

---

# Membership Inference Attacks against Large Language Models via Self-prompt Calibration

---

**Wenjie Fu**  
Huazhong University of  
Science and Technology  
wjfu99@outlook.com

**Huandong Wang**  
Tsinghua University  
wanghuandong@tsinghua.edu.cn

**Huandong Wang**  
Tsinghua University  
chgao96@gmail.com

**Guanghua Liu**  
Huazhong University of  
Science and Technology  
guanghualiu@hust.edu.cn

**Yong Li**  
Tsinghua University  
liyong07@tsinghua.edu.cn

**Tao Jiang**  
Huazhong University of  
Science and Technology  
taojiang@hust.edu.cn

## Abstract

Membership Inference Attacks (MIA) aim to infer whether a target data record has been utilized for model training or not. Existing MIAs designed for large language models (LLMs) can be bifurcated into two types: reference-free and reference-based attacks. Although reference-based attacks appear promising performance by calibrating the probability measured on the target model with reference models, this illusion of privacy risk heavily depends on a reference dataset that closely resembles the training set. Both two types of attacks are predicated on the hypothesis that training records consistently maintain a higher probability of being sampled. However, this hypothesis heavily relies on the overfitting of target models, which will be mitigated by multiple regularization methods and the generalization of LLMs. Thus, these reasons lead to high false-positive rates of MIAs in practical scenarios. We propose a Membership Inference Attack based on Self-calibrated Probabilistic Variation (SPV-MIA). Specifically, we introduce a self-prompt approach, which constructs the dataset to fine-tune the reference model by prompting the target LLM itself. In this manner, the adversary can collect a dataset with a similar distribution from public APIs. Furthermore, we introduce probabilistic variation, a more reliable membership signal based on LLM memorization rather than overfitting, where we rediscover neighbour attack with a more rigorous paradigm. Comprehensive evaluation conducted on three datasets and four exemplary LLMs shows that SPV-MIA raises the AUC of MIAs from 0.7 to a significantly high level of 0.9.<sup>1</sup>

## 1 Introduction

Large language models (LLMs) have been validated to have the ability to generate extensive, creative, and human-like responses when provided with suitable input prompts. Both commercial LLMs (e.g., ChatGPT [43]) and open-source LLMs (e.g., LLaMA [53]) can easily handle various complex application scenarios, including but not limited to chatbots [14], code generation [54], article co-writing [21]. Moreover, as the pretraining-finetuning paradigm becomes the mainstream pipeline in of LLM field, small-scale organizations and individuals can fine-tune pre-trained models over their private datasets for downstream applications [36], which further enhances the influence of LLMs.

---

<sup>1</sup>Our code and dataset are available at <https://github.com/wjfu99/MIA-LLMs>

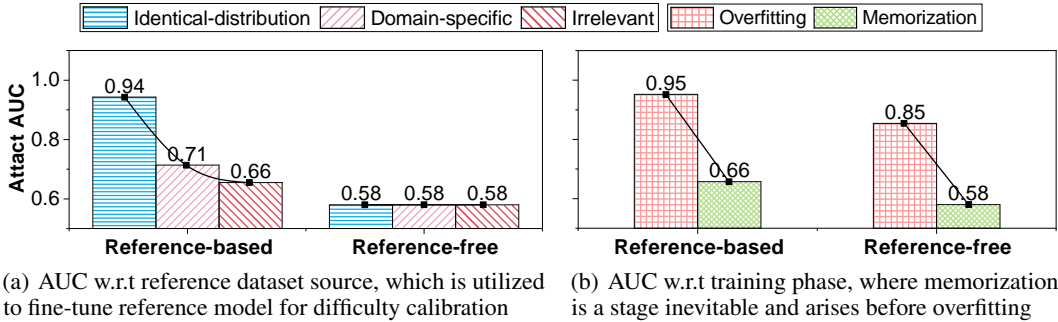


Figure 1: Attack performances of the reference-based MIA (LiRA [38, 39]) and reference-free MIA (LOSS Attack [60]) unsatisfy against LLMs in practical scenarios, where LLMs are in the memorization stage and only domain-specific dataset is available. (a) Reference-based MIA shows an exponential decrease in performance when the similarity between the reference and training datasets declines. (b) Existing MIAs are unable to pose privacy leakages on LLMs that only exhibit memorization.

However, while we enjoy the revolutionary benefits raised by the popularization of LLMs, we also have to face the potential privacy risks associated with LLMs. Existing work has unveiled that the privacy leakage of LLMs exists in almost all stages of the LLM pipeline [44]. For example, poisoning attacks can be deployed during pre-training, distillation, and fine-tuning [27, 55]. Moreover, data and model extraction attacks can be conducted through inference [4, 16]. Among these attacks, fine-tuning is widely recognized as the stage that is most susceptible to privacy leaks since the relatively small and often private datasets used for this process [61]. Therefore, this paper aims to uncover the underlying privacy concerns associated with fine-tuned LLMs through an exploration of the membership inference attack (MIA).

MIA is an adversary model that categorizes data records into two groups: member records, which are included in the training dataset of the target model, and nonmember records, which belong to a disjoint dataset [50]. MIAs have been well studied in classic machine learning tasks, such as classification, and reveal significant privacy risks [19]. Recently, some contemporaneous works attempt to utilize MIAs to evaluate the privacy risks of LLMs. For example, several studies have employed reference-free attack, especially LOSS attack [60], for privacy auditing [24] or more sophisticated attack [4]. Mireshghallah et al. introduce a reference-based attack, Likelihood Ratio Attacks (LiRA), on Masked Language Models (MLMs), which measure the calibrated likelihood of a specific record by comparing the discrepancy on the likelihood between the target LLM and the reference LLM. Following this concept, Mireshghallah et al. further adapt LiRA for analyzing memorization in Causal Language Models (CLMs). However, these methods heavily rely on several over-optimistic assumptions, including assuming the overfitting of target LLMs [34] and having access to a reference dataset from the same distribution as the training dataset [38, 39]. Thus, it remains inconclusive whether prior MIAs can cause considerable privacy risk in practical scenarios.

As illustrated in Fig. 1, LiRA [39] and LOSS Attack [60] are employed to represent reference-based and reference-free MIAs to explore their performance in practical scenarios. Firstly, as shown in Fig. 1(a), we evaluate LiRA and LOSS Attack with three reference datasets from different sources, i.e., the dataset with the identical distribution with the member records (identical-distribution), the dataset of the same domain with the member records (domain-specific), and the dataset irrelevant to the member records (irrelevant). The performance of LOSS attack is consistently low and independent of the source of the reference dataset. For LiRA, the attack performance will exponentially decline as the similarity between the reference dataset and the target dataset declines. Thus, the reference-based MIA can not pose critical privacy leakage on LLMs since similar datasets are usually not available to adversaries in real applications. Secondly, as shown in Fig. 1(b), two target LLMs are fine-tuned over the same pre-trained model but stop before and after overfitting, and the reference LLMs are fine-tuned on a different dataset from the same domain. We can observe that existing MIAs cannot effectively cause privacy leaks when the LLM is not overfitting. This phenomenon is addressed by the fact that the membership signal proposed by existing MIAs is highly dependent on overfitting in target LLMs. They assume that member records tend to have overall higher probabilities of being sampled than non-member ones, an assumption that is only satisfied in overfitting models [7].

In this work, to address the aforementioned two limitations of existing works, we propose a Membership Inference Attack based on Self-calibrated Probabilistic Variation (SPV-MIA) composed of two according modules. First, although existing reference-based MIAs are challenging to reveal actual privacy risks, they demonstrate the significant potential of achieving higher privacy risks with the reference model. Therefore, we design a self-prompt approach to extract the reference dataset by prompting the target LLMs themselves and collecting the generated texts. This approach allows us to acquire the significant performance improvement brought by the reference model while ensuring the adversary model is feasible on the practical LLMs. Second, prior work demonstrates that memorization is inevitable for machine learning models to arrive optimal [11] and can exist without overfitting in LLMs [52]. Thus, instead of utilizing probabilities as membership signals, we opt to identify member records based on memorization that will arise as an increased tendency in probability distribution around the member records [7]. We proposed a probabilistic variation metric that can detect local maxima points via a second partial derivative test [51] instantiated by a paraphrasing model. Furthermore, based on the new metric, we rediscover the neighbour attack [34] and reformulate it in a more rigorous manner. It is worth noting that our paraphrasing model does not rely on another LLM like the neighbour attack. Overall, our contributions are summarized as follows:

- We propose a self-prompt approach that collects reference datasets by prompting the target LLM to generate, which will have the closely resemble distribution as the training dataset. In this manner, the reference model fine-tuned on the reference dataset can significantly improve the attack performance without any unrealistic assumptions.
- We further design a probabilistic variation metric that detects the essential characteristics of member records memorized by LLMs. Then we provide a more rigorous principle and explanation of the neighbour attack [34] based on the proposed metric.
- We conducted extensive experiments to validate the effectiveness of SPV-MIA. The results suggest that SPV-MIA unveils significantly higher privacy risk across multiple fine-tuned LLMs and datasets compared with existing MIAs (about 23.6% improvement in AUC across four representative LLMs and three datasets).

## 2 Related Works

**Membership Inference Attack:** Initially, prior MIAs mainly focused on classical machine learning models, such as classification models [50, 6, 32, 5]. With the rapid development of other machine learning tasks, such as recommendation and generation tasks, MIAs against these task-specific models became a research direction of great value, and have been well investigated [62, 9, 13]. Meanwhile, ChatGPT released by OpenAI has propelled the attention towards LLMs to the peak over the past year, which promotes the study of MIAs against LLMs. Mireshghallah et al. proposed LiRA against MLMs via adopting pre-trained models as reference models. Following this study, Mireshghallah et al. further adapted LiRA for CLMs. Mattern et al. pointed out the unrealistic assumption of a reference model trained on similar data, then substituted it with a neighbourhood comparison method. Although MIAs against LMs have been studied by several works, the attack performance of existing MIAs in regard to LLMs with large-scale parameters and pre-trained on tremendous corpora is still not clear. We evaluate previous MIAs on LLMs in practical scenarios, and found that the revealed privacy breaches were far below expectations due to their strict requirements and over-optimistic assumptions. Then, we propose SPV-MIA, which discloses significant privacy risks on practical LLM applications.

**Large Language Models:** In the past year, LLMs have dramatically improved performances on multiple natural language processing (NLP) tasks and consistently attracted attention in both academic and industrial circles [36]. The widespread usage of LLMs has led to much other contemporaneous work on quantifying the privacy risks of LLMs [34, 40, 44]. In this work, we audit privacy leakages of LLMs by distinguishing whether or not a specific data record is used for fine-tuning the target LLM. The existing LLMs primarily fall into three categories: causal language modeling (CLM) (e.g. GPT), masked language modeling (MLM) (e.g. BERT), and Sequence-to-Sequence (Seq2Seq) approach (e.g. BART). Among these LLMs, CLMs such as GPT [45, 56] and LLaMA [53] have achieved the dominant position with the exponential improvement of model scaling [64]. Therefore, we select CLM as the representative LLM in this work for evaluation.

### 3 Preliminaries

#### 3.1 Causal Language Models

For a given text record  $\mathbf{x}$ , it can be split into a sequence of tokens  $[t_0, t_1, \dots, t_{|\mathbf{x}|}]$  with variable length  $|\mathbf{x}|$ . CLM is an autoregressive language model, which aims to predict the conditional probability  $p_\theta(t_i | \mathbf{x}_{<i})$  given the previous tokens  $\mathbf{x}_{<i} = [t_0, t_1, \dots, t_{i-1}]$ . During the training process, CLM calculates the probability of each token in a text with the previous tokens, then factorizes the joint probability of the text into the product of conditional token prediction probabilities. Therefore, the model can be optimized by minimizing the negative log probability:

$$\mathcal{L}_{\text{CLM}} = -\frac{1}{M} \sum_{j=1}^M \sum_{i=1}^{|\mathbf{x}^{(j)}|} \log p_\theta(t_i | \mathbf{x}_{<i}^{(j)}), \quad (1)$$

where  $M$  denotes the number of training records. In the process of generation, CLMs can generate coherent words by predicting one token at a time and producing a complete text using an autoregressive manner. Moreover, the pretraining-finetuning paradigm is proposed to mitigate the uncountable demands of training an LLM for a specific task [36]. Besides, multifarious parameters-efficient fine-tuning methods (e.g., LoRA [18], P-Tuning [31]) are introduced to further decrease consumption by only fine-tuning limited model parameters [8]. In this work, we concentrate on the fine-tuning phase, since the fine-tuning datasets are usually more private and vulnerable to the adversary [61].

#### 3.2 Threat Model

In this work, we consider an adversary who aims to infer whether a specific text record was included in the fine-tuning dataset of the target LLM. There are two mainstream scenarios investigated by previous research: white-box and black-box MIAs. White-box MIA assumes full access to the raw copy of the target model, which means the adversary can touch and modify each part of the target model [42]. For black-box, the adversary only approved to acquire the response results (e.g. generated texts, log probabilities) by requesting the provided service API [47], which is more realistic and aligned with practical application circumstances. Thus, we adopt the black-box scenario for evaluating existing works and our proposed method.  $D$  is a dataset collected for a specific task, which can be separated into two disjoint subsets:  $D_{mem}$  and  $D_{non}$ . The target LLM  $\theta$  is fine-tuned on  $D_{mem}$ , and the adversary has no prior information about which data records are utilized for fine-tuning. Besides, all reference-based MIA, including SPV-MIA, can at most fine-tune the reference model using a disjoint dataset  $D_{refer}$  from the same task. The adversary algorithm  $\mathcal{A}$  is designed to infer whether a text record  $\mathbf{x}^{(i)} \in D$  belong to the training dataset  $D_{mem}$ :

$$\mathcal{A}(\mathbf{x}^{(j)}, \theta) = \mathbb{1} \left[ P(m^{(j)} = 1 | \mathbf{x}^{(j)}, \theta) \geq \tau \right], \quad (2)$$

where  $m^{(j)} = 1$  indicates that the record  $\mathbf{x}^{(j)} \in D_{mem}$ ,  $\tau$  represents the threshold, and  $\mathbb{1}$  denotes the indicator function.

### 4 Membership Inference Attack via Self-calibrated Probabilistic Variation

In this section, we first introduce the general paradigm of Membership Inference Attack via Self-calibrated Probabilistic Variation (SPV-MIA) as illustrated in Fig 2. Then we discuss the detailed algorithm instantiations of this general paradigm by introducing practical difficulty calibration (PDC, refer to Section 4.2) and probabilistic variation assessment (PVA, refer to Section 4.3).

#### 4.1 General Paradigm

As formulated in Eq. 1, the objective of an LLM is to maximize the joint probability of the text in the training set. Thus, prior **reference-free MIAs** employ the joint probability of the target text being sampled as the membership signal [4, 24, 49]:

$$\mathcal{A}(\mathbf{x}, \theta) = \mathbb{1} [p_\theta(\mathbf{x}) \geq \tau], \quad (3)$$

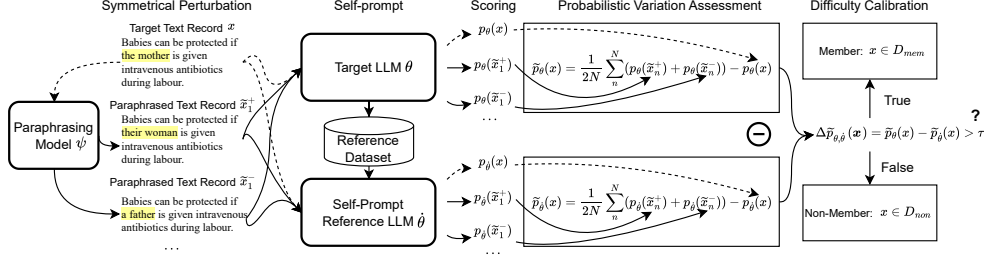


Figure 2: The overall workflow of SPV-MIA, where includes the probabilistic calibration via self-prompt reference model and the probabilistic variation assessment via paraphrasing model.

where  $p_\theta(\mathbf{x})$  denotes the probability measured on the target model  $\theta$ . Since some records are inherently over-represented, even non-member records can achieve high probability in the data distribution [57], which leads to a high FPR. Thus, **reference-based MIAs** adopt difficulty calibration [57], which further calibrates the probability by comparing it with the value measured on reference models [39, 38]:

$$\mathcal{A}_{exist}(\mathbf{x}, \theta) = \mathbb{1}[\Delta p_\theta(\mathbf{x}) \geq \tau] = \mathbb{1}[p_\theta(\mathbf{x}) - p_\phi(\mathbf{x}) \geq \tau], \quad (4)$$

where  $\Delta p_\theta(\mathbf{x})$  is the calibrated probability, and  $p_\phi(\mathbf{x})$  is estimated on the reference model  $\phi$ .

However, both reference-free and reference-based MIAs often encounter high FPR in practical scenarios. For reference-based MIAs, although they have the potential to offset the over-represented statuses of data records if the reference model can be trained on a dataset closely resembling the training dataset  $D_{mem}$ . Nevertheless, it is almost unrealistic for an adversary to obtain such a dataset, and adopting a compromising dataset will lead to the collapse of attack performance. We circumvent this by introducing a self-prompt reference model  $\hat{\theta}$ , which is trained on the generated text of the target model  $\theta$ . Besides, the probability signal adopted by existing MIAs is not reliable, since the confidence of the probability signal is notably declined when the target model is not overfitting only memorization [7]. Thus, we elaborately design a more stable membership signal, probabilistic variation  $\tilde{p}_\theta(\mathbf{x})$ , audited by a paraphrasing model and only rely on LLM memorization. Formally, as depicted in Fig. 2, **our proposed MIAs** can be formulated as:

$$\mathcal{A}_{our}(\mathbf{x}, \theta, \hat{\theta}) = \mathbb{1}[\Delta \tilde{p}_{\theta, \hat{\theta}}(\mathbf{x}) \geq \tau] = \mathbb{1}[\tilde{p}_\theta(\mathbf{x}) - \tilde{p}_{\hat{\theta}}(\mathbf{x}) \geq \tau], \quad (5)$$

where  $\tilde{p}_\theta(\mathbf{x})$  and  $\tilde{p}_{\hat{\theta}}(\mathbf{x})$  are probabilistic variations of the text record  $\mathbf{x}$  measured on the target model  $\theta$  and the self-prompt reference model  $\hat{\theta}$  respectively.

## 4.2 Practical Difficulty Calibration via Self-prompt Reference Model

Watson et al. has suggested that inferring the membership of a record by thresholding on a predefined metric (e.g. confidence [48], loss [60], and gradient norm [42]) will cause a high FPR. Since several non-member records may have high probabilities of being classified as member records simply because they are inherently over-represented in the data manifold. In other words, the metric estimated on the target model is inherently biased and has a high variance, which leads to a significant overlap in the metric distributions between members and non-members, making them more indistinguishable. To mitigate this phenomenon, Watson et al. propose difficulty calibration as a general approach for extracting a much more distinguishable membership signal, which can be adapted to most metric-based MIAs by constructing their calibrated variants [57, 38, 38]. Concretely, difficulty calibration assumes an ideal reference dataset  $D_{refer}$  drawn from the identical distribution as the training set  $D_{mem}$  of the target model  $\theta$ , and trains an ideal reference model  $\phi$  with a training algorithm  $\mathcal{T}$ . Then, it fabricates a calibrated metric by measuring the discrepancy between metrics on the target model and reference model, and this can offset biases on membership signals caused by some over-represented records. The calibrated metric is defined as:

$$\Delta m(\mathbf{x}) = m_\theta(\mathbf{x}) - \mathbb{E}_{\phi \leftarrow \mathcal{T}(D_{refer})}[m_\phi(\mathbf{x})], \quad (6)$$

where  $\Delta m(\mathbf{x})$  is the calibrated metric,  $m_\theta(\mathbf{x})$  and  $m_\phi(\mathbf{x})$  are metrics measured on target and reference models, respectively. The existing study has verified that reference-based MIA highly depends on the similarity of training and reference dataset [34]. A low-quality dataset will lead to an

exponential decrease in attack performance. However, the dataset used for fine-tuning an LLM is typically highly private, making extracting a high-quality reference dataset from the same distribution a non-trivial challenge.

We notice that LLMs possess revolutionary fitting and generalization capabilities, enabling them to generate a wealth of creative texts. Therefore, LLMs themselves have the potential to depict the distribution of the training data. Thus, we consider a self-prompt approach that collects the reference dataset from the target LLM itself by prompting it with few words. Concretely, we first collect a set of text chunks with an equal length of  $l$  from a public dataset from the same domain, where the domain can be easily inferred from the task of the target LLM (e.g., An LLM that serves to summary task has high probability using a summary fine-tuning dataset). Then, we utilize each text chunk of length  $l$  as the prompt text and request the target LLM to generate text. All the generated text can form a dataset of size  $N$ , which is used to fine-tune the proposed self-prompt reference model  $\phi$  over the pre-trained model. Accordingly, we can define the practical difficulty calibration as:

$$\Delta m(\mathbf{x}) = m_\theta(\mathbf{x}) - \mathbb{E}_{\hat{\theta} \leftarrow \mathcal{T}(\mathcal{D}_{self})}[m_{\hat{\theta}}(\mathbf{x})] \approx m_\theta(\mathbf{x}) - m_{\hat{\theta}}(\mathbf{x}), \quad (7)$$

where  $\mathcal{D}_{self} \sim p_\theta(\mathbf{x})$ ,  $m_\theta(\mathbf{x})$  and  $m_\phi(\mathbf{x})$  are membership metrics measured over the target model and the self-prompt reference model. Only one reference model is used for computational efficiency, which can achieve sufficiently high attack performance. It is worth noting that in some challenging scenarios where acquiring domain-specific datasets is difficult, our self-prompt method can still effectively capture the underlying data distribution, even when using completely unrelated prompt texts. The relevant experiments will be conducted and discussed in detail in Section 5.4.

### 4.3 Probabilistic Variation Assessment via Symmetrical Paraphrasing

In contrast to overfitting, memorization has been verified as an inevitable phenomenon for achieving optimal generalization on machine learning models [12], and it will exist before overfitting in LLMs [52]. Therefore, it will naturally be a more reliable signal for detecting member text. Memorization in generative models will cause member records to have a higher probability of being generated than neighbour records in the data distribution [7]. This principle can be shared with LLMs, as they can be considered generation models for texts. Thus, we suggest designing a more promising membership signal that can measure a value for each text record to identify whether this text is located on the local maximum in the sample distribution characterized by  $\theta$ . The second partial derivative test is an approach in multivariable calculus commonly employed to ascertain whether a critical point of a function is a local minimum, maximum, or saddle point [51]. For our objective of identifying maximum points, we need to confirm if the Hessian matrix is negative definite, meaning that all the directional second derivatives are negative. Thus, we define the probabilistic variation mentioned in Eq. 5 as the expectation of the directional derivative:

$$\tilde{p}_\theta(\mathbf{x}) := \mathbb{E}_{\mathbf{z}}(\mathbf{z}^\top H_p(\mathbf{x}) \mathbf{z}), \quad (8)$$

where  $H_p(\cdot)$  is the hessian matrix of the probability function  $p_\theta(\cdot)$ , then  $\mathbf{z}^\top H_p(\mathbf{x}) \mathbf{z}$  indicates the second-order directional derivative of  $p_\theta(\cdot)$  with respect to the text record  $x$  in the direction  $\mathbf{z}$ . Then, we further approximate the derivative with the symmetrical form [20]:

$$\mathbf{z}^\top H_p(\mathbf{x}) \mathbf{z} \approx \frac{p_\theta(\mathbf{x} + h\mathbf{z}) + p_\theta(\mathbf{x} - h\mathbf{z}) - 2p_\theta(\mathbf{x})}{h^2}, \quad (9)$$

where requires  $h \rightarrow 0$ , and  $\mathbf{z}$  can be interpreted as kind of "noise". Thus,  $\mathbf{x} \pm h\mathbf{z}$  can be considered as a pair of symmetrical adjacent text records of  $\mathbf{x}$  in the data distribution. Then we can reformulate Eq. 8 as follows by omitting coefficient  $h$ :

$$\tilde{p}_\theta(\mathbf{x}) \approx \frac{1}{2N} \sum_n^N \left( p_\theta(\tilde{\mathbf{x}}_n^+) + p_\theta(\tilde{\mathbf{x}}_n^-) \right) - p_\theta(\mathbf{x}). \quad (10)$$

where  $\tilde{\mathbf{x}}_n^\pm = \mathbf{x} \pm \mathbf{z}_n$  is a symmetrical text pair sampled by a paraphrasing model, which slightly paraphrases the original text  $\mathbf{x}$  in the high-dimension space (as Eq. 9 requires  $h \rightarrow 0$ , which means the paraphrasing should be modest.) Based on the aforementioned discussions, we designed two different paraphrasing models in the embedding domain and the semantic domain, respectively, to generate symmetrical paraphrased text embeddings or texts. For the embedding domain, we first

Method	Wiki					AG News					Xsum				
	GPT-2	GPT-J	Falcon	LLaMA	Avg.	GPT-2	GPT-J	Falcon	LLaMA	Avg.	GPT-2	GPT-J	Falcon	LLaMA	Avg.
Loss Attack	0.614	0.577	0.593	0.605	0.597	0.591	0.529	0.554	0.580	0.564	0.628	0.564	0.577	0.594	0.591
Neighbour Attack	0.647	0.612	0.621	0.627	0.627	0.622	0.587	0.594	0.610	0.603	0.612	0.547	0.571	0.582	0.578
DetectGPT	0.623	0.587	0.603	0.619	0.608	0.611	0.579	0.582	0.603	0.594	0.603	0.541	0.563	0.577	0.571
LiRA-Base	0.710	0.681	0.694	0.709	0.699	0.658	0.634	0.641	0.657	0.648	0.776	0.718	0.734	0.759	0.747
LiRA-Candidate	0.769	0.726	0.735	0.748	0.744	0.717	0.690	0.708	0.714	0.707	0.823	0.772	0.785	0.809	0.797
SPV-MIA	<b>0.975</b>	<b>0.929</b>	<b>0.932</b>	<b>0.951</b>	<b>0.938</b>	<b>0.949</b>	<b>0.885</b>	<b>0.898</b>	<b>0.903</b>	<b>0.909</b>	<b>0.944</b>	<b>0.897</b>	<b>0.918</b>	<b>0.937</b>	<b>0.924</b>

Table 1: AUC for detecting member texts from four LLMs across three datasets for SPV-MIA and five previously proposed methods. **Bold** and Underline respectively represent the best and the second-best results within each column (model-dataset pair).

embed the target text, then randomly sample noise following Gaussian distribution, and obtain a pair of symmetrical paraphrased texts by adding/subtracting noise. For the semantic domain, we randomly mask out 20% tokens in each target text, then employ T5-base to predict the masked tokens. Then, we compute the difference in the embeddings between the original tokens and predicted tokens to search for tokens that are symmetrical to predicted tokens with respect to the original tokens. We provide the detailed pseudo codes of both two paraphrasing models in Appendix A.3. In subsequent experiments, we default to paraphrasing in the semantic domain. Furthermore, we reformulate the neighbour attack and provide another explanation of its success based on the probabilistic variation metric with a more rigorous principle (refer to Appendix A.4). Additionally, supplementary experiments demonstrate that our proposed paraphrasing model in the embedding domain achieves considerable performance gains without relying on another MLM.

## 5 Experiments

### 5.1 Experimental Setup

Our experiments are conducted on four open-source LLMs: **GPT-2** [45], **GPT-J** [56], **Falcon-7B** [2] and **LLaMA-7B** [53], which are both fine-tuned over three dataset across multiple domains and LLM use cases: **Wikitext-103** [35], **AG News** [63] and **XSum** [41]. Each target LLM is fine-tuned with the training batch size of 16, and trained for 10 epochs. Each self-prompt reference model is trained for 4 epochs. We adopt LoRA [18] as the default fine-tuning method. The learning rate is set to 0.0001. We adopt the AdamW optimizer [33] and early stopping [58] to achieve the generalization of LLMs and avoid overfitting, the NLP performance of each LLM-dataset pair is provided in Appendix A.5.3. We compare SPV-MIA with six state-of-the-art MIAs designed for LMs, including three reference-free MIAs: **Loss Attack** [60], **Neighbour Attack** [34], **DetectGPT** [40] and two reference-based MIAs: **LiRA-Base** [39], **LiRA-Candidate** [39]. We defer the detailed setup information to Appendix A.6.

### 5.2 Overall Performance

As shown in Table 1, we first summarize the AUC scores [3] for all baselines and SPV-MIA against four LLMs across three datasets. Furthermore, we present receiver operating characteristic (ROC) curves for SPV-MIA and the top three best baselines on LLaMAs in Appendix A.5.2 for a more comprehensible presentation. The results demonstrate that SPV-MIA achieves the best overall attack performance with the highest average AUC of 0.924 over all scenarios. Furthermore, compared to the most competitive baseline, LiRA-Candidate, SPV-MIA has improved the AUC of the attack by 30%, even LiRA-Candidate assumes full access to the auxiliary dataset while SPV-MIA only needs some short text chunks from this dataset. This phenomenon indicates that our proposed self-prompt approach enables the reference model to gain a deeper understanding of the data distribution, thereby serving as a more reliable calibrator. Most baseline, especially reference-free attack methods, yield a low AUC, which is only slightly better than random guesses. Furthermore, their performances on larger-scale LLMs are worse. This phenomenon verifies the claim that existing MIAs designed for LMs can not handle LLMs with large-scale parameters. It is also worth noting that the privacy risks caused by MIAs are proportional to the overall parameter scale and language capabilities of LLMs. We interpret this phenomenon as follows: LLMs with stronger overall NLP performance have better learning ability, which means they are more likely to memorize records from the training set. Besides, MIAs fundamentally leverage the memorization abilities of machine learning models, making superior models more vulnerable to attacks.

### 5.3 How MIAs Rely on Reference Dataset Quality

In this work, a key contribution is introducing a self-prompt approach for constructing a high-quality dataset to fine-tune the reference model, which guides the reference model to become a better calibrator. Therefore, we conduct experiments to investigate how prior reference-based MIAs rely on the quality of the reference dataset, and evaluate whether our proposed method can build a high-quality reference dataset. In real-world scenarios, based on different prior information, adversaries can obtain datasets from different sources to fine-tune the reference model with uneven quality. We categorize the reference dataset into three types based on their relationship with the fine-tuning dataset of the target model and sort them in ascending order of difficulty in acquisition: 1) **Irrelevant** dataset, 2) **Domain-specific** dataset, and 3) **Identical distribution** dataset. Besides, the dataset extracted by the self-prompt approach is denoted as 4) **Self-prompt** dataset. The detailed information of these datasets is summarized in Appendix A.6. Then, we conduct MIAs with the aforementioned four data sources and summarize the results in Fig. 3. The experimental results indicate that the performance of MIA shows a noticeable decrease along the Identical, Domain, and Irrelevant datasets. This illustrates the high dependency of previous reference-based methods on the quality of the reference dataset. However, AUC scores on self-prompt reference datasets are only marginally below Identical datasets. It verifies that our proposed self-prompt method can effectively leverage the creative generation capability of LLMs, approximate sampling high-quality text records indirectly from the distribution of the target training set.

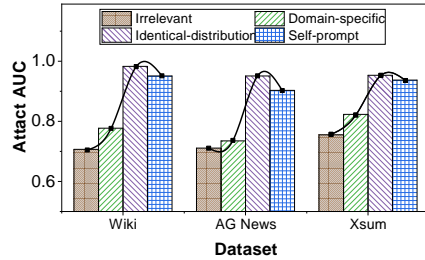


Figure 3: The performances of reference-based MIA on LLaMA while utilizing different reference datasets.

### 5.4 The Robustness of SPV-MIA in Practical Scenarios

We have verified that SPV-MIA can provide a high-quality reference model. However, the source and scale of self-prompt texts may face various limitations in practical scenarios. Therefore, we conducted experiments to verify the robustness of SPV-MIA performance in diverse practical scenarios.

**Source of Self-prompt Texts.** The sources of self-prompt texts available to attackers are usually limited by the actual deployment environment, and sometimes even domain-specific texts may not be accessible. Compared with using domain-specific text chunks for prompting, we also evaluate the self-prompt approach with irrelevant and identical-distribution text chunks. As shown in Fig. 4, the self-prompt method demonstrates an incredibly lower dependence on the source of the prompt texts. We found that even when using completely unrelated prompt texts, the performance of the attack only experiences a slight decrease (3.6% at most). This phenomenon indicates that the self-prompt method we proposed has a high degree of versatility across adversaries with different prior information.

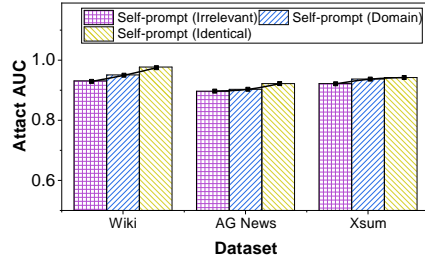


Figure 4: The performances of SPV-MIA on LLaMA while utilizing different prompt text sources.

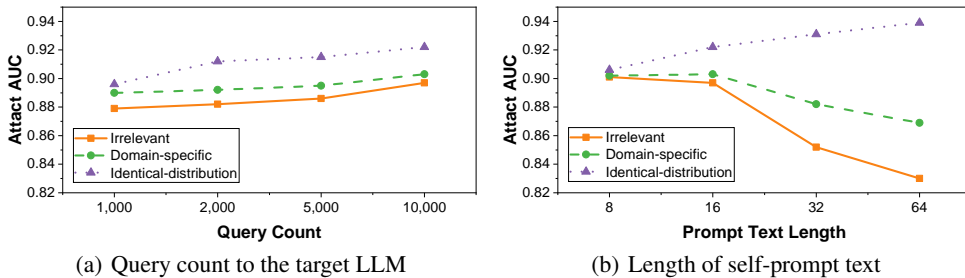


Figure 5: The performances of SPV-MIA on LLaMA while utilizing different query count to the target model and different prompt text lengths.



**Scale of Self-prompt Texts.** In real-world scenarios, the scale of self-prompt texts is usually limited by the access frequency cap of the LLM API and the number of available self-prompt texts. Thus, we set up two sets of experiments to verify the sensitivity of SPV-MIA to the aforementioned limitations. As shown in Fig. 5(a), our self-prompt reference model is minimally affected by the access frequency limitations of the target LLM. Even with only 1,000 queries, it achieves performance comparable to 10,000 queries. As shown in Fig. 5(b), even when the self-prompt texts are severely limited (with only 8 prompt tokens), the attack performance remains at a startlingly high level of 0.9. Besides, texts from different sources show varying attack performance trends based on text length. From the identical dataset, attack performance increases with text length. From the domain-specific dataset, it initially increases then decreases. From an irrelevant dataset, it decreases with longer texts. Therefore, we recommend setting smaller text lengths to allow LLMs to generate samples that are close to data distributions of training sets, unless adversaries can directly sample texts from the same data distribution as the training set. Overall, our proposed method can maintain stable attack performance in practical scenarios where the scale of self-prompt texts is limited.

### 5.5 Defending against SPV-MIAs

As privacy risks emerge from various attacks, including data extraction attack [4], model extraction attack [16], and membership inference attack [59, 50, 34], the research community actively promotes defending methods against these attacks [22, 37]. DP-SGD [1] is one of the most widely adopted defense methods based on differential privacy [10] to provide mathematical privacy guarantees. Through DP-SGD, the amount of information the parameters have about a single data record is bound. Thus, the privacy leakage will not exceed the upper bound, regardless of how many outputs we obtain from the target model. We follow the same manner as the existing study [29] and train LLaMA with DP-Adam on the three datasets. The results are summarized in Table 2, where we choose a set of appropriate  $\epsilon$  as existing works suggest that higher DP guarantees lead to a noticeable performance degradation [34, 15]. The performances of LLMs are supplemented in Appendix A.5.3. The results demonstrate that DP-SGD can reduce the privacy risk to a certain, but under moderate privacy budgets, SPV-MIA still poses a significant risk of privacy leakage.

Privacy Budget $\epsilon$	15	30	60	+ inf
Wiki	0.785	0.832	0.875	0.951
AG News	0.766	0.814	0.852	0.903
Xsum	0.771	0.827	0.867	0.937
Avg.	0.774	0.824	0.865	0.930

Table 2: The AUC performance of SPV-MIA against LLaMA fine-tuned with DP-SGD w.r.t different privacy budget  $\epsilon$ .

#### 5.5.1 Impact of Fine-tuning Methods

We further evaluated the generalizability of SPV-MIA under different Parameter-Efficient Fine-Tuning (PEFT) techniques. As shown in Table 3, SPV-MIA can maintain a high-level AUC across all PEFT techniques. Besides, the performance of MIA is positively correlated with the number of trainable parameters during the fine-tuning process. We hypothesize that this is because as the number of trainable parameters increases, LLMs retain more complete memory of the member records, making them more vulnerable to attacks.

PEFT	LoRA	Prefix Tuning	P-Tuning	(IA) <sup>3</sup>
# Parameters (M)	33.55	5.24	1.15	0.61
Wiki	0.951	0.943	0.922	0.914
Ag News	0.903	0.897	0.879	0.873
Xsum	0.937	0.931	0.924	0.911

Table 3: The AUC Performance of SPV-MIA across LLaMAs fine-tuned with different PEFT techniques over three datasets. We choose LoRA [18], Prefix Tuning [28], P-Tuning [31] and (IA)<sup>3</sup> [30] as four representative PEFT techniques.

## 6 Conclusion

In this paper, we reveal the under-performances of existing MIA methods against LLMs for practical applications and interpret this phenomenon from two perspectives. First, reference-based attacks seem to pose impressive privacy leakages by comparing the sampling probabilities of the target record between target and reference LLMs, but the inaccessibility of the appropriate reference dataset will be a big obstacle to deploying it in practice. Second, existing MIAs heavily rely on overfitting, which is usually avoided before releasing LLM for public access. To address these limitations, we propose a Membership Inference Attack based on Self-calibrated Probabilistic Variation (SPV-MIA), where we

propose a self-prompt approach to extract reference dataset from LLM itself in a practical manner, then introduce a more reliable membership signal based on memorization rather than overfitting. We conduct substantial experiments to validate the superiority of SPV-MIA over all baselines and verify its effectiveness in extreme conditions. One primary limitation of this study is that SPV-MIA is only designed for CLM, we leave the adaption on other LLMs as the future work.

## References

- [1] Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep Learning with Differential Privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, CCS '16*, pages 308–318, New York, NY, USA, October 2016. Association for Computing Machinery. ISBN 978-1-4503-4139-4.
- [2] Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, et al. Falcon-40b: an open large language model with state-of-the-art performance. Technical report, Technical report, Technology Innovation Institute, 2023.
- [3] Andrew P Bradley. The use of the area under the roc curve in the evaluation of machine learning algorithms. *Pattern recognition*, 30(7):1145–1159, 1997.
- [4] Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650, 2021.
- [5] Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramèr. Membership Inference Attacks From First Principles. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1897–1914, May 2022.
- [6] Christopher A. Choquette-Choo, Florian Tramèr, Nicholas Carlini, and Nicolas Papernot. Label-Only Membership Inference Attacks. In *Proceedings of the 38th International Conference on Machine Learning*, pages 1964–1974. PMLR, July 2021.
- [7] Gerrit J. J. Van den Burg and Chris Williams. On Memorization in Probabilistic Deep Generative Models. In *Advances in Neural Information Processing Systems*, November 2021.
- [8] Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, Jing Yi, Weilin Zhao, Xiaozhi Wang, Zhiyuan Liu, Hai-Tao Zheng, Jianfei Chen, Yang Liu, Jie Tang, Juanzi Li, and Maosong Sun. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235, March 2023. ISSN 2522-5839.
- [9] Jinhao Duan, Fei Kong, Shiqi Wang, Xiaoshuang Shi, and Kaidi Xu. Are Diffusion Models Vulnerable to Membership Inference Attacks? In *Proceedings of the 38th International Conference on Machine Learning, {ICML} 2023*. PMLR, February 2023.
- [10] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, pages 265–284. Springer, 2006.
- [11] Vitaly Feldman. Does learning require memorization? a short tale about a long tail. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing, STOC 2020*, pages 954–959, New York, NY, USA, June 2020. Association for Computing Machinery. ISBN 978-1-4503-6979-4.
- [12] Vitaly Feldman and Chiyuan Zhang. What Neural Networks Memorize and Why: Discovering the Long Tail via Influence Estimation. *Advances in Neural Information Processing Systems*, 33:2881–2891, 2020.

- [13] Wenjie Fu, Huandong Wang, Chen Gao, Guanghua Liu, Yong Li, and Tao Jiang. A probabilistic fluctuation based membership inference attack for diffusion models. *arXiv e-prints*, pages arXiv–2308, 2023.
- [14] Stephen Gilbert, Hugh Harvey, Tom Melvin, Erik Vollebregt, and Paul Wicks. Large language model ai chatbots require approval as medical devices. *Nature Medicine*, pages 1–3, 2023.
- [15] Jamie Hayes, Luca Melis, George Danezis, and Emiliano De Cristofaro. Logan: Membership inference attacks against generative models. *Proceedings on Privacy Enhancing Technologies*, (1):133–152, 2019.
- [16] Xuanli He, Lingjuan Lyu, Lichao Sun, and Qionghai Xu. Model extraction and adversarial transferability, your bert is vulnerable! In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2006–2012, 2021.
- [17] Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. Teaching machines to read and comprehend. *Advances in neural information processing systems*, 28, 2015.
- [18] Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021.
- [19] Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S. Yu, and Xuyun Zhang. Membership Inference Attacks on Machine Learning: A Survey. *ACM Computing Surveys*, 54 (11s):235:1–235:37, September 2022. ISSN 0360-0300.
- [20] John H Hubbard and Barbara Burke Hubbard. *Vector calculus, linear algebra, and differential forms: a unified approach*. Matrix Editions, 2015.
- [21] Maurice Jakesch, Advait Bhat, Daniel Buschek, Lior Zalmanson, and Mor Naaman. Co-writing with opinionated language models affects users’ views. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–15, 2023.
- [22] Jinyuan Jia, Ahmed Salem, Michael Backes, Yang Zhang, and Neil Zhenqiang Gong. Memguard: Defending against black-box membership inference attacks via adversarial examples. In *Proceedings of the 2019 ACM SIGSAC conference on computer and communications security*, pages 259–274, 2019.
- [23] Belveze Jules. Tldr news dataset, 2022. URL [https://huggingface.co/datasets/JulesBelveze/tldr\\_news](https://huggingface.co/datasets/JulesBelveze/tldr_news).
- [24] Nikhil Kandpal, Eric Wallace, and Colin Raffel. Deduplicating Training Data Mitigates Privacy Risks in Language Models. In *Proceedings of the 39th International Conference on Machine Learning*, pages 10697–10707. PMLR, June 2022.
- [25] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- [26] Mario Michael Krell, Matej Kosec, Sergio P Perez, and Andrew Fitzgibbon. Efficient sequence packing without cross-contamination: Accelerating large language models without impacting performance. *arXiv preprint arXiv:2107.02027*, 2021.
- [27] Keita Kurita, Paul Michel, and Graham Neubig. Weight poisoning attacks on pretrained models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2793–2806, 2020.
- [28] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, 2021.

- [29] Xuechen Li, Florian Tramer, Percy Liang, and Tatsunori Hashimoto. Large Language Models Can Be Strong Differentially Private Learners. In *International Conference on Learning Representations*, October 2021.
- [30] Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *Advances in Neural Information Processing Systems*, 35:1950–1965, 2022.
- [31] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt understands, too. *AI Open*, 2023.
- [32] Yiyong Liu, Zhengyu Zhao, Michael Backes, and Yang Zhang. Membership Inference Attacks by Exploiting Loss Trajectory. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security, CCS '22*, pages 2085–2098, New York, NY, USA, November 2022. Association for Computing Machinery. ISBN 978-1-4503-9450-5.
- [33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- [34] Justus Mattern, Fatemehsadat Mireshghallah, Zhijing Jin, Bernhard Schölkopf, Mrinmaya Sachan, and Taylor Berg-Kirkpatrick. Membership Inference Attacks against Language Models via Neighbourhood Comparison, May 2023.
- [35] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In *International Conference on Learning Representations*, 2016.
- [36] Bonan Min, Hayley Ross, Elinor Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. Recent Advances in Natural Language Processing via Large Pre-trained Language Models: A Survey. *ACM Computing Surveys*, 56(2): 30:1–30:40, September 2023. ISSN 0360-0300.
- [37] Fatemehsadat Mireshghallah, Huseyin Inan, Marcello Hasegawa, Victor Rühle, Taylor Berg-Kirkpatrick, and Robert Sim. Privacy regularization: Joint privacy-utility optimization in languagemodels. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3799–3807, 2021.
- [38] Fatemehsadat Mireshghallah, Kartik Goyal, Archit Uniyal, Taylor Berg-Kirkpatrick, and Reza Shokri. Quantifying Privacy Risks of Masked Language Models Using Membership Inference Attacks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8332–8347, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- [39] Fatemehsadat Mireshghallah, Archit Uniyal, Tianhao Wang, David Evans, and Taylor Berg-Kirkpatrick. An Empirical Analysis of Memorization in Fine-tuned Autoregressive Language Models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1816–1826, Abu Dhabi, United Arab Emirates, 2022. Association for Computational Linguistics.
- [40] Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. DetectGPT: Zero-Shot Machine-Generated Text Detection using Probability Curvature. In *Proceedings of the 38th International Conference on Machine Learning, {ICML} 2023*, June 2023.
- [41] Shashi Narayan, Shay B Cohen, and Mirella Lapata. Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807, 2018.
- [42] Milad Nasr, Reza Shokri, and Amir Houmansadr. Comprehensive Privacy Analysis of Deep Learning: Passive and Active White-box Inference Attacks against Centralized and Federated Learning. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 739–753, May 2019.

- [43] OpenAI. ChatGPT: Optimizing Language Models for Dialogue. <http://web.archive.org/web/20230109000707/https://openai.com/blog/chatgpt/>, January 2023.
- [44] Charith Peris, Christophe Dupuy, Jimit Majmudar, Rahil Parikh, Sami Smaili, Richard Zemel, and Rahul Gupta. Privacy in the Time of Language Models. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining, WSDM '23*, pages 1291–1292, New York, NY, USA, February 2023. Association for Computing Machinery. ISBN 978-1-4503-9407-9.
- [45] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [46] Samuel Reese, Gemma Boleda, Montse Cuadros, Lluís Padró, and German Rigau. Wikicorpus: A word-sense disambiguated multilingual Wikipedia corpus. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, Valletta, Malta, May 2010. European Language Resources Association (ELRA). URL [http://www.lrec-conf.org/proceedings/lrec2010/pdf/222\\_Paper.pdf](http://www.lrec-conf.org/proceedings/lrec2010/pdf/222_Paper.pdf).
- [47] Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, Yann Ollivier, and Herve Jegou. White-box vs Black-box: Bayes Optimal Strategies for Membership Inference. In *Proceedings of the 36th International Conference on Machine Learning*, pages 5558–5567. PMLR, May 2019.
- [48] Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. ML-Leaks: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models. In *Network and Distributed Systems Security (NDSS) Symposium 2019*, February 2019.
- [49] Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yingsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. In *The Twelfth International Conference on Learning Representations*, 2023.
- [50] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership Inference Attacks Against Machine Learning Models. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 3–18, 2017.
- [51] James Stewart. *Multivariable calculus: concepts and contexts*. Brooks/Cole, 2001.
- [52] Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. Memorization Without Overfitting: Analyzing the Training Dynamics of Large Language Models. *Advances in Neural Information Processing Systems*, 35:38274–38290, December 2022.
- [53] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [54] Priyan Vaithilingam, Tianyi Zhang, and Elena L Glassman. Expectation vs. experience: Evaluating the usability of code generation tools powered by large language models. In *Chi conference on human factors in computing systems extended abstracts*, pages 1–7, 2022.
- [55] Eric Wallace, Tony Zhao, Shi Feng, and Sameer Singh. Concealed data poisoning attacks on nlp models. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 139–150, 2021.
- [56] Ben Wang and Aran Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.
- [57] Lauren Watson, Chuan Guo, Graham Cormode, and Alexandre Sablayrolles. On the Importance of Difficulty Calibration in Membership Inference Attacks. In *International Conference on Learning Representations*, January 2022.

- [58] Yuan Yao, Lorenzo Rosasco, and Andrea Caponnetto. On early stopping in gradient descent learning. *Constructive Approximation*, 26(2):289–315, 2007.
- [59] Jiayuan Ye, Aadyaa Maddi, Sasi Kumar Murakonda, Vincent Bindschaedler, and Reza Shokri. Enhanced Membership Inference Attacks against Machine Learning Models. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security, CCS '22*, pages 3093–3106, New York, NY, USA, November 2022. Association for Computing Machinery. ISBN 978-1-4503-9450-5.
- [60] Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting. In *2018 IEEE 31st computer security foundations symposium (CSF)*, pages 268–282. IEEE, 2018.
- [61] Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A Inan, Gautam Kamath, Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, et al. Differentially private fine-tuning of language models. In *International Conference on Learning Representations*, 2021.
- [62] Minxing Zhang, Zhaochun Ren, Zihan Wang, Pengjie Ren, Zhunmin Chen, Pengfei Hu, and Yang Zhang. Membership Inference Attacks Against Recommender Systems. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security, CCS '21*, pages 864–879, New York, NY, USA, November 2021. Association for Computing Machinery. ISBN 978-1-4503-8454-4.
- [63] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28, 2015.
- [64] 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, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A Survey of Large Language Models, September 2023.

## A Appendix

### A.1 Ethic and Broader Impact Statements

For ethic consideration, this study proposes a membership inference attack algorithm, SPV-MIA, which can be maliciously utilized to infer whether a specific textual entry is fed to the target LLM during the training process. The extensive experiments reveal appreciable privacy leakage of LLMs through SPV-MIA, where the member records are identified with high confidence. We acknowledge that SPV-MIA can bring severe privacy risks to existing LLMs. Therefore, to prevent potential misuse of this research, all experimental findings are based on widely used public datasets. This ensures that every individual textual record we analyze has already been made public, and eliminates any further privacy violations.

For broader impact, we have made our code accessible to the public to allow additional research in the pursuit of identifying appropriate defense solutions. Thus, we posit that our article can inspire forthcoming research to not only focus on the linguistic ability of LLMs, but also take into account the dimensions of public data privacy and security. Besides, our research is more closely aligned with real-world LLM application scenarios, thereby revealing privacy risks in more realistic settings. Our proposed method is scalable, with many aspects left for further exploration and research, such as adapting it for other LLMs.

### A.2 Notations of This Work

Table 4: Notations and descriptions.

Notation	Description
$\mathbf{x}$	A specific data record.
$\tilde{\mathbf{x}}_n^\pm$	A pair of symmetrical paraphrasing text record of the target text record $\mathbf{x}$ .
$D_{mem}$	The training dataset utilized for LLM fine-tuning.
$D_{non}$	A disjoint dataset from the training dataset.
$D_{refer}$	The reference dataset that collected for fine-tuning reference LLM.
$m^{(j)}$	The membership of the data record $\mathbf{x}^{(j)}$ , 1 represents member, whereas 0 represents non-member.
$\theta$	The parameters of the target large language model (LLM).
$\phi$	The parameters of the reference LLM.
$\mathcal{A}(\mathbf{x}, \theta)$	The adversary algorithm for MIA.
$p_\theta(\mathbf{x})$	The probability of text record $\mathbf{x}$ being sampled by the LLM $\theta$ .
$p_\theta(\tilde{\mathbf{x}}_n)$	The probability of paraphrasing text $\tilde{\mathbf{x}}_n$ being sampled by the LLM $\theta$ .
$\Delta p_\theta(\mathbf{x})$	The calibrated probability of text record $\mathbf{x}$ .
$\tilde{p}_\theta(\mathbf{x})$	The probabilistic variation of $\mathbf{x}$ measured on the target LLM $\theta$ .
$\tilde{p}_\phi(\mathbf{x})$	The probabilistic variation of $\mathbf{x}$ measured on the self-prompt reference LLM $\phi$ .
$\Delta \tilde{p}_{\theta, \phi}(\mathbf{x})$	The calibrated probabilistic variation of $\mathbf{x}$ measured on both the target LLM $\theta$ and the self-prompt reference LLM $\phi$ .
$N$	The query times for estimating $\tilde{p}_\theta(\mathbf{x})$ .

### A.3 Detailed Pseudo Codes of Symmetrical Perturbation

---

**Algorithm 1** Symmetrical paraphrase in the embedding domain

---

**Input:** Target text set  $\{\mathbf{x}^{(i)}\}$ , Gaussian noise scale  $\sigma$ , paraphrasing number  $N$ , embedding matrix of tokens  $\mathbf{E}$ .

**Output:** Symmetrical paraphrased text embedding  $emb(\mathbf{x}^{(i)})^\pm$ .

```

1: for  $\mathbf{x}^{(i)} \in \{\mathbf{x}^{(i)}\}$  do
2:    $id(\mathbf{x}^{(i)}) \leftarrow tokenizer(\mathbf{x}^{(i)})$  ▷ Tokenize the text into token ids.
3:    $emb(\mathbf{x}^{(i)}) \leftarrow \mathbf{E}(id(\mathbf{x}^{(i)}))$  ▷ Convert token ids into embeddings.
4:   for  $n \in \{1, \dots, N\}$  do
5:      $\mathbf{z} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$  ▷ Sample noise from Gaussian distribution.
6:      $emb(\mathbf{x}^{(i)})^+ \leftarrow emb(\mathbf{x}^{(i)}) + \mathbf{z}$ 
7:      $emb(\mathbf{x}^{(i)})^- \leftarrow emb(\mathbf{x}^{(i)}) - \mathbf{z}$ 
8:     return  $emb(\mathbf{x}^{(i)})^\pm$ 
9:   end for
10: end for

```

---

**Algorithm 2** Symmetrical paraphrase in the semantic domain

---

**Input:** Target text set  $\{\mathbf{x}^{(i)}\}$ , paraphrasing percentage  $\lambda$ , paraphrasing number  $N$ , embedding matrix of tokens  $\mathbf{E}$ .

**Output:** Symmetrical paraphrased text  $\mathbf{x}^{(i)\pm}$ .

```

1: for  $\mathbf{x}^{(i)} \in \{\mathbf{x}^{(i)}\}$  do
2:    $id(\mathbf{x}^{(i)}) \leftarrow tokenizer.encode(\mathbf{x}^{(i)})$  ▷ Tokenize the text into token ids.
3:   for  $n \in \{1, \dots, N\}$  do
4:     for  $t_j \in id(\mathbf{x}^{(i)}) = [t_0, t_1, \dots, t_{|x|}]$  do
5:       if  $Rand() < \lambda$  then
6:          $t_j \leftarrow [MASK]$  ▷ Mask tokens with the percentage  $\lambda$ .
7:       else
8:          $t_j \leftarrow t_j$ 
9:       end if
10:    end for
11:     $\{t_j^\pm\} \leftarrow MLM(id(\mathbf{x}^{(i)}))$  ▷ Fill the mask tokens with MLM.
12:    for  $t_j^\pm \in \{t_j^\pm\}$  do
13:       $emb(t_j) \leftarrow \mathbf{E}(t_j)$  ▷ Extract the embedding of the original token.
14:       $emb(t_j)^+ \leftarrow \mathbf{E}(t_j^+)$  ▷ Extract the embedding of the paraphrased token.
15:       $\Delta emb(t_j) \leftarrow emb(t_j)^+ - emb(t_j)$  ▷ Measure the paraphrasing noise in the embedding domain.
16:       $emb(t_j)^- \leftarrow emb(t_j) - \Delta emb(t_j)$  ▷ Generate symmetrical embedding.
17:       $t_j^- \leftarrow \mathbf{SearchNearestToken}(emb(t_j)^-, \mathbf{E})$ 
18:    end for
19:     $id(\mathbf{x}^{(i)})^+ \leftarrow \mathbf{FillMaskToken}(\{t_j^+\})$ 
20:     $id(\mathbf{x}^{(i)})^- \leftarrow \mathbf{FillMaskToken}(\{t_j^-\})$ 
21:     $\mathbf{x}^{(i)+} \leftarrow tokenizer.decode(id(\mathbf{x}^{(i)})^+)$ 
22:     $\mathbf{x}^{(i)-} \leftarrow tokenizer.decode(id(\mathbf{x}^{(i)})^-)$ 
23:    return  $\mathbf{x}^{(i)\pm}$ 
24:  end for
25: end for

```

---



Table 5: The MIA performance of SPV-MIA while applied different paraphrasing methods.

Paraphrasing	Embedding	Semantic	Neighbour Comparing
Wiki	0.965	0.951	0.934
AG News	0.926	0.903	0.893
Xsum	0.949	0.937	0.928
Avg.	0.944	0.930	0.918

#### A.4 Reformulate the Neighbour Attack

In this work, we introduce a symmetrical paraphrasing method for assessing probabilistic variation, which is motivated by a rigorous principle: detect the memorization phenomenon rather than overfitting. Meanwhile, we found that the Neighbour comparing [34] has a similar form to our proposed probabilistic variation. Thus, we further consider reformulating neighbour attack based on our intuition, then provide another explanation and motivation for it. As shown in Eq. 10, the assessment of probabilistic variation requires a pair of symmetrical paraphrased text, thus we elaborately design two paraphrasing models on embedding and semantic domains. However, it is still non-trivial to define two neighboring samples with opposite paraphrasing directions for  $\mathbf{x}$ , we therefore consider directly ignoring the requirement for symmetry in the probabilistic variation. Thus, we simplify  $\tilde{\mathbf{x}}_n^\pm$  to be uniformly represented by  $\tilde{\mathbf{x}}_n$ . Then we can reformulate Eq. 10 to Neighbour comparing:

$$\tilde{p}_\theta(\mathbf{x}) \approx \frac{1}{2N} \sum_n \left( p_\theta(\tilde{\mathbf{x}}_n^+) + p_\theta(\tilde{\mathbf{x}}_n^-) \right) - p_\theta(\mathbf{x}) = \frac{1}{2N} \sum_n p_\theta(\tilde{\mathbf{x}}_n) - p_\theta(\mathbf{x}). \quad (11)$$

Therefore, we believe that the neighbour attack and our proposed probabilistic variation can share the same design motivation, namely, detecting special signals that indicate the LLM has memorized training set samples. Additionally, we compared the neighbour attack with our proposed symmetric paraphrasing methods. As shown in Table 5, paraphrasing in the embedding domain achieves considerable performance gains, while paraphrasing in the semantic domain yields a marginal advantage.

#### A.5 Supplementary Experimental Results

##### A.5.1 Ablation Study

Table 6: Results of Ablation Study on GPT-J and LLaMA across three datasets.

Target Model	Wiki		AG News		XSum	
	GPT-J	LLaMA	GPT-J	LLaMA	GPT-J	LLaMA
w/o PDC	0.648	0.653	0.632	0.641	0.653	0.661
w/o PVA	0.901	0.913	0.864	0.885	0.873	0.919
SPV-MIA	0.929	0.951	0.885	0.903	0.897	0.937

In the previous experiments, we have validated the superiority of our proposed SPV-MIA over existing algorithms, as well as its versatility in addressing various challenging scenarios. However, the specific contributions proposed by each module we proposed are still unknown. In this subsection, we conduct an ablation study to audit the performance gain provided by the two proposed modules. Concretely, we respectively remove the practical difficulty calibration (PDC) and probabilistic variation assessment (PVA) that we introduced in Section 4.2 and Section 4.3. The results are represented in Table 6, where each module contributes a certain improvement to our proposed method. Besides, the PVC approach seems to play a more critical role, which can still serve as a valid adversary without the PVA. Thus, in practical scenarios, we can consider removing the PVA to reduce the frequency of accessing public APIs.

##### A.5.2 AUC Curves

As a supplement to the main experimental results represented in Tab. 1, we further provide the raw ROC curve for a more comprehensive presentation in Fig. 6.

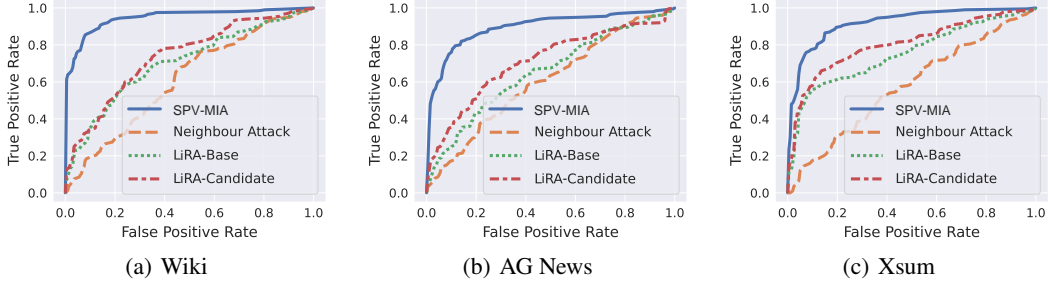


Figure 6: ROC curves of SPV-MIA and the top-three best baselines on LLaMAs fine-tuned over three datasets.

### A.5.3 Performance of Target LLMs

We supplemented the performance of all LLM-dataset pairs on both the training and test sets (estimated using PPL). As shown in Table 7, the experimental results indicate that none of the fine-tuned LLMs exhibit significant overfitting, which aligns with our claim in the main body. Additionally, we provided the performance of the LLM under different privacy budgets  $\epsilon$ , as shown in Table 8.

Table 7: The perplexity (PPL) of each LLM-dataset pair on training set and test set.

Target Model	GPT-2		GPT-J		Falcon		LLaMA	
	Training	Test	Training	Test	Training	Test	Training	Test
Wiki	25.34	27.47	12.35	12.44	7.26	7.73	7.00	7.43
AG News	23.86	26.34	12.49	13.24	8.92	9.54	9.03	9.13
Xsum	26.42	28.31	13.43	13.88	7.69	7.96	7.35	7.65

Table 8: The perplexity (PPL) of each LLM-dataset pair trained w.r.t different privacy budget  $\epsilon$ .

Privacy Budget $\epsilon$	15	30	60	+ inf
Wiki	8.45	8.16	7.76	7.43
AG News	10.84	9.68	9.32	9.13
Xsum	8.89	8.33	7.98	7.65

## A.6 Experimental Settings

In this subsection, we give a extensive introduction of experimental settings, including the datasets, target LLMs and baselines, as well as the implementation details.

### A.6.1 Datasets

Our experiments utilize six different datasets across multiple domains and LLM use cases, where we employ three datasets as the private datasets to fine-tune the target LLMs, and the remaining datasets as the public datasets from the exact domains. Specifically, we use the representative articles on Wikitext-103 dataset [35] to represent academic writing tasks, news topics from the AG News dataset [63] to represent news topic discussion task, and documents from the XSum dataset [41] to represent the article writing task. Besides, we utilize Wikicorpus [46], TLDR News [23], and CNNDM [17] datasets to respectively represent as the publicly accessible dataset from the same domain for each task.

### A.6.2 Target Large Language Models

To obtain a comprehensive evaluation result, we conduct our experiments over four well-known and widely adopted LLMs as the pre-trained models with different scales from 1.5B parameters to 7B parameters:

- **GPT-2 [45]:** It is a transformer-based language model released by OpenAI in 2019, which has 1.5 billion parameters and is capable of generating high-quality text samples.
- **GPT-J [56]:** It is an open-source LLM released by EleutherAI in 2021 as a variant of GPT-3. GPT-j has 6 billion parameters and is designed to generate human-like with appropriate prompts.
- **Falcon-7B [2]:** Falcon is a family of state-of-the-art LLMs created by the Technology Innovation Institute in 2023. Falcon has 40 billion parameters, and Falcon-7B is the smaller version with less consumption.
- **LLaMA-7B [53]:** LLaMA is one of the most state-of-the-art LLM family open-sourced by Meta AI in 2023, which has outperformed other open-source LLMs on various NLP benchmarks. It has 65 billion parameters and has the potential to accomplish advanced tasks, such as code generation. In this work, we utilize the lightweight version, LLaMA-7B.

### A.6.3 Baselines

We choose six MIAs designed for LMs to comprehensively evaluate our proposed method, including three reference-free attacks and one reference-based attack with one variant.

- **Loss Attack [60]:** A standard metric-based MIA that distinguishes member records simply by judging whether their losses are above a preset threshold.
- **Neighbour Attack [34]:** The Neighbour Attack avoids using a reference model to calibrate the loss scores and instead utilizes the average loss of plausible neighbor texts as the benchmark.
- **DetectGPT [40]:** A zero-shot machine-generated text detection method. Although DetectGPT is specially designed for LLMs-generated text detection, but has the potential to be adapted for identifying the text utilized for model training.
- **Likelihood Ratio Attack (LiRA-Base) [39]:** A reference-based attack, which adopts the pre-trained model as the reference model to calibrate the likelihood metric to infer membership.
- **LiRA-Candidate [39]:** A variant version of LiRA, which utilizes a publicly available dataset in the same domain as the training set to fine-tune the reference model.

### A.6.4 Detailed Information for Reproduction

Table 9: Detailed split and other information of datasets.

Dataset	Relative Datasets		Target Model		Reference Model	
	Domain-specific	Irrelevant	# Member	# Non-member	# Member	# Non-member
Wikitext-103	Wikicorpus	AG News	10,000	1,000	10,000	1,000
AG News	TLDR News	Xsum	10,000	1,000	10,000	1,000
Xsum	CNNM	Wikitext-103	10,000	1,000	10,000	1,000

All experiments are compiled and tested on a Linux server (CPU: AMD EPYC-7763, GPU: NVIDIA GeForce RTX 3090), Each set of experiments for the LLM-dataset pairs took approximately 8 hours, and we spent around 14 days completing all the experiments. For each dataset, we pack multiple tokenized sequences into a single input, which can effectively reduce computational consumption without sacrificing performance [26]. Besides, the packing length is set to 128 tokens. Then, we use 10,000 samples for fine-tuning over pre-trained LLMs and 1,000 samples for evaluation. The detailed information of datasets is summarized in Tab. 9. For each target LLM, we let it fine-tuned with the training batch size of 16, and trained for 10 epochs. The learning rate is set to 0.0001. We adopt the AdamW optimizer [33] to achieve the generalization of LLMs, which is composed of the Adam optimizer [25] and the L2 regularization. For GPT-2, which has a relatively small scale, we adopt the full fine-tuning, which means all parameters are trainable. For other LLMs that are larger, we utilize a parameter-efficient fine-tuning method, Low-Rank Adaptation (LoRA) [18], as the default fine-tuning method. For the paraphrasing model in the embedding domain, the Gaussian noise scale is set to  $\sigma = 0.05$ . For the paraphrasing model in the semantic domain, the paraphrasing percentage is set to  $\lambda = 0.2$ . For both of the two paraphrasing models, we generate 10 symmetrical paraphrased text pairs for each target text record. For the reference LLM fine-tuned with our proposed self-prompt approach, we utilize the domain-specific data as the default prompt text source. Then, we collect

10,000 generated texts from target LLMs with an equal length of 128 tokens to construct reference datasets. We fine-tune the reference LLM for 4 epochs and the training batch size of 16.