  into generator training set  $Seq_{\text{train}}$  and generator inference set  $Seq_{\text{test}}$
  - 12: Optimize generator  $g$  on  $Seq_{\text{train}}$  with next-token prediction loss
  - 13: **for** each  $seq_i \in Seq$  **do**
  - 14:    $pre_i, suf_i = \text{Divide}(seq_i)$
  - 15:    $suf'_i = g(pre_i)$
  - 16:    $\mathcal{D}_{\text{sus}} = \mathcal{D}_{\text{sus}} \cup \{(suf_i, 0)\}$
  - 17:    $\mathcal{D}_{\text{val}} = \mathcal{D}_{\text{val}} \cup \{(suf'_i, 1)\}$
  - 18: **end for**
- 

---

### Algorithm 2 Post-hoc Calibration

---

**Require:** Target model  $f$

**Require:** Suspect set  $\mathcal{D}_{\text{sus}}$  and held-out set  $\mathcal{D}_{\text{val}}$  are nearly IID.

- 1: Randomly split  $\mathcal{D}_{\text{sus}}$  into suspect training set  $\mathcal{D}_{\text{sus}}^{\text{train}}$  and suspect test set  $\mathcal{D}_{\text{sus}}^{\text{test}}$
  - 2: Randomly split  $\mathcal{D}_{\text{val}}$  into held-out training set  $\mathcal{D}_{\text{val}}^{\text{train}}$  and held-out test set  $\mathcal{D}_{\text{val}}^{\text{test}}$
  - 3: Optimize a text classifier  $c_{\text{text}}(x)$  on  $\mathcal{D}_{\text{sus}}^{\text{train}} \cup \mathcal{D}_{\text{val}}^{\text{train}}$
  - 4: Optimize a combined classifier  $c_{\text{comb}}(x, \text{MIA}(f(x)))$  on  $\mathcal{D}_{\text{sus}}^{\text{train}} \cup \mathcal{D}_{\text{val}}^{\text{train}}$
  - 5:  $\mathcal{D}_{\text{text}}^{\text{diff}} = \{\}$
  - 6:  $\mathcal{D}_{\text{comb}}^{\text{diff}} = \{\}$
  - 7: **for**  $x_{\text{sus}}^{\text{test}}, x_{\text{val}}^{\text{test}} \in \mathcal{D}_{\text{sus}}^{\text{test}}, \mathcal{D}_{\text{val}}^{\text{test}}$  **do**
  - 8:    $\mathcal{D}_{\text{comb}}^{\text{diff}} = \mathcal{D}_{\text{comb}}^{\text{diff}} \cup \{c_{\text{comb}}(x_{\text{val}}^{\text{test}}, \text{MIA}(f(x_{\text{val}}^{\text{test}}))) - c_{\text{comb}}(x_{\text{sus}}^{\text{test}}, \text{MIA}(f(x_{\text{sus}}^{\text{test}})))\}$
  - 9:    $\mathcal{D}_{\text{text}}^{\text{diff}} = \mathcal{D}_{\text{text}}^{\text{diff}} \cup \{c_{\text{text}}(x_{\text{val}}^{\text{test}}) - c_{\text{text}}(x_{\text{sus}}^{\text{test}})\}$
  - 10: **end for**
  - 11: Compare and  $\mathcal{D}_{\text{comb}}^{\text{diff}}$  and  $\mathcal{D}_{\text{text}}^{\text{diff}}$  with t-test
-