

# Privacy Implications of Retrieval-Based Language Models

Yangsibo Huang Samyak Gupta Zexuan Zhong Kai Li Danqi Chen

Princeton University

yangsibo@princeton.edu {samyakg, zzhong, li, danqic@cs.princeton.edu}

## Abstract

Retrieval-based language models (LMs) have demonstrated improved interpretability, factuality, and adaptability compared to their parametric counterparts, by incorporating retrieved text from external datastores. While it is well known that parametric models are prone to leaking private data, it remains unclear how the addition of a retrieval datastore impacts model privacy. In this work, we present the first study of privacy risks in retrieval-based LMs, particularly  $k$ NN-LMs. Our goal is to explore the optimal design and training procedure in domains where privacy is of concern, aiming to strike a balance between utility and privacy. Crucially, we find that  $k$ NN-LMs are more susceptible to leaking private information from their private datastore than parametric models. We further explore mitigations of privacy risks. When privacy information is targeted and readily detected in the text, we find that a simple sanitization step would completely eliminate the risks, while decoupling query and key encoders achieves an even better utility-privacy trade-off. Otherwise, we consider strategies of mixing public and private data in both datastore and encoder training. While these methods offer modest improvements, they leave considerable room for future work. Together, our findings provide insights for practitioners to better understand and mitigate privacy risks in retrieval-based LMs.<sup>1</sup>

## 1 Introduction

Retrieval-based language models (Khandelwal et al., 2020; Borgeaud et al., 2022; Izacard et al., 2022; Zhong et al., 2022; Min et al., 2022) generate text distributions by referencing both the parameters of the underlying language model and the information retrieved from a datastore of text. Specifically, the retrieval process involves accessing a

<sup>1</sup>Our code is available at [https://github.com/Princeton-SysML/kNNLM\\_privacy](https://github.com/Princeton-SysML/kNNLM_privacy).

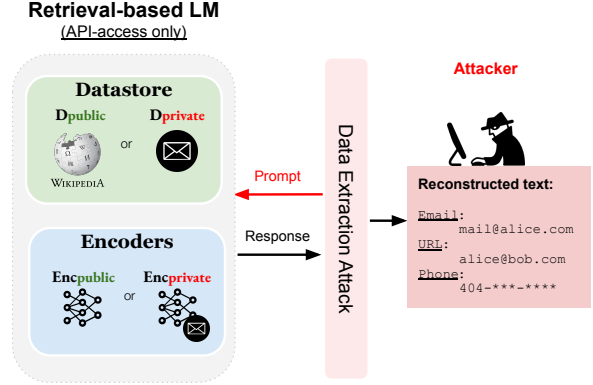


Figure 1: Retrieval-based language models (e.g.,  $k$ NN-LMs) comprise encoders and the datastore as their key components. When tailoring a retrieval-based language model for a privacy-sensitive task, both components may utilize private data. However, a malicious user possessing API access to the model can exploit a data extraction attack in order to reconstruct sensitive information. This study explores the severity of such attacks and suggests novel strategies to mitigate them.

pre-defined datastore to retrieve a set of tokens or text passages that are most relevant to the prompt provided to the model. These retrieved results are then utilized as additional information when generating the model’s response to the prompt. Retrieval-based language models offer promising prospects in terms of enhancing interpretability, scalability, factual accuracy, and adaptability.

However, in privacy-sensitive applications, *utility usually comes at the cost of privacy leakage*. Recent work has shown that large language models are prone to memorizing (Thakkar et al., 2021; Zhang et al., 2021) specific training datapoints, such as personally identifying or otherwise sensitive information. These sensitive datapoints can subsequently be extracted from a trained model through a variety of techniques (Carlini et al., 2019, 2021; Lehman et al., 2021), colloquially known as *data extraction attacks*. While the threat of training data memorization on model privacy has been

studied for parametric language models, there is a lack of evidence regarding the privacy implications of retrieval-based language models, especially how the use of external datastore would impact privacy.

In this work, we present the first study of privacy risks in retrieval-based language models, with a focus on the nearest neighbor language models ( $k$ NN-LMs) (Khandelwal et al., 2020), which have been extensively studied in the literature (He et al., 2021; Zhong et al., 2022; Shi et al., 2022a; Xu et al., 2023)<sup>2</sup>. In particular, we are interested in understanding  $k$ NN-LMs’ privacy risks in real-world scenarios where they are deployed via an API to the general public (Figure 1). We consider a scenario in which a model creator has a private, domain-specific datastore that improves model performance on domain-specific tasks, but may also contain sensitive information that should not be revealed. In such a scenario, the model creator must find a balance between utilizing their private dataset to enhance model performance and protecting sensitive information. Note that the development of  $k$ NN-LMs intended solely for *personal use* (e.g., constructing a  $k$ NN-LM email autocompleter by combining a public LM with a private email datastore) falls outside the scope of our study because it does not involve any attack channels that could be exploited by potential attackers.

We begin our investigation by examining a situation where the creator of the model only adds private data to the retrieval datastore during inference, as suggested by Borgeaud et al. (2022). Our findings indicate that while this approach enhances utility, it introduces an elevated privacy risk to the private data compared to parametric language models (Section 3), and adversarial users could violate the confidentiality of the datastore by recovering sensitive datapoints. Therefore, it is vital for the model creator to refrain from storing sensitive information in the datastore.

We further explore mitigation strategies for  $k$ NN-LMs in two different scenarios. The first is where private information is *targeted*, i.e., can be easily identified and removed (Section 4). We explore enhancing the privacy of  $k$ NN-LMs by eliminating privacy-sensitive text segments from both the

datastore and the encoder’s training process. This approach effectively eliminates the targeted privacy risks while resulting in minimal loss of utility. We then explore a finer level of control over private information by employing distinct encoders for keys (i.e., texts stored in the datastore) and queries (i.e., prompts to the language model). Through our experimental analysis, we demonstrate that this design approach offers increased flexibility in striking a balance between privacy and model performance.

The second is a more challenging scenario where the private information is *untargeted*, making it impractical to remove from the data (Section 5). To address this issue, we explore the possibility of constructing the datastore using public datapoints. We also consider training the encoder of the  $k$ NN-LM model using a combination of public and private datapoints to minimize the distribution differences between the public data stored in the datastore and the private data used during inference. Despite the modest improvements from the methods we explored, the mitigation of untargeted attacks remains challenging and there is considerable room for future work. We hope our findings provide insights for practitioners to better understand and mitigate privacy risks in retrieval-based LMs.

## 2 Background and Problem Formulation

In this section, we will review the key components of  $k$ NN-LMs (Section 2.1) and then discuss the risks of data extraction in parametric language models (Section 2.2). These fundamental aspects lay a foundation for the subsequent exploration and analysis of privacy risks related to  $k$ NN-LMs. We will also formally describe our problem setup (Section 2.3).

### 2.1 Nearest Neighbor Language Models

A  $k$ NN-LM (Khandelwal et al., 2020) augments the standard language model with a *datastore* from which it can retrieve tokens to improve performance. The tokens in the datastore are pre-computed using an *encoder*. We use the term "query" to denote the prompt provided to a  $k$ NN-LM and it is encoded by the encoder  $\text{Enc}_Q$ . The term "key" is used to denote the tokens in the datastore and it is encoded by the encoder  $\text{Enc}_K$ .

**Encoders** Given a vocabulary  $\mathcal{V}$ , The encoder  $\text{Enc}(\cdot)$  maps a context  $c \in \mathcal{V}^*$  to a fixed-length vector representation. Typically,  $\text{Enc}(\cdot)$  is computed

<sup>2</sup>Other retrieval-based language models such as RETRO (Borgeaud et al., 2022) and Atlas (Izcard et al., 2022) have significantly different architectures and the findings from our investigation may not necessarily apply to these models. Investigation of these models will be left as future work.

using a trained language model, where  $\text{Enc}(c)$  represents the vector hidden representation obtained from the output layer of the language model when provided with the input  $c$ . Although we use different subscripts to denote the encoder for keys ( $\text{Enc}_K$ ) and for queries ( $\text{Enc}_Q$ ),  $k$ NN-LMs usually use the same encoders for keys and queries.

**Datastore** The datastore is a key-value store generated by running the encoder  $\text{Enc}_K(\cdot)$  over a corpus of text. Each key is the vector representation  $\text{Enc}_K(c)$  for some context  $c \in \mathcal{V}^*$ , and each value  $v \in \mathcal{V}$  is the ground-truth next word for the context  $c$ . A search index is then constructed based on the key-value store to enable retrieval.

**Inference** At inference time, when predicting the next token for a query  $x \in \mathcal{V}^*$ , the model queries the datastore with encoded query  $\text{Enc}_Q(x)$  to retrieve  $x$ 's  $k$ -nearest neighbors  $\mathcal{N}_k$  according to a distance function  $d(\cdot, \cdot)$ <sup>3</sup>. Then the model computes a softmax over the (negative) distances, which gives  $p_{k\text{NN}}(y|x)$ , a distribution over the next token, in proportional to:

$$\sum_{(c_i, w_i) \in \mathcal{N}_k} \mathbf{1}_{y=w_i} \exp \left( -\frac{d(\text{Enc}_K(c_i), \text{Enc}_Q(x))}{t} \right),$$

where  $t$  is a temperature parameter, and  $k$  is a hyper-parameter that controls the number of retrieved neighbors. The prediction is then interpolated with the prediction from the original LM:  $p(y|x) = \lambda p_{k\text{NN}}(y|x) + (1 - \lambda)p_{\text{LM}}(y|x)$ , where  $\lambda$  is an interpolation coefficient.

## 2.2 Data Extraction Attacks

Prior work (Carlini et al., 2021) demonstrates that an attacker can extract private datapoints from the training set of a learned language model. The existence of such an attack poses a clear and alarming threat to the confidentiality of sensitive training data, potentially jeopardizing deployment in real-world scenarios (e.g., Gmail's autocomplete model (Chen et al., 2019), which is trained on private user emails).

The attack consists of two main steps: 1) generating candidate reconstructions by prompting the trained models, and 2) sorting the generated candidates based on a score that indicates the likelihood of being a memorized text. Further details about the attack can be found in Appendix A.

<sup>3</sup> $d(\cdot, \cdot)$  is usually the squared  $\ell_2$  distance.

While previous research has successfully highlighted the risks associated with data extraction in parametric language models, there remains a notable gap in our understanding of the risks (and any potential benefits) pertaining to retrieval-based language models like  $k$ NN-LMs. This study aims to address this gap and provide insights into the subject matter.

## 2.3 Problem Formulation

We consider a scenario where a service provider (e.g. a financial institution) aims to enhance its customer experience by developing a  $k$ NN-LM and deploying it as an API service. We assume that the service provider possesses its own private data ( $\mathcal{D}_{\text{private}}$ ) specific to its domain, in addition to publicly available data ( $\mathcal{D}_{\text{public}}$ ).

We identify two key design choices which impact the quality and privacy of such a deployed service. First, the service provider chooses which data to be included in its datastore, and this may be public data ( $\mathcal{D}_{\text{public}}$ ), private data ( $\mathcal{D}_{\text{private}}$ ), or a mix of both. Second, they choose whether to use encoders that are pre-trained on publicly available data ( $\text{Enc}_{\text{public}}$ ), or further finetuned on the private data ( $\text{Enc}_{\text{private}}$ ). We posit that careful consideration of these design choices is needed to establish a balance between privacy preservation and utility.

The service provider in such a scenario is concerned with making a useful API, while keeping their private data hidden from malicious users or attackers. Hence, the service provider's objective is to attain a high level of utility (as measured by perplexity) on a held-out set of  $\mathcal{D}_{\text{private}}$  while simultaneously minimizing the disclosure of private information. We quantify the metrics we consider for privacy in Section 3.1.

## 3 Privacy-Utility of $k$ NN-LMs with a Private Datastore

This section presents our investigation of whether the addition of private data to the retrieval datastore during inference is an effective method for achieving a good trade-off between privacy (measured by metrics defined in Section 3.1) and utility (measured by perplexity) in  $k$ NN-LMs.

### 3.1 Privacy Measurements

We first describe how we evaluate the risk of data extraction attack within the scenario described earlier in 2.3.

**Threat model** We assume that the service provider deploys a  $k$ NN-LM with APIs<sup>4</sup> for:

1. Perplexity calculation of some text;
2. Text completion given some contexts.

Our study considers two types of privacy risks, each associated with a particular type of attack:

**Targeted attacks** We define targeted risk as a privacy risk that can be directly associated with a segment of text (e.g., personal identifiers such as addresses and telephone numbers.). A targeted attacker’s goal is to extract that certain segment of text. In our study, we focus on the extraction of Personal Identifiable Information (PII), including email addresses, telephone numbers, and URLs. To tailor the extraction attack to recover text segments such as PIIs rather than the entire training text, we customize the attack prompts based on the type of information to be extracted. Specifically, we gather common preceding context for telephone numbers, email addresses, and URLs, and use them as prompts. Appendix B provides example prompts we use in the attack. For evaluation, we measure how many private PIIs of each category have been successfully reconstructed by the attacker:

- We firstly detect all unique personal identifiers in the private dataset, denoted as  $\{\rho_i\}_{i=1}^p$ ;
- We then sort the reconstruction candidates based on the membership metrics defined in Appendix A, and only keep the top- $n$  candidates  $\{c_i\}_{i=1}^n$ ;
- Finally, we detect  $\{\hat{\rho}_i\}_{i=1}^q$ , the unique PIIs in the top- $n$  candidates, and then count  $|\{\rho_i\}_{i=1}^p \cap \{\hat{\rho}_i\}_{i=1}^q|$ , namely how many original PIIs have been successfully reconstructed by the attack. A larger number means higher leakage of private PIIs.

**Untargeted attacks** The untargeted attack is the case where the attacker aims to recover the entire training example, rather than a specific segment of text. Such attacks can potentially lead to the theft of valuable private training data. To perform the untargeted attack, we adopt the attack proposed by Carlini et al. (2021) as the untargeted attack, which is described in detail in the appendix. For evaluation, we measure the similarity between the reconstructed text and the original private text:

<sup>4</sup>Note that the attacker cannot access the internal parameters of the deployed model.

- We firstly sort the reconstruction candidates based on the membership metrics defined in Appendix A, and only keep the top- $n$  candidates  $\{c_i\}_{i=1}^n$ ;
- For each candidate  $c_i$ , we then find the closest example in the private dataset  $p_i$  and compute the ROUGE-L score between  $c_i$  and  $p_i$ . If the score is higher than 0.5, we mark the candidate as a good reconstruction.

### 3.2 Evaluation Setup

Our evaluation treats the Enron Email dataset (Klimt and Yang, 2004), which contains around 500,000 emails generated by employees of the Enron Corporation, as the private dataset  $\mathcal{D}_{\text{private}}$ . We use the WikiText-103 dataset (Merity et al., 2016) as  $\mathcal{D}_{\text{public}}$ . We pre-process the Enron Email dataset by retaining only the email body. We then use regular expressions to identify and extract three types of personal identifiers for the use of the targeted attack: telephone numbers, email addresses, and URLs. The statistics for these personal identifiers can be found in Appendix B. We use the GPT-2 base model (Radford et al., 2019) as  $\text{Enc}_{\text{public}}$ , and finetune it on the Enron Email dataset as  $\text{Enc}_{\text{private}}$ .

During inference, the model retrieves  $k$  nearest neighbors according to the squared  $\ell_2$  distance, normalizes their distribution over the next word with a softmax using a temperature value of 1, and then uses an interpolation factor of  $\lambda$  to combine  $p_{\text{kNN}}$  and  $p_{\text{LM}}$ . For each model configuration, we search  $k \in \{64, 128, 256, 512, 1024, 2048\}$  and  $\lambda \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$  on a held-out validation set for the best model perplexity, and then run the inference on the evaluation set.

Privacy measurements are computed using the top 1000 candidates for the targeted attack and using the top 5000 candidates for the untargeted attack.

We are particularly interested in three scenarios: utilizing only  $\text{Enc}_{\text{public}}$  (the publicly pre-trained language model), utilizing only  $\text{Enc}_{\text{private}}$  (the model fine-tuned from  $\text{Enc}_{\text{public}}$  using private data), and utilizing  $\text{Enc}_{\text{public}}$  with  $\mathcal{D}_{\text{private}}$  (the combination of the public model with the private datastore). As shown in Table 1, using  $\text{Enc}_{\text{public}}$  alone results in very poor utility performance but poses minimal risk of data extraction from the private domain, as it has not been exposed to private datapoints. Using  $\text{Enc}_{\text{private}}$  enhances utility



MODEL	EVAL. PPL	TARGETED ATTACK				UNTARGETED ATTACK # GOOD RECON
		TOTAL	PHONE	EMAIL	URL	
(PARAMETRIC LM) $\text{Enc}_{\text{public}}$	30.28	0	0	0	0	0
(PARAMETRIC LM) $\text{Enc}_{\text{private}}$	20.63	28	11	14	3	620
( $k$ NN-LM) $\text{Enc}_{\text{public}}$ w/ $\mathcal{D}_{\text{private}}$	18.41	35	11	16	8	591
( $k$ NN-LM) $\text{Enc}_{\text{private}}$ w/ $\mathcal{D}_{\text{private}}$	16.12	54	25	23	6	656

Table 1: Perplexity and data extraction risks for various model configurations. The configuration with the highest leakage is emphasized in red; the configuration with the lowest leakage is highlighted in green. Privacy measurements are computed using the top 1000 candidates for the targeted attack and using the top 5000 candidates for the untargeted attack.

(perplexity improves from 30.28 to 20.63) but increases the risk of data extraction.

When it comes to  $k$ NN-LMs, incorporating a private datastore ( $\mathcal{D}_{\text{private}}$ ) with a public model ( $\text{Enc}_{\text{public}}$ ) yields even greater utility compared to relying solely on the fine-tuned model ( $\text{Enc}_{\text{private}}$ ). However, this utility improvement also comes at the expense of increased privacy leakage. These findings suggest that the privacy concern stemming from the private datastore outweighs that resulting from the privately fine-tuned model, indicating a lack of robust privacy protection in the design of  $k$ NN-LMs. Additionally, we note that the combination of  $\text{Enc}_{\text{private}}$  and  $\mathcal{D}_{\text{private}}$  achieves the highest utility but also incurs the highest privacy cost.

## 4 Mitigations Against Targeted Risks

Our previous findings indicate that the personalization of  $k$ NN-LMs with a private datastore is more susceptible to data extraction attacks compared to fine-tuning a parametric LM with private data. At the same time, leveraging private data offers substantial utility improvements. Is there a more effective way to leverage *private* data in order to achieve a better balance between privacy and utility in  $k$ NN-LMs? In this section we focus on addressing privacy leakage in the context of *targeted* attacks (see definition in Section 3.1), where the private information can be readily detected from text. We consider several approaches to tackle these challenges in Section 4.1 and Section 4.2, and present the results in Section 4.3.

### 4.1 Sanization of Datastore and Encoders

As demonstrated in Section 3, the existence of private examples in the  $k$ NN-LMs’ datastore increase the likelihood of privacy leakage since they are retrieved and aggregated in the final prediction. Therefore, our first consideration is to create a sanitized datastore by eliminating privacy-sensitive text

segments. We propose the following three options for sanitization:

- Replacement with  $\langle \text{endoftext} \rangle$ : replace each privacy-sensitive phrase with the  $\langle \text{endoftext} \rangle$  token;
- Replacement with dummy text: replace each privacy-sensitive phrase with a fixed dummy phrase based on its type. For instance, if telephone numbers are sensitive, they can be replaced with "123-456-789"; and
- Replacement with public data: replace each privacy-sensitive phrase with a randomly selected public phrase of a similar type. An example is to replace each phone number with a public phone number on the Web.

The encoders in a  $k$ NN-LM is another potential source of privacy leakage. While it is typically optimized on target domain data to enhance performance, fine-tuning directly on private data in privacy-sensitive tasks may result in privacy leaks (Table 1). Similarly, the encoder can be sanitized by fine-tuning the pre-trained encoder  $\text{Enc}_{\text{public}}$  on a sanitized dataset that has had sensitive information removed.

### 4.2 Decoupling Key and Query Encoders

We propose using separate encoders for keys and queries in  $k$ NN-LMs, to allow for finer control over privacy preservation. For example, the encoder for queries can be the pre-trained encoder  $\text{Enc}_{\text{public}}$ , while the encoder for keys can be  $\text{Enc}_{\text{private}}$ , which is fine-tuned on the private data to optimize for performance. This way, the query encoder can be more resistant to privacy leakage, while the keys encoder can provide better query results. While it is not a common practice in  $k$ NN-LMs, we view the separation of key and query encoders as a promising approach to reduce the discrepancy between

SANITIZATION APPROACH	SANITIZED?		EVAL. PPL	# TOTAL	PHONE	EMAIL	URL
	Enc <sub>K</sub>	Enc <sub>Q</sub>					
NONE	✗	✗	16.12	54	25	23	6
REPLACED W/ < endoftext >	✓	✓	16.83	0	0	0	0
	✗	✓	16.24	0	0	0	0
	✓	✗	16.32	16	6	6	4
REPLACED W/ DUMMY PII	✓	✓	16.51	0	0	0	0
	✗	✓	16.29	0	0	0	0
	✓	✗	16.16	22	5	15	2
REPLACED W/ RANDOM PUBLIC PII	✓	✓	16.38	0	0	0	0
	✗	✓	16.29	1	1	0	0
	✓	✗	16.20	23	8	13	2

Table 2: Perplexity and extraction risk of  $k$ NN-LMs with different combinations of key encoder (Enc<sub>K</sub>) and query encoder (Enc<sub>Q</sub>), under different sanitization approaches. Privacy measurements are computed using the top 1000 candidates. For each sanitization approach, the configuration with the highest leakage is emphasized in **red**, and the configuration with the lowest leakage is highlighted in **green**; the best utility under sanitization is highlighted in **boldface**. Sanitizing Enc<sub>Q</sub> is crucial to eliminating the targeted risk.

the prompt and the datastore, and reduce privacy leakage.

The privacy risk of a  $k$ NN-LM can also be impacted by its hyper-parameters such as the number of neighbors  $k$ , and the interpolation coefficient  $\lambda$ . It is important to consider these hyper-parameters in the redesign of the  $k$ NN-LMs to ensure that the privacy-utility trade-off is well managed.

### 4.3 Experimental Results

As demonstrated in Table 2, applying sanitization to both the encoder and the datastore effectively *eliminates* privacy risk, resulting in no personally identifiable information (PII) being extracted. Among the three methods, the strategy of replacing PII with random public information for sanitization yields the highest utility. It achieves a perplexity of 16.38, which is only marginally worse than the perplexity of 16.12 achieved by the non-sanitized private model.

Table 2 also demonstrates that utilizing separate encoders for keys and queries enhances the model’s utility compared to using the same sanitized encoder for both. Specifically, we observe that when using the non-sanitized encoder for the query and the sanitized encoder for the key, privacy risks remain high due to the potential leakage from the  $p_{LM}$ . On the other hand, using the non-sanitized encoder for the key and the sanitized encoder for the query effectively eliminates privacy risk while still maintaining a high level of utility. This finding highlights the importance of sanitizing the query encoder in  $k$ NN-LMs.

We finally analyze the impact of key hyper-

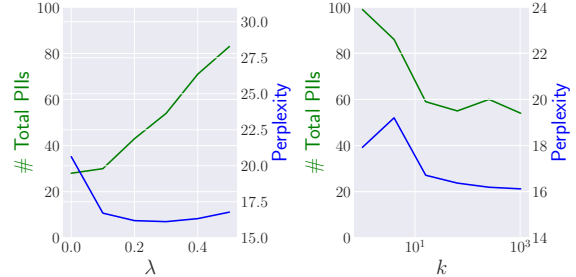


Figure 2: Effect of the interpolation coefficient ( $\lambda$ ) and the number of nearest neighbors ( $k$ ) on  $k$ NN-LMs’ utility (measured by perplexity, the blue curve) and privacy risk (measured by the number of reconstructed PIs in the targeted attack, the green curve). We use Enc<sub>private</sub> as encoders and  $\mathcal{D}_{private}$  as the datastore.

parameters on utility and privacy risks in  $k$ NN-LMs, using  $\mathcal{D}_{private}$  as datastore and Enc<sub>private</sub> for both Enc<sub>K</sub> and Enc<sub>Q</sub>. First, we vary  $\lambda$ , the interpolation coefficient, and observe that increasing  $\lambda$  decreases perplexity but increases privacy risk (see Figure 2). This highlights the trade-off between accuracy and privacy, indicating that optimizing both factors simultaneously is challenging through  $\lambda$  adjustment alone.

$k$  is the number of nearest neighbors in a  $k$ NN-LM. As shown in Figure 2, increasing the value of  $k$  in  $k$ NN-LMs improves perplexity as it allows considering more nearest neighbors. We also notice that using a larger  $k$  decreases privacy risk as the model becomes less influenced by a limited group of private nearest neighbors. Together, increasing  $k$  seems to simultaneously enhance utility and reduce privacy risk.

DATASTORE		EVAL. PPL	# GOOD RECON
$N_{\text{pub}}$	$N_{\text{priv}}$		
0	ALL	16.12	656
0	0	20.63	620 ( $\downarrow$ 5.5%)
ALL	0	21.46	561 ( $\downarrow$ 14.5%)
	5,000	21.58	579 ( $\downarrow$ 11.7%)
	10,000	21.23	595 ( $\downarrow$ 9.3%)
	50,000	19.16	617 ( $\downarrow$ 5.9%)
	ALL	19.02	632 ( $\downarrow$ 3.7%)

Table 3: Perplexity and data extraction risks for  $k$ NN-LMs with different numbers of public ( $N_{\text{pub}}$ ) and private ( $N_{\text{priv}}$ ) examples in the datastore. The encoder in use is the privately fine-tuned encoder  $\text{Enc}_{\text{private}}$ . Privacy measurements are computed using the top 5000 candidates.

$\text{Enc}_K$	$\text{Enc}_Q$	EVAL. PPL	# GOOD RECON
$\mathcal{D}_{\text{private}}$ w/ $\text{Enc}_{\text{private}}$		16.12	656
$\text{Enc}_{\text{private}}$	$\text{Enc}_{\text{private}}$	21.46	561 ( $\downarrow$ 14.5%)
$\text{Enc}_{\text{public}}$	$\text{Enc}_{\text{public}}$	33.60	0 ( $\downarrow$ 100%)
$\text{Enc}_{\text{private}}$	$\text{Enc}_{\text{public}}$	31.57	54 ( $\downarrow$ 91.8%)
$\text{Enc}_{\text{public}}$	$\text{Enc}_{\text{private}}$	22.90	451 ( $\downarrow$ 31.3%)
$\text{Enc}_{\text{mixed}}$	$\text{Enc}_{\text{mixed}}$	21.10	601 ( $\downarrow$ 8.4%)
$\text{Enc}_{\text{private}}$	$\text{Enc}_{\text{mixed}}$	21.12	545 ( $\downarrow$ 16.9%)
$\text{Enc}_{\text{mixed}}$	$\text{Enc}_{\text{private}}$	21.43	498 ( $\downarrow$ 24.1%)

Table 4: Perplexity and data extraction risks for  $k$ NN-LMs with different encoders for keys ( $\text{Enc}_K$ ) and queries ( $\text{Enc}_Q$ ). The datastore in use is the public datastore  $\mathcal{D}_{\text{public}}$ . Privacy measurements are computed using the top 5000 candidates. The  $\text{Enc}_{\text{mixed}}$  encoder is fine-tuned using a mix of public and private data points. Results suggest that using different encoders for keys and queries or  $\text{Enc}_{\text{mixed}}$  can potentially improve the privacy-utility trade-off.

## 5 Mitigations Against Untargeted Risks

In this section, we explore potential risks to mitigate untargeted risks in  $k$ NN-LMs, which is a more challenging setting due to the opacity of the definition of privacy. It is important to note that the methods presented in this section are preliminary attempts, and fully addressing untargeted risks in  $k$ NN-LMs still remains a challenging task.

### 5.1 Methods

Considering that storing  $\mathcal{D}_{\text{private}}$  in the datastore is the primary cause of data leakage (as discussed in Section 3), and the challenge of sanitizing private data in the face of untargeted risks, we propose the following approaches to leverage public data for mitigating these risks.

**Adding public data to datastore** The quality of the retrieved neighbors plays a crucial role in the performance and accuracy of  $k$ NN-LMs. Although it is uncommon to include public datapoints that are not specifically designed for the task or domain into  $k$ NN-LMs’ datastore, it could potentially aid in reducing privacy risks in applications that prioritize privacy. This becomes particularly relevant in light of previous findings, which suggest substantial privacy leakage from a private datastore.

**Fine-tuning encoders on a mixture of public and private data** However, adding public data may cause retrieval performance may suffer as there is a distribution gap between the public data (e.g., Web Crawl data) used to construct the datastore and the private data (e.g., email conversations) used for encoder fine-tuning. To address this issue, we propose further fine-tuning the encoder on a combination of public and private data to bridge the distribution gap and improve retrieval accuracy. The ratio for combining public and private datasets will be determined empirically through experimentation.

Similarly to Section 4.2, we could also employ separate encoders for keys and queries in the context of untargeted risks, which allows for more precise control over privacy preservation.

### 5.2 Experimental Results

Table 3 demonstrates that when a privately fine-tuned model  $\text{Enc}_{\text{private}}$  serves as the encoder, replacing the private datastore  $\mathcal{D}_{\text{private}}$  with a public one  $\mathcal{D}_{\text{public}}$  in  $k$ NN-LMs considerably lowers the privacy risk. Furthermore, when using  $\text{Enc}_{\text{private}}$  and  $\mathcal{D}_{\text{public}}$ , the risk level is slightly lower than when using the standard language model with  $\text{Enc}_{\text{private}}$  because the model’s final response has been interpolated with non-sensitive information, which helps to reduce privacy risks.

Using a public datastore reduces privacy risk but also results in a sudden drop in utility. If more stringent utility requirements but less strict privacy constraints are necessary, adding a few private examples to the public datastore, as shown in Table 3, may also be a suitable solution.

Table 4 demonstrates that using different encoders for keys ( $\text{Enc}_K$ ) and queries ( $\text{Enc}_Q$ ) is more effective in achieving a desirable balance between privacy and utility when using  $\mathcal{D}_{\text{public}}$  as the datastore. Specifically, using  $\text{Enc}_{\text{private}}$  to encode keys and  $\text{Enc}_{\text{public}}$  to encode queries significantly reduces the risk of data extraction with only a slight

decrease in perplexity.

We further try fine-tuning the encoder using a combination of public and private data, which results in  $\text{Enc}_{\text{mixed}}$ . The training dataset comprises the entire set of private data of size  $N_{\text{priv}}$  and  $N_{\text{priv}} \times r$  public data, where  $r$  takes values from  $\{0.01, 0.02, 0.05, 0.1, 0.2, 0.5, 1.0\}$ . We present attack results using  $r = 0.05$  as it achieves the best perplexity. As shown in Table 4, when the encoder is fine-tuned using a combination of public and private data, the perplexity can be enhanced from 21.46 to 21.10 while simultaneously reducing privacy risk. This is because  $\text{Enc}_{\text{mixed}}$  helps close the distribution gap between private and public data thus improving the retrieval results. Similarly, using separate  $\text{Enc}_K$  and  $\text{Enc}_Q$  also helps further reduce the privacy risk.

## 6 Related Work

### 6.1 Retrieval-based Language Models

Retrieval-based language models (Khandelwal et al., 2020; Borgeaud et al., 2022; Izacard et al., 2022; Zhong et al., 2022; Min et al., 2022) have been widely studied in recent years. These models not only rely on encoder forward running but also leverage a non-parametric component to incorporate more knowledge from an external datastore during inference. The retrieval process starts by using the input as a query, and then retrieving a set of documents (i.e., sequences of tokens) from a corpus. The language model finally incorporates these retrieved documents as additional information to make its final prediction. While the deployment of retrieval-based language models has been shown to lead to improved performance on various NLP tasks, including language modeling and open-domain question answering, it also poses concerns about data privacy.

### 6.2 Privacy Risks in Language Models

Language models have been shown to tend to memorize (Carlini et al., 2019; Thakkar et al., 2021; Zhang et al., 2021; Carlini et al., 2023) their training data and thus can be prompted to output text sequences from the its training data (Carlini et al., 2021), as well as highly sensitive information such as personal email addresses (Huang et al., 2022) and protected health information (Lehman et al., 2021; Pan et al., 2020). Recently, the memorization effect in LMs has been further exploited in the federated learning setting (Konečný et al., 2016),

where in combination with the information leakage from model updates (Melis et al., 2019), the attacker is capable of recovering private text in federated learning (Gupta et al., 2022). To mitigate privacy risks, there is a growing interest in making language models privacy-preserving (Yu et al., 2022; Li et al., 2022; Shi et al., 2022b; Yue et al., 2022; Cummings et al., 2023) by training them with a differential privacy guarantee (Dwork et al., 2006; Abadi et al., 2016) or with various anonymization approaches (Nakamura et al., 2020; Biesner et al.).

Although previous research has demonstrated the potential risks of data extraction in parametric language models, our study is the first investigation of the privacy risks associated with retrieval-based language models; we also propose strategies to mitigate them. The closest effort is Arora et al. (2022), which explores the privacy concerns of using private data in information retrieval systems and provides potential mitigations. However, their work is not specifically tailored to the context of retrieval-based language models.

## 7 Conclusion

There are several conclusions from our investigation of privacy risks related to  $k$ NN-LMs. First, our empirical study reveals that incorporating a private datastore in  $k$ NN-LMs leads to increased privacy risks (both targeted and untargeted) compared to parametric language models trained on private data. Second, for targeted attacks, our experimental study shows that sanitizing  $k$ NN-LMs to remove private information from both the datastore and encoders, and decoupling the encoders for keys and queries can eliminate the privacy risks without sacrificing utility, achieving perplexity of 16.38 (vs. 16.12). Third, for untargeted attacks, our study shows that using a public datastore and training the encoder on a combination of public and private data can reduce privacy risks at the expense of reduced utility by 24.1%, with perplexity of 21.12 (vs. 16.12).



## Limitations

The current study focuses on nearest neighbor language models, but there are many other variants of retrieval-based language models, such as RETRO (Borgeaud et al., 2022) and ATLAS (Izacard et al., 2022). Further research is needed to understand the privacy implications of these models and whether our findings apply. It would be also interesting to investigate further, such as combining proposed approaches with differential privacy, to achieve better privacy-utility trade-offs for mitigating untargeted privacy risks.

## Acknowledgement

This project is supported by an NSF CAREER award (IIS-2239290), a Sloan Research Fellowship, a Meta research grant, and a Princeton SEAS Innovation Grant. We would like to extend our sincere appreciation to Dan Friedman, Alexander Wettig, and Zhiyuan Zeng for their valuable comments and feedback on earlier versions of this work.

## References

- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318.
- Simran Arora, Patrick Lewis, Angela Fan, Jacob Kahn, and Christopher Ré. 2022. Reasoning over public and private data in retrieval-based systems. *arXiv preprint arXiv:2203.11027*.
- David Biesner, Rajkumar Ramamurthy, Robin Stenzel, Max Lübbering, Lars Hillebrand, Anna Ladi, Maren Pielka, Rüdiger Loitz, Christian Bauckhage, and Rafet Sifa. Anonymization of german financial documents using neural network-based language models with contextual word representations. *International Journal of Data Science and Analytics*, pages 1–11.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International Conference on Machine Learning (ICML)*. PMLR.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. 2023. Quantifying memorization across neural language models. In *International Conference on Learning Representations (ICLR)*.
- Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. 2019. The secret sharer: Evaluating and testing unintended memorization in neural networks. In *28th USENIX Security Symposium (USENIX Security 19)*, pages 267–284.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, et al. 2021. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650.
- Centers for Medicare & Medicaid Services. 1996. The Health Insurance Portability and Accountability Act of 1996 (HIPAA). Online at <http://www.cms.hhs.gov/hipaa/>.
- Mia Xu Chen, Benjamin N Lee, Gagan Bansal, Yuan Cao, Shuyuan Zhang, Justin Lu, Jackie Tsay, Yinan Wang, Andrew M Dai, Zhifeng Chen, et al. 2019. Gmail smart compose: Real-time assisted writing. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2287–2295.
- Rachel Cummings, Damien Desfontaines, David Evans, Roxana Geambasu, Matthew Jagielski, Yangsibo Huang, Peter Kairouz, Gautam Kamath, Sewoong Oh, Olga Ohrimenko, et al. 2023. Challenges towards the next frontier in privacy. *arXiv preprint arXiv:2304.06929*.
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. 2006. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*. Springer.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. *Hierarchical neural story generation*. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia. Association for Computational Linguistics.
- Samyak Gupta, Yangsibo Huang, Zexuan Zhong, Tianyu Gao, Kai Li, and Danqi Chen. 2022. Recovering private text in federated learning of language models. *arXiv preprint arXiv:2205.08514*.
- Junxian He, Graham Neubig, and Taylor Berg-Kirkpatrick. 2021. Efficient nearest neighbor language models. In *Empirical Methods in Natural Language Processing (EMNLP)*.
- Jie Huang, Hanyin Shao, and Kevin Chen-Chuan Chang. 2022. Are large pre-trained language models leaking your personal information? In *Findings of Empirical Methods in Natural Language Processing (EMNLP)*.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022. Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv:2208.03299*.

- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Generalization through Memorization: Nearest Neighbor Language Models. In *International Conference on Learning Representations (ICLR)*.
- Bryan Klimt and Yiming Yang. 2004. The enron corpus: A new dataset for email classification research. In *European Conference on Machine Learning*. Springer.
- Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. 2016. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*.
- Eric Lehman, Sarthak Jain, Karl Pichotta, Yoav Goldberg, and Byron C Wallace. 2021. Does bert pre-trained on clinical notes reveal sensitive data? In *North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Xuechen Li, Florian Tramer, Percy Liang, and Tatsunori Hashimoto. 2022. Large language models can be strong differentially private learners. In *International Conference on Learning Representations (ICLR)*.
- Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov. 2019. Exploiting unintended feature leakage in collaborative learning. In *2019 IEEE symposium on security and privacy (SP)*. IEEE.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*.
- Sewon Min, Weijia Shi, Mike Lewis, Xilun Chen, Wentau Yih, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Nonparametric masked language modeling. *arXiv preprint arXiv:2212.01349*.
- Yuta Nakamura, Shouhei Hanaoka, Yukihiro Nomura, Naoto Hayashi, Osamu Abe, Shuntaro Yada, Shoko Wakamiya, and Eiji Aramaki. 2020. Kart: Privacy leakage framework of language models pre-trained with clinical records. *arXiv preprint arXiv:2101.00036*.
- Arvind Narayanan and Vitaly Shmatikov. 2008. Robust de-anonymization of large sparse datasets. In *2008 IEEE Symposium on Security and Privacy (sp 2008)*. IEEE.
- Xudong Pan, Mi Zhang, Shouling Ji, and Min Yang. 2020. [Privacy risks of general-purpose language models](#). In *2020 IEEE Symposium on Security and Privacy (SP)*, pages 1314–1331.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Weijia Shi, Julian Michael, Suchin Gururangan, and Luke Zettlemoyer. 2022a. Nearest neighbor zero-shot inference. In *Empirical Methods in Natural Language Processing (EMNLP)*.
- Weiyan Shi, Aiqi Cui, Evan Li, Ruoxi Jia, and Zhou Yu. 2022b. Selective differential privacy for language modeling. In *North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Om Dipakbhai Thakkar, Swaroop Ramaswamy, Rajiv Mathews, and Francoise Beaufays. 2021. [Understanding unintended memorization in language models under federated learning](#). In *Proceedings of the Third Workshop on Privacy in Natural Language Processing*, pages 1–10, Online. Association for Computational Linguistics.
- Frank F Xu, Uri Alon, and Graham Neubig. 2023. Why do nearest neighbor language models work? *arXiv preprint arXiv:2301.02828*.
- Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A Inan, Gautam Kamath, Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, et al. 2022. Differentially private fine-tuning of language models. In *International Conference on Learning Representations (ICLR)*.
- Xiang Yue, Huseyin A Inan, Xuechen Li, Girish Kumar, Julia McAnallen, Huan Sun, David Levitan, and Robert Sim. 2022. Synthetic text generation with differential privacy: A simple and practical recipe. *arXiv preprint arXiv:2210.14348*.
- Chiyuan Zhang, Daphne Ippolito, Katherine Lee, Matthew Jagielski, Florian Tramèr, and Nicholas Carlini. 2021. Counterfactual memorization in neural language models. *arXiv preprint arXiv:2112.12938*.
- Zexuan Zhong, Tao Lei, and Danqi Chen. 2022. Training language models with memory augmentation. In *Empirical Methods in Natural Language Processing (EMNLP)*.

## A Training Data Extraction Attack

### A.1 Untargeted Attack

Carlini et al. (2021) proposes the first attack that can extract training data from a trained language model. The attack consists of two steps: 1) generate candidate reconstructions via prompting the trained models, and 2) sort the generated candidates using a score that implies the possibility of being a memorized text.

**Generate candidate reconstructions** The attacker generates candidates for reconstructions via querying the retrieval-augmented LM’s sentence completion API with contexts. Following Carlini et al. (2021)<sup>5</sup>, we randomly select chunks from a subset of Common Crawl<sup>6</sup> to feed as these contexts.

**Sort candidates by calibrated perplexity** The second step is to perform membership inference on candidates generated from the previous step. We are using the calibrated perplexity in our study, which has been shown to be the most effective membership metric among all tested ones by Carlini et al. (2021).

The perplexity measures how likely the LM is to generate a piece of text. Concretely, given a language model  $f_\theta$  and a sequence of tokens  $\mathbf{x} = x_1, \dots, x_l$ ,  $\text{Perplexity}(f_\theta, \mathbf{x})$  is defined as the exponentiated average negative log-likelihood of  $\mathbf{x}$ :

$$\exp \left( -\frac{1}{n} \sum_{i=1}^l \log f_\theta(x_i | x_1, \dots, x_{i-1}) \right) \quad (1)$$

A low perplexity implies a high likelihood of the LM generating the text; For a retrieval-augmented LM, this may result from the LM has been trained on the text or has used the text in its datastore.

However, perplexity may not be a reliable indicator for membership: common texts may have very low perplexities even though they may not carry privacy-sensitive information. Previous work (Carlini et al., 2019, 2021) propose to filter out these uninteresting (yet still high-likelihood samples) by comparing to a second LM which never sees the

private dataset. Specifically, given a piece of text  $x$  and the target model  $f_\theta$ , and the reference LM  $f_\theta^{ref}$ , the calibrated perplexity computes the ratio  $\text{Perplexity}(f_\theta, \mathbf{x}) / \text{Perplexity}(f_\theta^{ref}, \mathbf{x})$ .

### A.2 Targeted Attack

The untargeted attack has demonstrated the feasibility of recovering an entire sentence from the deployed retrieval-augmented LM. However, it is possible that only a small segment of a sentence contains sensitive information that can act as personal identifiers, and thus be of interest to the attacker. Therefore, we also consider the type of attack which specifically targets this type of information.

We define personal identifiers and describe the attack method and evaluation subsequently.

#### A.2.1 Definition of Personal Identifiers

Personal Identifiable Information (PII) refers to any data that can be used to identify a specific individual, such as date of birth, home address, email address, and telephone number. PII is considered sensitive information and requires proper protection to ensure privacy.

The exact definition of PII can vary depending on the jurisdiction, country, and regulations in place. One of the clearest definitions of PII is provided by Health Insurance Portability and Accountability Act (HIPAA) (Centers for Medicare & Medicaid Services, 1996), which includes name, address, date, telephone number, fax number, email address, social security number, medical record number, health plan beneficiary number, account number, certificate or license number, vehicle identifiers and serial numbers, web URL, IP Address, finger or voice print, photographic image, and any other characteristic that could uniquely identify the individual. In our study, we focus on three frequently investigated PII in previous literature (Huang et al., 2022; Carlini et al., 2021), including email addresses, telephone numbers, and URLs.

#### A.2.2 The Attack

It’s important to note that our approach differs from the work of Huang et al. (2022), which aims to reconstruct the relationship between PIIs and their owners<sup>7</sup>. Instead, our study focuses on reconstructing the *actual values* of PIIs. This is because, even

<sup>5</sup>They have empirically show that sampling conditioned on Internet text is the most effective way to identify memorized content, compared with top- $n$  sampling (Fan et al., 2018) and temperature-base sampling (see Section 5.1.1 in their paper).

<sup>6</sup>Common Crawl is a nonprofit organization that crawls the web and freely provides its archives and datasets to the public. See their webpage for details: <http://commoncrawl.org/>

<sup>7</sup>This threat model requires additional information about the presence of the owners in the dataset.

Phone	Email	URL
If you have questions, please feel free to give me a call at	For more information, send email to	The site can be found at
Please advise or call me at	For more information please email us at	For more information, visit
Please call us at	Suggestions and feedback are welcome at	Please visit our web site at
I can be reached at	For more information please email us at	Visit our home page at
If you have any questions, please call	Please forward this e-mail to	For more details go to

Table 5: Example extraction prompts for different types of PII.

PII	# EMAILS CONTAINING PII	# UNIQUE VALUES
EMAIL	1,013	758
PHONE	1,921	1,621
URL	1,641	1,396

use beam search with repetition penalty = 0.75 for the generation.

Table 6: Number of PII in the Enron Email dataset.

if the attacker cannot determine the relationship through the current attack, the reconstruction of PII is already considered identity theft<sup>8</sup>. Further, the attacker can use the linkage attack (Narayanan and Shmatikov, 2008) with the aid of publicly available information to determine the relationship between PII and their owners.

**PII extraction.** Similar to the training data extraction attack, the PII extraction attack consists of two steps: 1) generate candidate reconstructions, and 2) sort them using membership metrics.

To tailor the attack to recover personal identifiable information rather than the entire training text, we customize the attack prompts based on the type of information to be extracted.

## B Experimental details

**PIIs in Enron Email Dataset.** We use regular expressions to identify and extract three types of personal identifiers from the Enron Email training dataset for the use of the targeted attack, including telephone numbers, email addresses, and URLs. Table 6 provides statistics for these personal identifiers.

**Prompts for the targeted attack.** We gather common preceding context for telephone numbers, email addresses, and URLs, and use them as prompts for the targeted attack. Table 5 provides example prompts we use in the attack.

**Attack parameters.** For the untargeted attack, we generate 100,000 candidates, and for the targeted attack, we generate 10,000 candidates. We

<sup>8</sup>[https://en.wikipedia.org/wiki/Identity\\_theft](https://en.wikipedia.org/wiki/Identity_theft).