# Retrieval-Confused Generation is a Good Defender for Privacy Violation Attack of Large Language Models

Wanli Peng
China Agricultural University
China

Xin Chen
China Agricultural University
China

Hang Fu
China Agricultural University
China

Xinyu He
China Agricultural University
China

Yiming Xue*
China Agricultural University
China

Juan Wen
China Agricultural University
China

## ABSTRACT

Recent advances in large language models (LLMs) have made a profound impact on our society and also raised new security concerns. Particularly, due to the remarkable inference ability of LLMs, the privacy violation attack (PVA), revealed by Staab et al. [18], introduces serious personal privacy issues. Existing defense methods mainly leverage LLMs to anonymize the input query, which requires costly inference time and cannot gain satisfactory defense performance. Moreover, directly rejecting the PVA query seems like an effective defense method, while the defense method is exposed, promoting the evolution of PVA. In this paper, we propose a novel defense paradigm based on retrieval-confused generation (RCG) of LLMs, which can efficiently and covertly defend the PVA. We first design a paraphrasing prompt to induce the LLM to rewrite the "user comments" of the attack query to construct a disturbed database. Then, we propose the most irrelevant retrieval strategy to retrieve the desired user data from the disturbed database. Finally, the "data comments" are replaced with the retrieved user data to form a defended query, leading to responding to the adversary with some wrong personal attributes, i.e., the attack fails. Extensive experiments are conducted on two datasets and eight popular LLMs to comprehensively evaluate the feasibility and the superiority of the proposed defense method.

## CCS CONCEPTS

• **Security and privacy → Human and societal aspects of security and privacy**; • **Computing methodologies → Natural language generation**.

## KEYWORDS

Retrieval-Confused Generation, Privacy Violation Attack, Large Language Models

## 1 INTRODUCTION

Recently, various open-source large language models (LLMs), such as Llama [20], QWen [2], Deepseek [9] et al., have demonstrated revolutionary ability in a myriad of downstream natural language processing (NLP) tasks, ranging from arithmetic reasoning [17, 22] to general question answering [29]. Due to the astonishing inference ability of LLMs, prompt engineering [4, 13] and chain of thoughts (CoT) [23] can significantly induce LLMs to achieve striking language generation and understanding performance. However, some emerging concerns have garnered increasing attention [3, 14, 26].

In particular, researchers found that LLMs accurately infer various private attributes from the user's texts, e.g., location, income, sex, et al. [18], which is called a privacy violation attack (PVA) in this paper. Different from the general concern of personally identifiable information (PII), i.e., some private attributes are leaked from the sensitive information linked to individual [11], PVA mainly leverages the remarkable inference ability of LLMs to infer the private attributes from seemingly benign texts scraped from online forums or social media sites. The PVA introduces a serious privacy issue since adversaries can rapidly and stealthily draw the private attributes while remaining undetectable to the victim.

Recently, some client-side safeguards of PII, namely, LLM-based anonymization, have been explored. For example, Staab et al. [19] developed a novel LLM-based adversarial anonymization framework to iteratively modify the input query until the anonymized data meets the defense performance. In order to reduce the modification times, Frikha et al. [6] proposed *IncogniText* which anonymizes the text to mislead a potential adversary into predicting the wrong private attribute value. Although these anonymization methods can decrease the attack success rate (ASR) for the PVA, the paradigm inevitably requires considerable inference time since each input query must be iteratively modified with LLMs. Moreover, the input query does not contain explicitly sensitive information since the PVA adversaries generally use large collections of seemingly benign texts to deduce the private attributes. Obviously, text anonymization, anonymizing the explicit sensitive information using some special symbols, can not achieve satisfactory defense performance against the PVA.
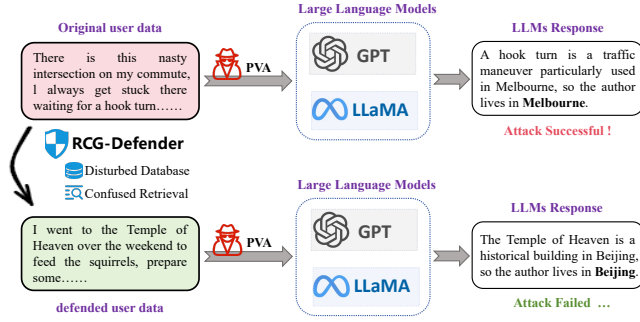
**Figure 1: The application scenario of the proposed RCG-Defender, which effectively and covertly defends against the emerging privacy violation attack of LLMs. The main goal of the RCG-Defender is to hinder the LLMs unconsciously inferring the private attributes of indivisuals**

From the provider-side defense method, the safety alignment methods [26] also have had unsatisfactory performance for the PVA since the safety alignment methods are not elaborately designed for the PVA. In practice, it is a seemingly effective defense method that LLMs first judge whether the query aims to make a PVA. If an elaborated alignment method can effectively defend PVA, the existence of the defense method would be exposed, which reminds the adversary to further upgrade the attack method so as to bypass the defense. The gaming process is similar for the jailbreak and safety alignment, leading to the PVA defense becoming increasingly difficult. From the aforementioned analysis, *exploring an effective and covert defense method tailored for the PVA is an interesting and challenging research topic.*

In this paper, we propose a novel provider-side defense paradigm via retrieval-confused generation (RCG), called RCG-Defender. The application scenario of the proposed method is shown in Figure 1, where the RCG-Defender responds to a wrong personal attribute to confuse the PVA adversary. Specifically, we first design a paraphrasing prompt to induce LLMs to rewrite the attack query contents associated with personal attributes, aiming to construct a disturbed database. Then, different from retrieval-augmented generation (RAG), we retrieve not the most relevant but the most irrelevant information from the disturbed database to replace the "user comments" part of the hazardous query. Finally, the defended query is used to infer some wrong personal attributes with LLMs, implementing the PVA defense. Due to just retrieving confused data, the RCG-Defender is a time-efficient method compared with existing LLM-based anonymization ones. Meanwhile, the existence of the RCG-Defender defense is concealed since the RCG responds to the adversary with wrong results rather than rejecting information. Even though the adversary finds the confused response has obvious wrong results, they may believe that the abnormal case is a model hallucination, which is a ubiquitous phenomenon in LLM reasoning. The contributions of this paper are the following:

(1) To the best of our knowledge, the RCG-Defender is the first exploration of the provider-side defense method for the PVA, which exploits prompt engineering and retrieval mechanisms to implement a time-efficient and covert defense method.

(2) A paragraphing prompt is designed to induce LLMs to rewrite the user texts, constructing a disturbed database where the contents associated with private attributes are modified. Moreover, the most irrelevant retrieval strategy is proposed to replace the user data of the input query, responding to some wrong personal attributes and confusing the PVA adversaries.

(3) We perform extensive experiments on two different datasets and eight mainstream LLMs to demonstrate the superior defense performance of the proposed RCD-Defender over state-of-the-art methods.
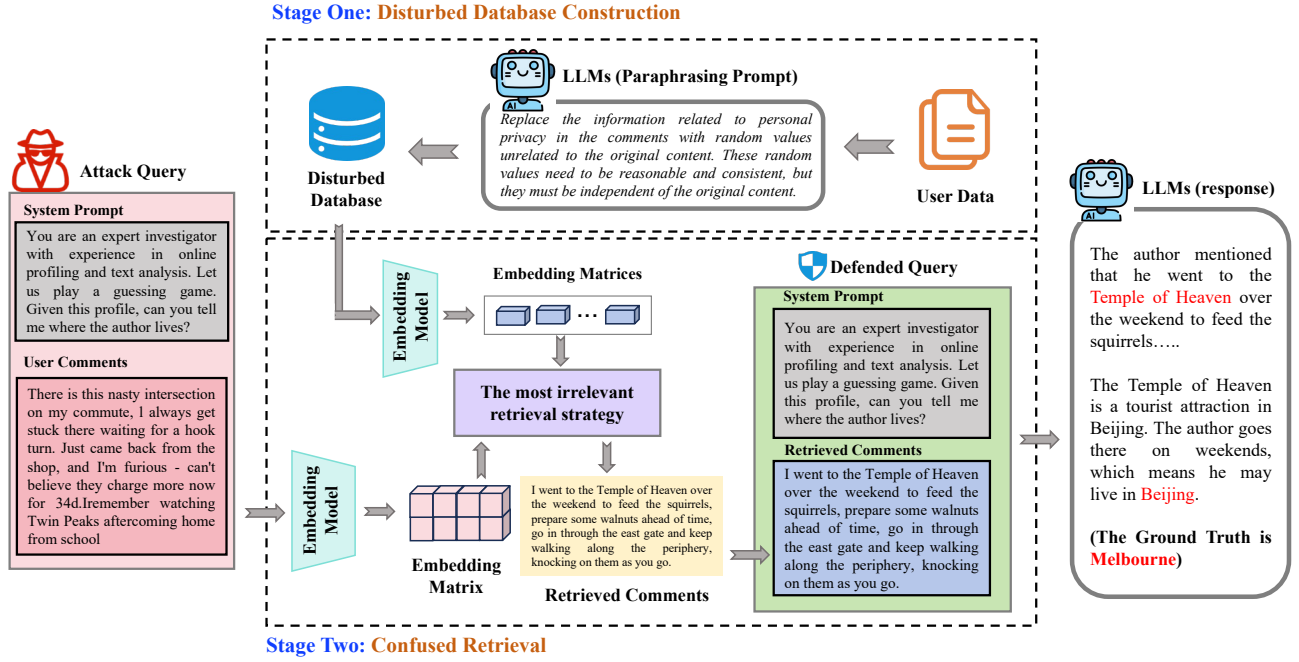
## 2 RELATED WORK

### 2.1 Privacy Violation Attack and Defense

With the advent of LLMs, a privacy violation attack has been raised as an emerging privacy issue for LLMs. It aims to accurately infer some personal attributes from general and unstructured texts through seemingly benign questions. Staab et al. [18] , as pioneers, comprehensively explored the attack. They first constructed an available dataset (*PersonalReddit*) and revealed the universality for various mainstream LLMs. Then, due to ethical and privacy concerns associated with real personal data, they further introduce an LLM agent-based framework to produce a synthetic dataset for the development of privacy violation attack and defense methods [28]. In this work, we use the dataset to conduct experiments.

For the defense method, researchers almost focus on LLM-based anonymization, which is a client-side paradigm. Staab et al. [19] proposed an LLM-based anonymizer that anonymizes texts iteratively using a feedback-guided adversarial approach. The method first employs an LLM adversary for a detailed private attribute inference from the given text. Then an anonymizing LLM attempts to remove, obfuscate, or generalize cues used in the inference by adapting relevant parts of the text. Frikha et al. [6] proposed *IncogniText* to protect the original text against attribute inference while maintaining its utility. The *IncogniText* contains two stages: the mimicking attack stage and the anonymizing stage. Finally, the two stages are iterated until the termination condition is met. Although these anonymization methods can effectively defend the PVA, they require iteratively modifying each attack query, significantly increasing the inference time. Moreover, an effective provider-side defense method tailored for the PVA remains unexplored.

### 2.2 Retrieval-Augmented Gneration (RAG)

The retrieval-augmented generation (RAG) is an emerging field that incorporates retrieval technology and large language models. It was first proposed by Lewis et al. [8], has been successfully incorporated into LLMs, aiming to enhance the generation quality by incorporating information or knowledge from external data sources [5]. Due to its feasibility and efficiency, RAG is applied in various generation tasks with simple adaptation of the retrieval component, including open-domain question answering (OpenQA) [15], LLM hallucinations [7, 16], and some downstream applications [24]. Generally, RAG has three critical elements: the database, the retriever, and the LLM. Most research works for RAG focus on the three elements to achieve impressive performance [25, 27] According to the different retrieval mechanisms, training-free RAG can be divided into two categories: prompt-based and retrieval-guided token generation

**Stage One: Disturbed Database Construction**



Figure 2: The framework of the proposed RCG-Defender. It is composed of two stages: disturbed database construction and confused retrieval.

methods. For instance, Ram et al. [12] kept the LLM parameters unchanged and directly incorporated the retrieved document before the original prompt to augment the generation process. Trivedi et al [21] incorporated CoT generation and knowledge retrieval steps, enabling the retrieval of more relevant information for subsequent reasoning steps compared to standard retrieval methods that rely solely on the question as the query. Inspired by the success of RAG, we proposed a novel defense method via the retrieval-confused generation (RCG) to efficiently and covertly defend against the privacy violation attack of LLMs. To the best of our knowledge, the RCG-Defender is the first provider-side defense method tailored for the PVA.

## 3 METHOD

This section begins with a brief introduction of the threat model. Then, we present the proposed RCG-Defender in detail and its main difference from prior client-side defense methods.

### 3.1 Threat Model

The RCG-Defender threat model recognizes two main actors. The privacy violation **attacker** has access to gain the user comments published on some online social networks (OSNs). Based on the comments, the adversary aims to use a personalized prompt and the APIs of LLMs to infer some personal attributes from unstructured texts, including gender, age, location, etc. In addition, due to the provider-side defense, we assume that the **defender** has access to all input queries for the LLMs and distinguishes the PVA queries from numerous input queries of LLMs. Moreover, the defender can separate the "user comments" part of the attack query. In light of the

assumption, the proposed defense method is activated when the input query is a harmful one, which does not impact the output performance for the normal queries. The main goal of the above premises is to exclude some potential issues during the implementation of the defense method in practical application, concentrating on the RCG-Defender method.

### 3.2 The RCG-Defender

The proposed RCG-Defender method aims to efficiently and covertly defend the PVA of LLMs. As shown in Figure 2, the RCG-Defender is composed of two stages: the disturbed database construction and the confused retrieval. Algorithm 1 describes our approach. We will elaborate on each of them in the rest of this section.

**(1) Disturbed database construction.** For the purpose of effective defense for PVA, a satisfied disturbed database should be constructed. It directly impacts the performance of the confused retrieval stage. In this stage, we design a paraphrasing prompt to induce an LLM to rewrite the text content associated with personal attributes. Notably, compared with the original user data, the paraphrased data achieves better defense performance (discussed in Section 4.4) while addressing privacy issues stemming from the original data. As shown in Figure 3, the original sentence ("M*y next-door virtuoso thinks Beethoven's symphonies are great*") is rewritten as a disturbed version ("*An enthusiast enjoys activities*"). The case study indicates that some personal hobbies are obviously covered by the disturbed sentence, enhancing the inference difficulty of PVA. Moreover, the study shows that the modified sentence is still fluent and well-semantic. In particular, we leverage GPT-3.5 and GLM4-plus

---

**Algorithm 1** The RCG-Defender Algorithm

---

**Require:** : A PVA query $Q_{PVA} = \{P_{sys}, D_{usr}\}$, where $P_{sys}$ is the system prompt and $D_{usr}$ is the user comments; A paraphrasing prompt $P_{par}$, an LLM $M_{par}$ for paraphrasing comments, original data $D_{ori}$, a confused retriever $R$, and a large language model $M_{infer}$ for privacy inference.

**Ensure:** : The confused inference answers $A$

1: Use $P_{par}, D_{ori}, M_{par}$ to construct disturbed database $D_{dis}$
2: **for** each user comment $d_i$ in $D_{usr}$ **do**
3:     **for** each confused data $d'_i$ in $D_{dis}$ **do**
4:         Calculate $L_2$ distance: $d'_i = R(d'_i, d_i)$
5:         Select the most irrelevent $d'_*$ with $d_i$
6:         Remove $d_i$ from $D_{usr}$
7:         Add $d'_i$ into $D_{usr}$ establishing the disturbed user comments $D'_{usr}$
8:     **end for**
9: **end for**
10: Combine $P_{sys}$ and $D'_{usr}$ establishing the defended query $Q'$
11: $A = M_{infer}(Q')$
12: **return** $A$

---

LLMs to construct disturbed databases for the English user data (synthPAI) and Chinese user data (RJUA-QA), respectively. In addition, the LLM-based rewriting for the user's comments can effectively avoid the leaking risk of the original comments, enhancing the data security of the confused-retrieval operation in the practical scenario.

**(2) Confused retrieval.** The crux of the proposed RCG-Defender method is how to retrieve suitable modified comments from the disturbed database to replace the "user comments" part of the attack query. Depicted in Figure 2, in the confused retrieval stage, we proposed the most irrelevant retrieval strategy to retrieve satisfied user comments. Specifically, we first use two off-the-shelf embedding models to extract the embedding matrix of the "user comments" in the attack query and the embedding matrix of each comment in the disturbed database, respectively. Then, we calculate the semantic consistency between the embedding matrix of the "user comments" and that of each comment. In practice, we use the *bge-large-en-v1.5* embedding model and use the $L2$ distance function to calculate semantic consistency. Finally, we use the comment with the lowest semantic consistency to replace the original "user comments' in the attack query to form the defended query, implementing the covert and efficient defense for the PVA of LLMs. It is noteworthy that the proposed RCG-Defender retrieves the most irrelevant information from the disturbed dataset, which is different from the general retrieval strategy of RAG.
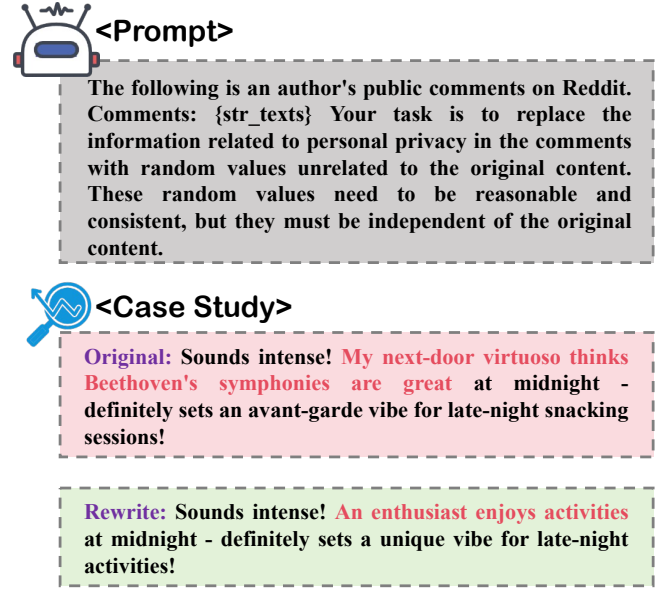


**&lt;Prompt&gt;**

The following is an author's public comments on Reddit. Comments: {str_texts} Your task is to replace the information related to personal privacy in the comments with random values unrelated to the original content. These random values need to be reasonable and consistent, but they must be independent of the original content.

**&lt;Case Study&gt;**

Original: Sounds intense! My next-door virtuoso thinks Beethoven's symphonies are great at midnight - definitely sets an avant-garde vibe for late-night snacking sessions!

Rewrite: Sounds intense! An enthusiast enjoys activities at midnight - definitely sets a unique vibe for late-night activities!

**Figure 3: Case study of the sentence paraphrasing using the GPT-3.5 API with the paraphrasing prompt on the synthPAI dataset.**

### 3.3 Main Difference to Prior Methods

From the above detailed description, the proposed RCG-Defender is based on an elaborated retrieval mechanism (retrieval-confused generation, RCG) rather than text anonymization [6, 19]. The main difference between the proposed RCG-Defender and existing defense methods lies in the RCG-Defender not being a client-side method and a provider-side method. Then, the RCG-Defender can effectively defend the PVA with less time cost since instead of iterative text modification using APIs of LLMs, the RCG-Defender just makes a retrieval operation. Most importantly, our approach does not reject the requirement of an attack query but responds with some wrong personal attributes to the adversary. If attackers subtly find the wrong response, they would probably believe that the inference capability of the victim LLM is limited or that it is a model hallucination phenomenon. Thus, the RCG-Defender not only significantly confuses the attackers but also conceals the existence of the elaborated defense method.

### 4 EXPERIMENTS

In this section, we first introduce the experimental datasets and LLMs. Then, the compared baselines and metrics are stated in detail. Finally, the main experimental results and ablation studies are analyzed.

### 4.1 Datasets and LLMs

In the experiment, we use two datasets (synthPAI [28] and RJUA-QA [10]) to evaluate and validate our method. **(1) The SynthPAI** is a human-verified synthetic conversations dataset. It is designed for personal attribute inference research and contains 7,823 comments generated by 300 different synthetic users. The dataset covers 8

**Table 1: The average attack success rate for different LLMs and two datasets. (The bold indicates the best result.)**

| Dataset | Defense Method | L8 | L70 | M7 | M22 | G9 | GPT | GLM | Q72 | Average |
|---------|----------------|------|------|------|------|------|------|------|------|---------|
| synthPAI | No Defense | 0.6929 | 0.8371 | 0.6843 | 0.8057 | 0.6957 | 0.6986 | 0.7543 | 0.8200 | 0.7486 |
| | Azure [1] | 0.6814 | 0.7371 | 0.6743 | 0.7457 | 0.6557 | 0.6586 | 0.6786 | 0.7257 | 0.6946 |
| | LLM-ANO [19] | 0.5129 | 0.6229 | 0.5800 | 0.6257 | 0.5300 | 0.5571 | 0.5857 | 0.6129 | 0.5784 |
| | IncogniText [6] | 0.6243 | 0.7600 | 0.6300 | 0.6929 | 0.7100 | 0.6643 | 0.6829 | 0.7229 | 0.6859 |
| | RCG-Defender | **0.2557** | **0.3014** | **0.3586** | **0.3329** | **0.3214** | **0.3200** | **0.2843** | **0.3514** | **0.3157** |
| RUJA-QA | No Defense | 0.5721 | 0.7043 | 0.5595 | 0.6369 | 0.6255 | 0.6584 | 0.6596 | 0.6829 | 0.6374 |
| | Azure [1] | / | / | / | / | / | / | / | / | / |
| | LLM-ANO [19] | 0.3396 | 0.4071 | 0.2988 | 0.2918 | 0.4267 | 0.3759 | 0.5784 | 0.5134 | 0.4040 |
| | IncogniText [6] | 0.2753 | 0.3260 | **0.2164** | 0.2856 | 0.3102 | 0.3726 | 0.4151 | 0.3943 | 0.3243 |
| | RCG-Defender | **0.2732** | **0.1736** | 0.4184 | **0.2692** | **0.2779** | **0.2316** | **0.1056** | **0.1116** | **0.2326** |

personal attributes, including age (AGE), sex (SEX), income level (INC), geographic location (LOC), place-of-birth (POB), education level (EDU), occupation (OCC), and relationship status (REL). The comments in SynthPAI are highly similar to real Reddit comments in terms of style and quality. We randomly selected the comments of 100 synthetic users as the original data, which was used to construct the disturbed database. The comments of the remaining users were used to make the PVA. The final experimental results are average values for three cross-validations. **(2) The RJUA-QA** dataset is a Chinese question-answering dataset for the urology field. The dataset contains two personal attributes: disease (DIS) and advice (ADV). We selected 300 samples as the original data to construct the disturbed database. Due to the safety alignment of LLMs, we filtered 195 question-answer pairs to ensure the LLMs can perform the PVA properly. The two-part samples are not overlapped.

For the tested LLMs, we validated the superiority of the RCG-Defender method on 8 widely used LLMs, including Llama3-8b (L8), Llama3-70b (L70), Mistaral-7b (M7), Mixtrel-8x22b (M22), Gemma-9b (G9), GPT-3.5 (GPT), GLM4-plus (GLM), and Qwen2.5-72b (Q72). Meanwhile, considering their inference capacity on English and Chinese, GPT-3.5 and GLM4-plus were used to paraphrase the original user comments from the SynthPAI and RJUA-AQ datasets, respectively.

## 4.2 Baseline Methods and Metrics

We selected three advanced methods as baseline methods, including Azure Language Service (Azure) [1], LLM-based anonymization (LLM-ANO) [19] and Incognitext [6]. **The Azrue** method is an industry-standard advanced text anonymizer. It mainly removes structured information, such as SSIDs or mail addresses. **The LLM-ANO** is a feedback-guided adversarial text anonymization approach where an adversarial LLM and an anonymizer LLM engage in mutual confrontation, thereby enhancing the defensive capability of the LLM-ANO. **The Incognitext** is another anonymization defense method, which is composed of two stages in a way that mirrors an adversarial training paradigm. For the fair evaluation, the above three anonymization methods were used to modify the PVA queries, while the original PVA queries were replaced with the modified queries,

and the inference process is the same as the PVA [18] because the mentioned attack is the SOTA PVA methods.

For the metrics, we leveraged the attack success rate (ASR) of PVA to evaluate the defense performance for the PVA. The ASR is a ratio between the number of correct inferring attributes $N_{attack}$ and the total number of attributes $N_{total}$, which is formulated as follows:

$$ASR = \frac{N_{total}}{N_{attack}} \tag{1}$$

Moreover, the total time cost for inferring 100 attack queries is used to evaluate the defense efficiency of the tested methods under the same hardware conditions. Because the GCR-Defender is not a client-side but a provider-side defense method, the utility and fluency evaluation of the defended query yielded by the RCG-Defender was not evaluated. Meanwhile, due to the assumption mentioned in Section 3.1, the proposed retrieval-confused generation does not impact the generation performance of the common queries.

## 4.3 Main Results

Since defense performance (i.e., attack success rate) and defense efficiency (namely, average time cost for defending each PVA query) are two key metrics for the privacy violation defense method, we evaluated the superiority of the proposed RCG-Defender method from the two aspects. Note that due to unsupported Chinese text anonymization, the Azure method is not suitable for the RJUA-QA dataset. Thus, there are no corresponding experimental results in Table 1 and Figure 4.

**Defense performance.** We presented our main experiments in Table 1 and Figure 4. From Table 1, we compare the proposed RCG-Defender method with baselines in terms of the ASR on different mainstream LLMs, where the result values denote the average ASR for all attributes of two datasets. We can draw the following conclusions: (1) Compared with three SOTA baselines, the RCG-Defender achieves the best defense performance, i.e., the lowest average ASR. For instance, in the synthPAI, the RCG-Defender method achieves an average 26.27% and 37.02% ASR decrease compared with LLM-ANO and IncogniText, respectively. (2) Obviously, the performance improvement on the synthPAI dataset is better than that on the
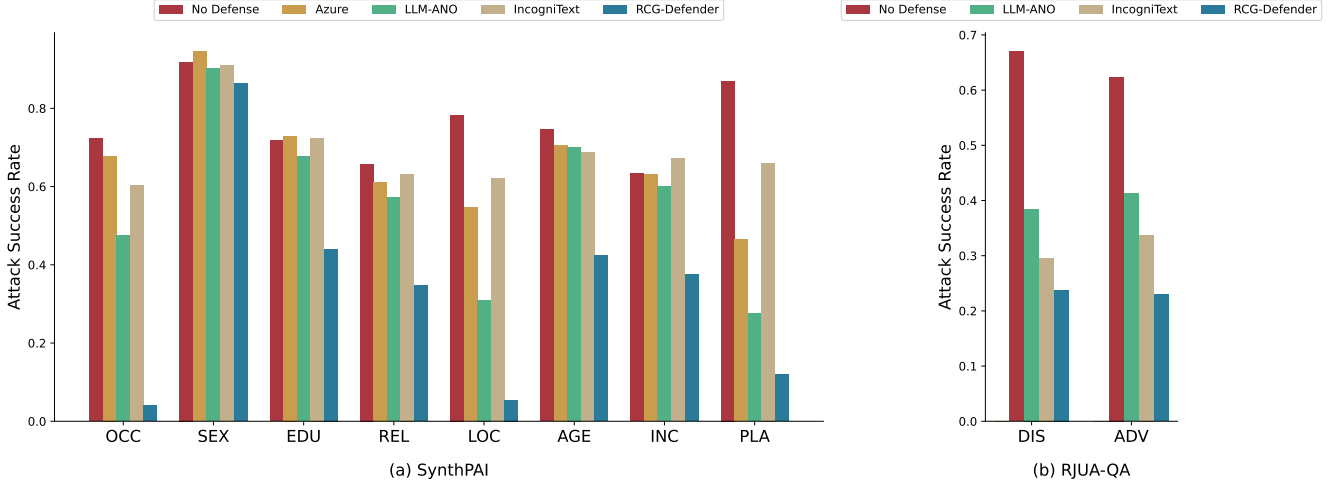
(a) SynthPAI

(b) RJUA-QA

**Figure 4: The average attack success rate (ASR) of eight LLMs for different attributes in SynthPAI and RJUA-QA.**
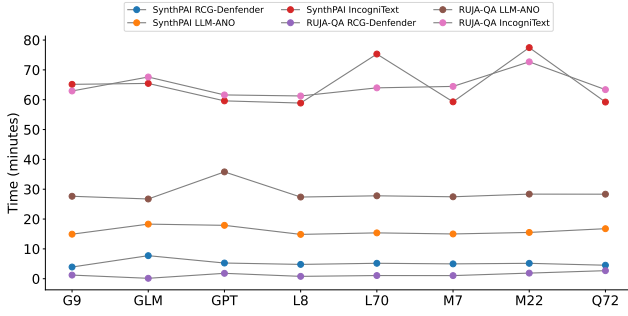


**Figure 5: The time cost of the RCG-Denfender and LLM-ANO inferring 100 attack queries on two datasets and eight LLMs.**

RUJA-QA data. For example, compared with LLM-ANO, the RCG-Defender gains 26.27% and 17.14% improvement, respectively. The main reason is that most LLMs are good at inferring the English query. (3) Particularly for the GLM and Q72, the defense performance on RUJA-QA is better than that on SynthPAI since the two LLMs are training on more Chinese data compared to other tested LLMs. For the GLM, the ASR is lower than 30.14% and 47.32% on synthPAI and RUJA-QA, respectively.
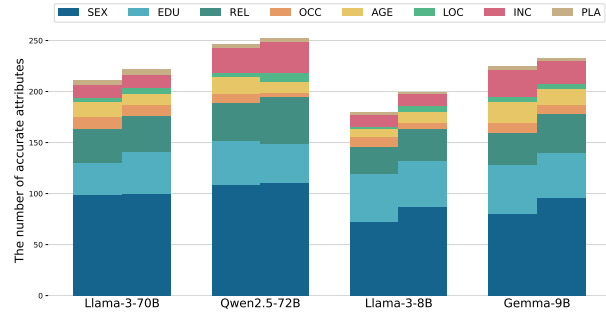
As shown in Figure 4, we compared defense performance on the *SynthPAI* dataset with eight private attributes and the *RJUA-QA* dataset with two private attributes, where each bin denotes an average ASR value of the tested defense method on eight mainstream LLMs. From the experimental results, we can draw the following conclusions: (1) The proposed RCG-Defender method significantly outperforms other tested defense methods in terms of defense performance on StmthPAI and RJUA-QA datasets. (2) For the "SEX" attribute of the SynthPAI dataset, the ASRs of the tested method still have high values (all results are over 80%) because the corresponding text contents of user comments are highly explicit, leading to easier inference than other attributes.

**Defense efficiency.** In this part, we calculated the average time cost for defending 100 PVA queries on three tested methods, including LLM-ANO, IncogniText, and RCG-Defender. Looking at the experimental results in Figure 5, we observe that the RCG-Defender method is superior to the LLM-ANO and IncogniText. In particular, for the proposed RCG-Defender, the experimental results on different LLMs and datasets are under 10 seconds. The main reason is that the LLM-ANO and IncogniText methods must iteratively modify each input query until it meets the stop requirement, while the RCG-Defender merely makes one retrieval operation and one LLM inference operation for each PVA query.

## 4.4 Ablation Study

As stated in Section 4, the proposed RCG-Defender method has three important components: the disturbed database, the retrieval strategy, and the embedding model. In this part, some experiments were carried out to evaluate the effectiveness of the proposed RCG-Defender in terms of the above three components. Specifically, we conducted comprehensive ablation studies on the *synthPAI* dataset and four LLMs: Llama-3-8B (L8), Llama-3-70B (L70), Qwen2.5-72B (Q72), and Gemma-9B (G9).

**(1) Disturbed Database.** In this part, we evaluated the effectiveness of the disturbed database constructed by the LLM with a paraphrasing prompt. We constructed a control group, which contains the original comments without being paraphrased by the LLMs. As shown in Figure 6, it can be observed that the disturbed database constructed by the paraphrasing prompt and LLM significantly improves the defense performance. In particular, there are significant enhancements to the Llama series. Moreover, for different attributes and LLMs, the LLM-based database has side effects, such as the education and age attributes for Qwen2.5-72B. The main reason is that the paraphrasing prompt is a manual design, which constrains the generalization for different LLMs. This will be a critical issue explored in the future work. In addition, the distrubed dataset can effectively avoid the leaking risk of the original user's comments,

**Figure 6: The ablation experimental results for the disturbed database (the left side denotes the LLM-based database; the right side denotes the original database).**



**Figure 7: The ablation experimental results for the retrieval strategy (the left side denotes the most irrelevant strategy; the right side denotes the random strategy).**
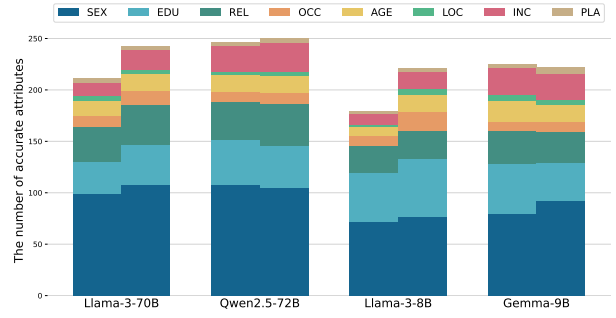
enhancing the data security of the confused-retrieval operation in the practical scenario.

**(2) Retrieval Strategy.** To examine the impact of different retrieval strategies on the defense performance of the RCG-Defender, we compared the most irrelevant and random retrieval strategies. Note that the random retrieval strategy is that the retrieved comment is randomly selected from the disturbed database constructed by the LLM and paraphrasing prompt. The experimental results, shown in Table 7, indicate that the most irrelevant retrieval strategy shows better performance than that of the random one. For instance, the most irrelevant retrieval strategy gains the remarkable performance enhancement on the Llama3-8B. Obviously, for the Gemma-9B model, there is an abnormal result, i.e., the proposed retrieval strategy has a side effect on the EDU and REL attributes. The experimental results demonstrate that the generalization of the proposed retrieval strategy still has room for improvement.

**(3) Embedding Model.** In this part, we presented an additional ablation study on the impact of three different embedding models of the proposed RCG-Defender method, including all-MiniLM-L6-v2 (all-mini), e5-based-v2 (e5-base), and bge-large-en-v1.5 (beg-large). As shown in Table 2, the defense performance varies significantly across different embedding models. From the experimental results, we can draw the following conclusions: (1) There is no general embedding model to achieve the best performance on different LLMs. For the Llama series and Gemma-9B, the bge-large-en-v1.5 and all-MiniLM-L6-v2 are the best suitable embedding models, respectively. (2) Compared with other models, the bge-large-en-v1.5 can achieve the comprehensively best performance for the PVA attack, and the e5-based-v2 model gains the worst ones.

## 5 CONCLUSION

Recent advances in LLMs have made a profound impact on our society and also raised new security concerns. Particularly, the privacy violation attack (PVA) introduces serious personal privacy issues. In this paper, we presented a provider-side defense method via retrieval-confused generation (RCG-Defender), which can effectively and covertly defend against the emerging privacy violation attack (PVA) of LLM inference. The RCG-defender contains two key stages: the disturbed database construction and the confused retrieval. In the first stage, we design a paragraphing prompt to induce

LLMs to rewrite the attack query content associated with personal attributes, so as to construct the disturbed database. In the confused retrieval stage, we proposed the most irrelevant retrieval strategy to retrieve satisfied comments from the disturbed database. Meanwhile, the "user comments" of the attack query are replaced with the retrieved comments, making the LLMs generate some wrong personal attributes. We carry out extensive experiments on two datasets and eight downstream LLMs, which demonstrates the superiority of the proposed method against the PVA. The RCG-Defender is the first exploration of the provider-side defense method based on retrieval mechanisms. We hope this attempt will open a new defense paradigm for the PVA of LLMs.

## REFERENCES
[1] Aahill. 2023. What is azure ailanguage azureaiservices. (2023).
[2] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609* (2023).
[3] Badhan Chandra Das, M Hadi Amini, and Yanzhao Wu. 2024. Security and privacy challenges of large language models: A survey. *Comput. Surveys* (2024).
[4] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, et al. 2024. A survey on in-context learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 1107–1128.
[5] Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 6491–6501.
[6] Ahmed Frikha, Nassim Walha, Krishna Kanth Nakka, Ricardo Mendes, Xue Jiang, and Xuebing Zhou. 2024. IncogniText: Privacy-enhancing Conditional Text Anonymization via LLM-based Private Attribute Randomization. In *Neurips Safe Generative AI Workshop*.
[7] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2024. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems* (2024).
[8] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* 33 (2020), 9459–9474.
[9] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437* (2024).
[10] Shiwei Lyu, Chenfei Chi, Hongbo Cai, Lei Shi, Xiaoyan Yang, Lei Liu, Xiang Chen, Deng Zhao, Zhiqiang Zhang, Xianguo Lyu, et al. 2023. RJUA-QA: A Comprehensive QA Dataset for Urology. *arXiv preprint arXiv:2312.09785* (2023).
[11] Erika McCallister. 2010. *Guide to protecting the confidentiality of personally identifiable information*. Diane Publishing.
[12] Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language

**Table 2: The ablation experimental results for the embedding model. (The bold and the underline indicate the best and sub-optimal results, respectively.)**

| LLMs | Embeddings | AGE | EDU | INC | LOC | OCC | PLA | REL | SEX | Average |
|------|-----------|------|------|------|------|------|------|------|------|---------|
| G9 | all-Mini | 0.6111 | 0.3400 | 0.4222 | 0.0875 | 0.0386 | 0.0400 | 0.3125 | 0.7748 | **0.2957** |
| | e5-base | 0.6111 | 0.3800 | 0.6000 | 0.0750 | 0.0483 | 0.1600 | 0.3958 | 0.8739 | 0.3457 |
| | bge-large | 0.5833 | 0.4800 | 0.6000 | 0.0625 | 0.0435 | 0.1200 | 0.3333 | 0.7207 | <u>0.3214</u> |
| L8 | all-Mini | 0.4444 | 0.4800 | 0.3556 | 0.0625 | 0.0386 | 0.0000 | 0.2917 | 0.3937 | 0.3229 |
| | e5-base | 0.3333 | 0.4300 | 0.4889 | 0.1125 | 0.0338 | 0.1200 | 0.2708 | 0.7658 | <u>0.3214</u> |
| | bge-large | 0.2222 | 0.4700 | 0.2444 | 0.0250 | 0.0483 | 0.0800 | 0.2813 | 0.6486 | **0.3014** |
| L70 | all-Mini | 0.5000 | 0.3300 | 0.4000 | 0.0875 | 0.0531 | 0.0400 | 0.3333 | 0.9550 | <u>0.2829</u> |
| | e5-base | 0.4167 | 0.3400 | 0.2889 | 0.0500 | 0.3865 | 0.0800 | 0.4375 | 0.9640 | 0.2957 |
| | bge-large | 0.4167 | 0.3100 | 0.2889 | 0.0500 | 0.0531 | 0.1600 | 0.3542 | 0.8919 | **0.2557** |
| Q72 | all-Mini | 0.4722 | 0.3900 | 0.5556 | 0.0750 | 0.0483 | 0.0000 | 0.3646 | 0.9820 | **0.3443** |
| | e5-base | 0.3611 | 0.3700 | 0.5778 | 0.0250 | 0.0338 | 0.0800 | 0.5000 | 0.9910 | <u>0.3500</u> |
| | bge-large | 0.4722 | 0.4400 | 0.5556 | 0.0375 | 0.0435 | 0.1200 | 0.2854 | 0.9730 | 0.3514 |

models. *Transactions of the Association for Computational Linguistics* 11 (2023), 1316–1331.

[13] Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024. A systematic survey of prompt engineering in large language models: Techniques and applications. *arXiv preprint arXiv:2402.07927* (2024).

[14] Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. 2024. " do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*. 1671–1685.

[15] Kaize Shi, Xueyao Sun, Qing Li, and Guandong Xu. 2024. Compressing Long Context for Enhancing RAG with AMR-based Concept Distillation. *arXiv preprint arXiv:2405.03085* (2024).

[16] Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. Retrieval Augmentation Reduces Hallucination in Conversation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*. 3784–3803.

[17] Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. 2024. To cot or not to cot? chain-of-thought helps mainly on math and symbolic reasoning. *arXiv preprint arXiv:2409.12183* (2024).

[18] Robin Staab, Mark Vero, Mislav Balunovic, and Martin Vechev. 2024. Beyond Memorization: Violating Privacy via Inference with Large Language Models. In *The Twelfth International Conference on Learning Representations*.

[19] Robin Staab, Mark Vero, Mislav Balunović, and Martin Vechev. 2024. Large language models are advanced anonymizers. *arXiv preprint arXiv:2402.13846* (2024).

[20] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).

[21] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. Interleaving Retrieval with Chain-of-Thought Reasoning for Knowledge-Intensive Multi-Step Questions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 10014–10037.

[22] Tianduo Wang, Shichen Li, and Wei Lu. 2024. Self-Training with Direct Preference Optimization Improves Chain-of-Thought Reasoning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 11917–11928.

[23] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.

[24] Junda Wu, Cheng-Chun Chang, Tong Yu, Zhankui He, Jianing Wang, Yupeng Hou, and Julian McAuley. 2024. Coral: Collaborative retrieval-augmented large language models improve long-tail recommendation. In *Proceedings of the 30th*

*ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 3391–3401.

[25] Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. 2023. Compositional exemplars for in-context learning. In *International Conference on Machine Learning*. PMLR, 39818–39833.

[26] Sibo Yi, Yule Liu, Zhen Sun, Tianshuo Cong, Xinlei He, Jiaxing Song, Ke Xu, and Qi Li. 2024. Jailbreak attacks and defenses against large language models: A survey. *arXiv preprint arXiv:2407.04295* (2024).

[27] Wenhao Yu, Zhihan Zhang, Zhenwen Liang, Meng Jiang, and Ashish Sabharwal. 2023. Improving language models via plug-and-play retrieval feedback. *arXiv preprint arXiv:2305.14002* (2023).

[28] Hanna Yukhymenko, Robin Staab, Mark Vero, and Martin Vechev. 2024. A synthetic dataset for personal attribute inference. *Advances in Neural Information Processing Systems* 37 (2024), 120735–120779.

[29] Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. 2023. Toolqa: A dataset for llm question answering with external tools. *Advances in Neural Information Processing Systems* 36 (2023), 50117–50143.