

Protecting Privacy in Classifiers by Token Manipulation

Anonymous ACL submission

Abstract

Using language models as a remote service entails sending private information to an untrusted provider. In addition, potential eavesdroppers can intercept the messages, thereby exposing the information. In this work, we explore the prospects of avoiding such data exposure at the level of text manipulation. We focus on text classification models, examining various token mapping and contextualized manipulation functions in order to see whether classifier accuracy may be maintained while keeping the original text unrecoverable. We find that although some token mapping functions are easy and straightforward to implement, they heavily influence performance on the downstream task, and via a sophisticated attacker can be reconstructed. In comparison, contextualized manipulation provides an improvement in performance.

1 Introduction

Large language models (LLMs) have greatly advanced the field of NLP in recent years, exhibiting exceptional proficiency across a wide spectrum of tasks, including dependency parsing (Duong et al., 2015), natural language understanding (Dong et al., 2019), automatic question-answering (OpenAI, 2021; Ouyang et al., 2022), machine translation (Dabre et al., 2020), text classification (Minaee et al., 2021), and many more (Li et al., 2022). However, this success comes with potential privacy risks, as the models process vast amounts of data that might contain personal or sensitive information and may abuse or leak it. For instance, information can be leaked by model inversion (Li et al., 2017), re-identification techniques (Lison et al., 2021; Ben Cheikh Larbi et al., 2023), exploitation of feature memorization within the LLM (Carlini et al., 2021), and more. Offering LLMs as cloud services, such as ChatGPT (Ouyang et al., 2022), might also impose potential threats to privacy if the server exhibits a semi-honest stance, actively

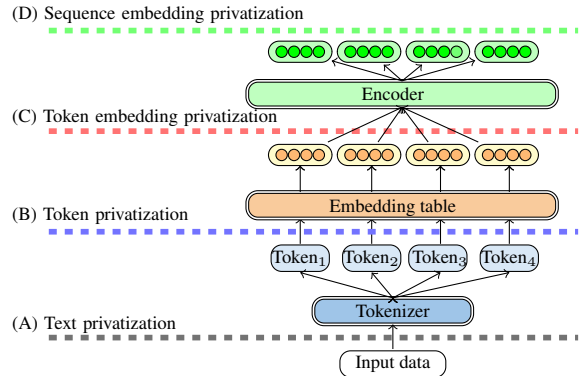


Figure 1: A schematic of the various stages where differential privacy techniques can be applied in an LLM. This work focuses on level (B).

seeking to glean more insights from the input than is appropriate or by a possible eavesdropper intercepting the input sent to the server.

In order to safeguard privacy, many privacy-preserving techniques have been proposed, based on the local differential privacy framework (LDP; Arachchige et al., 2019). In this framework, the user applies a differential privacy mechanism, which can be hosted on a local server, and then sends the privatized data to the remote server. This approach doesn't require trust from the remote server, and protects the data against potential eavesdroppers. In general, any privacy mechanism can be applied at one or several components of the LLM pipeline. Figure 1 depicts these components: at the text level (*text privatization*), after the tokenization process (*token privatization*), after the initial embedding lookup (*token embedding privatization*), or after applying several layers of the encoder (*sequence embedding privatization*).

Currently, most privacy-preserving strategies focus on incorporating noise into sequence embedding vectors. The rationale behind this strategy is to minimize the privacy-preserving technique's impact on the downstream task. Specifically, most systems first obtain a sequence embedding repre-

067 presentation, either by assuming partial access to the
068 remote model (Zhou et al., 2022; Lyu et al., 2020;
069 Qu et al., 2021) or by using a dedicated model to
070 create these embeddings (Li et al., 2018; Coavoux
071 et al., 2018; Mosallanezhad et al., 2019; Plant et al.,
072 2021; Zhou et al., 2023). Afterwards, random noise
073 is incorporated into the embeddings, thus conceal-
074 ing the original input. However, this approach re-
075 lies on partial access to the remote model, on the
076 ability to provide input to the remote model in vec-
077 tor form, or on sufficient computational and mem-
078 ory resources on the user’s end. These are often not
079 the case. In addition, Kugler et al. (2021) showed
080 that publishing a model’s encoder along with the
081 contextualized embeddings allows an adversary to
082 generate data to train a decoder with a high level of
083 reconstruction accuracy, making these approaches
084 highly susceptible to violation of privacy.

085 We propose **a secure way to use LLMs with-**
086 **out assuming access to their parameters.** In our
087 framework, both input and output for the privacy-
088 providing mechanism must be given in a token
089 sequence format, eliminating the need to intervene
090 with the LLM’s pre-training procedure or text pro-
091 cessing. We focus on applying privacy preserva-
092 tion techniques at the token level, corresponding to
093 layer (B) in Figure 1.

094 Specifically, we propose two privacy-preserving
095 techniques based on manipulating **the input token**
096 **sequence.** The first set of techniques relies on naïve
097 rules of token substitution. The second is based on
098 leveraging contextual information to strategically
099 replace tokens, aiming to retain as much actionable
100 information as possible for the classifier to mini-
101 mize the impact on the performance of the down-
102 stream task. We test these techniques both for their
103 impact on the downstream task accuracy and for
104 their resilience against reconstruction attacks. We
105 find that replacing tokens based on simple rules is
106 easy for a knowledgeable attacker to reverse, while
107 manipulating tokens based on contextual informa-
108 tion can enhance privacy without sacrificing much
109 of the performance.

110 2 Lossy Mapping

111 In order to protect against potential eavesdropping
112 by a middle party, under the assumption that the
113 layers of LLMs are inaccessible to the local device,
114 we start by employing several mapping functions
115 on the tokens of the input text available at the lo-
116 cal device. Our initial, naïve mapping functions

117 introduce a random noise component that follows
118 a specific rule: the vocabulary is partitioned into
119 pairs of tokens (u, v) , or triplets (u, v, z) , and when
120 encountered in an input text to be manipulated,
121 all tokens are mapped to a single representative
122 token of their tuple, without loss of generality u .
123 This strategy produces outputs that are inherently
124 ambiguous, blocking any potential eavesdroppers
125 from recovering the original input text determi-
126 nistically, given that a many-to-one mapping is not
127 invertible. The only available recourse for an at-
128 tacker is a statistical strategy, which imposes as-
129 sumptions on the properties of the input, for ex-
130 ample that it was grammatical English text written
131 by a speaker with high proficiency. Indeed, even
132 if an eavesdropper obtains full information of the
133 privacy system, i.e. the partition into token tuples
134 and each tuple’s representative token, each mapped
135 sequence of length m still generates a candidate set
136 of 2^m or 3^m possible permutations (depending on
137 tuple size) through which the attacker must search.
138 We will examine the practical implications of this
139 large search space later in the section.

140 For our stated use case of manipulating text be-
141 ing input into a sequence classifier operating atop
142 an LLM, there are two distinct scenarios depend-
143 ing on when we may apply our manipulation. The
144 first scenario involves applying the manipulation
145 process only during the inference phase of a model
146 trained on regular, unmanipulated text, which we
147 will refer to as the TEST case. This operation mode
148 simulates a query sent by a user to an already-
149 trained model, such as a user interacting with Chat-
150 GPT or another model allowing only inference text
151 interaction via user interface or an API. In the sec-
152 ond scenario, which we call ALL, we also apply
153 the manipulation during the training phase, pro-
154 tecting sensitive information in the training data,
155 hoping that the inference phase will now leverage
156 the model’s ability to handle manipulated input as
157 expected and produce better results. In this scenario
158 the model does not inadvertently learn or memorize
159 the sensitive data during the training process, nor
160 does it spend learning resources on tokens never to
161 be seen during inference, but since it is not always
162 possible to assume its availability, we perform our
163 experiments in both settings.

164 When protecting the original input data, it is es-
165 sential for the mapper to have minimal impact on
166 the performance of the downstream task, defining
167 the fundamental trade-off in our study. Therefore,

Dataset	Mapper	TEST	ALL	Unchanged Tokens
SST2	Plain text	94.5%	94.5%	100%
	2-Random	75.0%	85.0%	51.0%
	3-Random	62.0%	80.0%	34.0%
	High-freq	90.0%	91.0%	93.0%
	Low-freq	60.0%	78.0%	7.0%
IMDb	Plain text	95.0%	95.0%	100%
	2-Random	75.0%	90.0%	50.0%
	3-Random	68.0%	85.0%	32.0%
	High-freq	93.0%	94.0%	94.0%
	Low-freq	60.0%	80.0%	6.0%

Table 1: The mapping strategy accuracy on SST2 and IMDb datasets and the percentage of unchanged tokens after applying the mappers to the training and test sets.

the selection process for grouping tokens and selecting each tuple’s representative token is crucial, as it aims to both minimize the mapping’s effect on the downstream task and hinder the attacker’s ability to uncover the original text. We consider the following mapping functions:

Purely random mapping the selection of the token pairs tuples from the vocabulary and of each tuple’s representative is uniformly random.

High-frequency mapping token pairs are selected based on their frequency of occurrence in a tokenized corpus, such as Wikipedia (Foundation, 2023). This involves pairing a higher-frequency token with a lower-frequency token, with the higher-frequency token being designated as the representative. In our mapper, given a vocabulary of even size V , sorted by descending frequency, each token with rank $1 \leq k \leq \frac{V}{2}$ is paired with the token of rank $k + \frac{V}{2}$. While selecting the high-frequency token as the representative may have a lesser impact on the downstream task, it could potentially weaken the privacy-preserving characteristics, depending on the knowledge possessed by the attacker.

Low-frequency mapping the process is similar to that of the higher-frequency mapper, except that the lower-frequency token is chosen as the representative. Opting for less-frequent tokens as representatives can aid in preserving privacy, but it will likely harm the downstream task.

Due to the simplicity of these mapping strategies, we consider them baselines for further research and developing better, potentially language-aware strategies. In addition, these mapping functions can easily be generalized to larger tuples, expanding the search space even further, but greatly harming

Mapper	Text		
Plain Text	no	apparent	joy
2-Random	his	buffers	University
High-freq	no	apparent	joy
Noise(150)	non	evident	joyful
STEN(9, 0.8)	No	evident	joyful
STEN _p (9, 1.0)	apparent	No	joyful

Table 2: Examples of the privatized textual sequences obtained with different privacy-preserving techniques.

downstream task performance as a result of a much more restricted active vocabulary.

2.1 Task Performance

To assess the impact of the baseline models on downstream task performance, we use two datasets for sequence classification: SST2 (Socher et al., 2013) and IMDb (Maas et al., 2011). The base model chosen was RoBERTa (Liu et al., 2019), a state-of-the-art encoder language model known for its strong performance in sequence classification tasks. In Table 1, we present the results of four baselines on the two datasets, compared with the null mapping results labeled “Plain text”. Perhaps unsurprisingly, the high-frequency baseline achieved the highest accuracy, most likely due to the fact that retaining high-frequency tokens while removing low-frequency ones results in a relatively small number of tokens altered in the datasets. In both datasets this number is roughly 6%, compared with low-frequency mapping’s complement of 94% and with the randomly-selected sets’ 50% and 67%, giving a correlative relationship between this number and the performance level: the fewer tokens are altered, the better the model performs. This effect is much more pronounced when only the test set is affected, and the model is dealing not only with loss of information but also with out-of-distribution behavior. In absolute terms, we find it remarkable that this alteration of a non-negligible portion of tokens causes only a 1–2 percentage point reduction in performance for the IMDb dataset and still under 5 points for SST2.

In Table 2, we present an example of the outcome of applying the 2-Random and the High-freq privatization techniques on a random phrase (“no apparent joy”) from the SST2 dataset. As expected, the 2-random baseline produces a random sequence of words, whereas the high-frequency mapper leaves the phrase unchanged as the tokens in the original sequence are frequent.

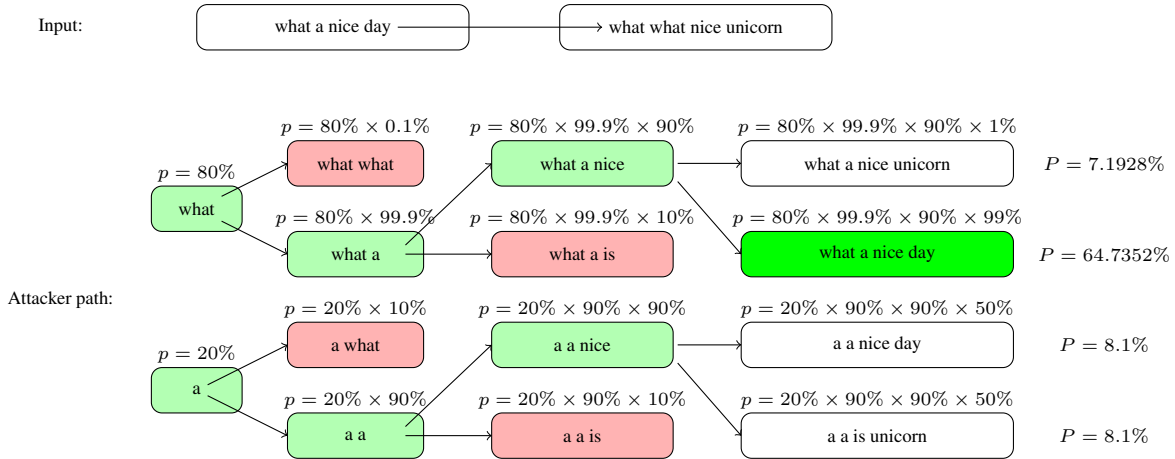


Figure 2: Schematic overview of the proposed heuristic oracle attacking scenario path over trying to reconstruct the sentence “what a nice day” which is remapped to “what what nice unicorn”. The red boxes indicate that the probability (presented above the box) of the candidate is low enough to be dropped in the next step, while the green boxes are the candidates that will be expanded in the next step.

2.2 Brute-force Attacker

Although the many-to-one mapping function introduces some form of protection against data leakage, in practice, reconstructing the original text might be relatively straightforward under certain circumstances. In particular, if an “oracle” attacker has access to the token pairings, it can theoretically determine the original text from the pool of 2^m possible permutations by applying a generative LLM such as GPT (Radford et al., 2019) and picking the most probable sequence. However, generating and evaluating all 2^m permutations is impractical even for small values of m due to the computational complexity involved. To mitigate this challenge, alternative approaches, such as employing heuristics or utilizing statistical methods, can be explored to narrow down the potential candidates for the original text.

To cope with this task, we describe a heuristic approach to reducing the search space based on **beam search** (Eisenstein, 2019, §11.3.1) and **nucleus sampling** (Holtzman et al., 2019). In each step of the process, candidates are generated based on the prefixes of tokens that were produced in the previous steps. In the case of token pairs, each prefix sequence is followed by one of two candidate tokens for the next step based on the known (oracle) token pair that the observed representative token belongs to. Unlike conventional beam search, where a fixed number of candidates is retained following each step, we opt for a dynamic approach inspired by nucleus sampling, made possible since the scores for each of the two tokens

reflect a generative probabilistic process where the relative probability of each interim token sequence on the beam can be estimated and used for dropping highly unlikely sequence prefixes. This means that the number of candidates remaining on the beam varies at each step, adapting to their likelihood and ensuring flexibility in the selection process. We estimate the likelihood of each candidate prefix using a language model.¹ After all prefixes on the beam have been scored, we remove the least probable candidates such that the total probability of the remaining candidates exceeds a certain threshold π set by computational constraints but maintaining discoverability. Since the probability of a sequence cannot exceed that of its prefix, the process guarantees that complete sequences that are likely are not being discarded before getting the chance to be fully generated. Overall, this process effectively eliminates highly unlikely candidates, dramatically reducing the search space during its application and streamlining the computational efforts.

This process is illustrated in Figure 2. The “oracle” attacker gains access to the remapped words: (what, a) → a, (nice, is) → nice, (day, unicorn) → unicorn. In the first step, two initial candidates (what and a) are generated based on the first observed token (what). Following the described process, each prefix is evaluated via an LLM to determine its probability, for instance, the probability of what being the first word is 80% when considering the possible set $\{[s] \text{ what,}$

¹<https://github.com/simonepri/lm-scorer>

Dataset	Mapper	MRR (↓)	Pr@5 (↓)	Edit dist (↑)
SST2	2-Random	0.89	0.97	1.32
	3-Random	0.81	0.92	1.35
	High-freq	0.86	0.98	1.33
IMDb	2-Random	0.48	0.59	1.60
	3-Random	0.45	0.53	1.70
	High-freq	0.63	0.72	1.60

Table 3: The three random mappings’ capability of preserving privacy against an “oracle” attacker. Edit distance is calculated at the token level.

[s] a}. This process is repeated, and the candidates with low probability are removed, such that the total probability of the remaining candidates is above 85%, as indicated by the red boxes. Finally, the probability of the sequence `what a beautiful day` is the highest, thus the “oracle” attacker returns it as the inferred original text. We note that the low-frequency and high-frequency mappers, despite their differences in representative token selection, will demonstrate equivalent safeguarding mechanisms against this attacker since the attacker does not factor in the choice of the representative token and examines all potential candidates in its effort to uncover the original text.

2.3 Resilience Against Reconstruction Attacks

In Table 3, we present the outcomes of the attacker’s endeavors to reveal the original text from the three techniques: 2-Random, 3-Random, and High-freq (equivalent to Low-freq for a knowledgeable attacker). We report the mean reciprocal rank (MRR) of the correct sequences, the rate of the actual input sequence ranking among the top 5 predictions (Pr@5), and the token-level edit distance between the produced top prediction and the original sequence. The relative success of the mappers in thwarting the oracle attacks on the IMDb dataset compared to SST2 can be attributed to the average token sequence length (\bar{m}), which is 65 and 12, respectively. As sequence length increases, the attacker’s task of uncovering the original text becomes more challenging.

Our results indicate that the naïve baselines are overly simplistic and allow an easy and straightforward reconstruction, even within a vast search space (although attacker knowledge of the mapping specifications is required). In cases where performance on the task remains close to that of unmapped text, the recovery price is too high to ne-

glect. Having said that, the computational complexity of applying the naïve baselines is relatively low, and the greatly reduced active vocabulary brings great savings in parameter budgets, which embedding tables often dominate. In a less powerful attack environment, this would make them an efficient choice for preserving privacy on low-resource devices. We expect future work on more principled many-to-one static mappings would be able to improve both task performance and resilience to attackers, while work on attack strategies can present challenges hitherto unseen.

3 STENCIL Privacy Preservation

In the context of protecting privacy within NLP practices, a widely adopted approach for implementing local differential privacy involves introducing a controlled level of *noise* into different components of the model, effectively concealing the original input. These components may include sequence embeddings, token embeddings, or the tokens themselves (Mosallanezhad et al., 2019; Feyisetan et al., 2020; Lyu et al., 2020; Qu et al., 2021; Zhou et al., 2022). However, in essence, the success of models in most NLP tasks is primarily attributed to their effective utilization of contextual information. Moreover, our study focuses on token-level privacy preservation, i.e., we assume that the parameters of the LLMs are inaccessible, making the importance of contextual information more pronounced. Therefore, a fundamental limitation associated with incorporating noise is the exclusion of contextual information when defining the noise. This omission may hinder the potential benefits contextual details can offer for maintaining the performance of the downstream tasks.

Given this limitation, we propose a new privacy preservation technique, which we call STENCIL.² With this technique, a mapped token in a sequence “absorbs” information from adjacent tokens to form a new context-aware token, effectively concealing the original token while retaining information beneficial for maintaining task performance.

In order to generate the new contextualized token $t_k \rightarrow t'_k$, we first retrieve an embedding vector representation of the neighborhood, of size $n + 1$, containing the tokens $t_i, \forall i \in \{k - n/2 \dots k + n/2\}$ using some embedding lookup table $\mathbf{E} \in \mathbb{R}^{V \times d}$,

²This term hails from numerical analysis (Spotz, 1995), where it denotes a computation that involves the surrounding values.

which can be trained independently in a preliminary step or obtained from an available model such as the target model itself. We then subject the $n + 1$ embedding vector representations to a weighted transformation and incorporate them to form a new “quasi-embedding” vector $\sum_{i=k-n/2}^{k+n/2} f_i \cdot \mathbf{E}[t_i]$. Finally, we return the token t'_k that is closest to the quasi-embedding vector in the embedding space, based on cosine-similarity or euclidean distance computation, as an output. To further enhance privacy, we ensure that the new token is different from the original one. Formally, the process can be defined as follows:

$$t'_k = \arg \min_{t_j \in \mathcal{V}} \left\| \mathbf{E}[t_j] - \sum_{i=k-\frac{n}{2}}^{k+\frac{n}{2}} f_i \cdot \mathbf{E}[t_i] \right\|, \quad (1)$$

where \mathcal{V} is the vocabulary and f_i is the weighted transformation function of the tokens such that $\sum_{i=k-\frac{n}{2}}^{k+\frac{n}{2}} f_i = 1$.

The level of privacy enhancement and its impact on the downstream task by employing the STENCIL method can be managed by adjusting the window size and the properties of the weighted function f . In our study, we use the gaussian smoothing function as the weighted function. Consequently, the standard deviation, σ , plays a crucial role in the performance and amount of privacy achieved.

In our experiments, we compared our STENCIL mechanism to two other privacy-preserving techniques. The first, technique was proposed by [Qu et al. \(2021\)](#)’s, namely the NOISE mapper. In contrast to our proposed technique, this approach does not consider context but rather incorporates random noise into token embeddings to enhance privacy. Similar to our proposed method, the new token is the closest to the quasi-embedding vector in the embedding space. The random noise is obtained by multiplying a sample from a Gamma distribution $\Gamma(d, 1/\eta)$ and a uniform sample from a unit hypersphere, where η corresponds to the amount of noise introduced to the original token and d is the dimension of the embedding space.

For the second technique, we include [Chen et al. \(2023\)](#)’s CUSTEXT⁺ privacy-preserving mechanism. The CUSTEXT⁺ mechanism consists of a mapping procedure and a sampling function. The mapping procedures generate a list of the top K tokens for each token, selecting those with the highest semantic relevance to the original token. Similar to the NOISE and STENCIL mechanisms, semantic

relevance is determined by calculating either the cosine similarity or Euclidean distance of the quasi-embedding vectors. Then, each token is remapped to one of the K candidates using an exponential sampling function.

We note that the most time-intensive operation in all mechanisms is searching for the closest token to the perturbed quasi-embedding vector, whereas all other operations are negligible in comparison. Overall, the average computational cost per token is 0.005 seconds on two 16-core 3.2 GHz AMD EPYC 7343 Milan processors.

3.1 STENCIL⁺ and STENCIL_p Mechanisms

Identifying sensitive words, such as those involved in named entity recognition (e.g., names, addresses, workplaces), is a challenging task typically approached using statistical methods ([Liu et al., 2017](#); [Cohn et al., 2019](#); [Poostchi et al., 2018](#); [Friedrich et al., 2019](#)). As a result, our mechanism treats all tokens as sensitive since we cannot reliably distinguish sensitive from non-sensitive tokens. However, since treating stopwords as non-sensitive may pose a low privacy risk ([Chen et al., 2023](#)), we propose STENCIL⁺, which applies the STENCIL mechanism to all words except stopwords, thereby enhancing accuracy while maintaining privacy.

An additional variation of STENCIL, namely STENCIL_p, can be obtained by excluding the target token from the computation of the quasi-embedding vector in (1) by setting f_k to zero. This exclusion significantly improves the privacy of each token and diminishes the attacker’s ability to reconstruct the original token at the expense of performance.

3.2 Downstream Task Performance

To evaluate the impact of the STENCIL, NOISE and CUSTEXT⁺ methods on the model performance, we repeat the methodology outlined in §2: we use RoBERTa as the base model; SST2 and IMDb as the datasets; and the two distinct application cases: manipulating tokens on inference data only (TEST), and applying the technique during the training phase as well (ALL). However, as these privacy techniques exhibit a realistic case, we also test it on an encoder-decoder model T5-small ([Raffel et al., 2020](#)) on the QNLI task from the GLUE dataset ([Wang et al., 2019](#)). As in [Raffel et al. \(2020\)](#), we concatenate the question and its corresponding sentence to form a single sequence that serves as the input, while the target prediction is

Dataset	Mapper	TEST (\uparrow)	ALL (\uparrow)	Pr@5 (\downarrow)
SST2	Plain Text	94.5%	94.5%	-
	NOISE(100)	80.0%	87.8%	70.0%
	NOISE(150)	83.0%	90.0%	75.0%
	CUSTEXT ⁺	79.4%	82.5%	70.0%
	STEN(9, 0.8)	83.5%	89.3%	49.0%
	STEN ⁺ (9, 0.8)	85.3%	89.5%	47.0%
	STEN _p (9, 1.0)	84.7%	87.0%	0.0%
	STEN _p ⁺ (9, 1.0)	85.0%	89.4%	0.0%
IMDb	Plain Text	95.0%	95.0%	-
	NOISE(100)	89.0%	92.6%	86.0%
	NOISE(150)	90.0%	93.5%	90.0%
	CUSTEXT ⁺	88.9%	91.1%	90.0%
	STEN(9, 0.8)	90.2%	93.1%	67.0%
	STEN ⁺ (9, 0.8)	92.4%	94.0%	69.0%
	STEN _p (9, 1.0)	89.7%	92.4%	0.0%
	STEN _p ⁺ (9, 1.0)	89.7%	92.4%	0.0%
QNLI	Plain Text	88.1%	88.1%	-
	NOISE(100)	80.0%	84.0%	93.0%
	NOISE(150)	81.1%	84.4%	93.0%
	CUSTEXT ⁺	78.5%	81.5%	85.0%
	STEN(9, 0.8)	74.8%	83.1%	54.0%
	STEN ⁺ (9, 0.8)	81.4%	84.8%	52.3%
	STEN _p (9, 1.0)	67.9%	82.5%	0.0%
	STEN _p ⁺ (9, 1.0)	71.4%	83.8%	0.0%

Table 4: The best results achieved by the different STENCIL mapper variations, the NOISE mapper and CUSTEXT⁺ considering the Test and All cases on the SST2, IMDb, and QNLI datasets. Pr@5 represents the average token hit managed by the nearest-neighbor attacker.

either “entailment” or “not_entailment”, thus forming a classification task.

We report three distinct manipulations based on STENCIL, STENCIL⁺, and STENCIL_p. The weighting function f_i , for all three variations, is derived from a gaussian smoothing. For the STENCIL and STENCIL⁺ mechanism, we consider a standard deviation of $\sigma = 0.8$, and the number of adjacent tokens considered is set to nine (four from each side, as well as the target token). The standard deviation we consider for the STENCIL_p approach is $\sigma = 1.0$, with a window width of nine. In all approaches, to preserve model performance, the tokenizer and embedding lookup table used to derive the new tokens were sourced directly from the model being trained. For the NOISE mechanism, we report the two best η values: $\eta = 100, 150$. For the CUSTEXT⁺ mechanism, the K parameter was set to 20 with the privacy parameter $\epsilon = 3$, which yields the overall best results.

The results are presented in Table 4. The overall best accuracy is achieved by STENCIL and STENCIL⁺, demonstrating the advantage of uti-

lizing contextual information to achieve privacy and maintain high performance. Nevertheless, in the SST2 dataset, NOISE ($\eta = 150$) achieves the highest accuracy. NOISE with $\eta = 150$ introduces minimal noise, resulting in negligible alterations to the original tokens. However, the NOISE method comes at great cost in discoverability, to be presented in §3.3.

Compared to the sentiment analysis tasks (SST2 and IMDb), the QNLI task presents greater challenges, primarily due to the complex logical connections required for the model to discern entailment between the given sentence and question. Therefore, despite its instance sizes being very similar to those of IMDb (62 vs. 65), the fact that noise-based perturbations disrupt contextual and semantic information leads to a significant decrease in the model’s ability to discern the logical connections between the parts of the input. This results in a more pronounced performance degradation compared to the long-sequenced IMDb on the TEST case. In contrast, training the model on the noisy data (the ALL setup) proves effective in overcoming this effect, leading to improved results for T5-small.

In Table 2, we present an example of the outcome of applying STENCIL, STENCIL_p, and the NOISE mapper on a random phrase from the SST2 dataset. The NOISE mapper with a value of $\eta = 150$ introduces negligible noise, thus producing a similar sequence to the original one. The STENCIL-based techniques also produce a similar sequence, although STENCIL_p swaps the positions of some tokens as a direct result of excluding the target token from the obfuscation process.

3.3 Nearest-neighbor Reconstruction

An attacker can potentially exploit the fact that these techniques utilize contextualized tokens and the selection of the nearest token as the quasi-embedding vector (Qu et al., 2021). Specifically, given the new perturbed token t' , the attacker can obtain the embedding vector representation $\mathbf{E}[t']$. Afterward, the attacker can calculate the cosine similarity between $\mathbf{E}[t']$ and the other embedding vector representations ($\mathbf{E}[t]$ where $t \in \mathcal{V} \setminus \{t'\}$) and statistically determine the original token.

Additionally, since STENCIL incorporates information from its neighboring tokens, the new perturbed token t'_k might resemble the original token of a neighboring token, for instance, t_{k+1} . This

563	allows the attacker to attempt to determine the original token by analyzing and comparing the neighbors' most similar tokens.	613
564		614
565		615
566	To test the resilience of these techniques against token inversion attacks, we implement the described attacker and report whether the original token was found to be one of the nearest five (Pr@5), or its neighbor's nearest five.	616
567		617
568		618
569		619
570		620
571	The success rate of the attacker for the four techniques is presented in Table 4. While the minor alterations in the original tokens contributed to performance improvement in the NOISE mapper, it is found to be highly vulnerable to simple reconstruction attacks. Taking into account both accuracy and resilience against reconstruction attacks, the STENCIL method demonstrates the best results, with a marginal trade-off in performance. The STENCIL _p demonstrates the best privacy protection against the attacker, highlighting the effectiveness of excluding the target token from the computation of the new token.	621
572		622
573		623
574		624
575		625
576		626
577		627
578		628
579		629
580		630
581		631
582		632
583		633
584	4 Conclusion	634
585	In this paper, we propose several token manipulation methods to preserve privacy under the assumption that the model parameters are inaccessible. We first introduce four simple mappers that offer distinct advantages compared to existing privacy-preserving techniques. Notably, these mappers operate independently of the LLM and the specific downstream task, resulting in a high degree of versatility. Additionally, their computational complexity is relatively low, making them efficient choices for privacy preservation on local, low-resource devices. However, it is essential to acknowledge that these mappers harm the performance of the downstream tasks and can be easily reconstructed by a knowledgeable attacker.	635
586		636
587		637
588		638
589		639
590		640
591		641
592		642
593		643
594		644
595		645
596		646
597		647
598		648
599		649
600	The second mapper class we propose is based on utilizing contextualized information to maintain performance while obfuscating the original input text. This technique achieves higher privacy measures and has less impact on the downstream task, which makes it more applicable for cases where the downstream task is important. Nevertheless, opting for different weighted functions, such as ones based on a trained model, can further help improve both accuracy and privacy.	650
601		651
602		652
603		653
604		654
605		655
606		656
607		657
608		658
609		659
610	An inherent problem with existing privacy-preserving techniques is their inability to maintain linguistic properties such as grammar and readability (as seen in Table 2) that are crucial for the performance of the model. Therefore, an additional avenue we plan to explore is application of these and similar rules in differential privacy techniques. For instance, following the application of random perturbations to an embedding vector, instead of simply returning the nearest token to the perturbed vector, one could consider returning a token with similar syntactic attributes, such as part of speech, or verbs with similar causative meanings or stable subcategorization frames.	660
611		661
612		662
	Lastly, our experiments were limited to classification tasks in the English language. In future research, we intend to explore the effectiveness of these methods in generative tasks, across languages, and in multilingual settings.	663
	Limitations	664
	We demonstrated the privacy achieved by our methods empirically under one attacking scenario. Further comprehensive testing or mathematical proofs would enhance our understanding of the extent of privacy achieved.	665
	An additional limitation of our proposed mechanism is the unchanged sentence length. This imposes a privacy breach in which an author who prefers writing longer or shorter sentences can be re-identified even when introducing random perturbations. Hence, another avenue in this research is reducing the amount of tokens by introducing, for example, a stride parameter to the STENCIL family of mappers. This parameter will determine how often tokens will be output, thus reducing the amount of tokens.	666
	References	667
	Pathum Chamikara Mahawaga Arachchige, Peter Bertok, Ibrahim Khalil, Dongxi Liu, Seyit Camtepe, and Mohammed Atiquzzaman. 2019. Local differential privacy for deep learning. <i>IEEE Internet of Things Journal</i> , 7(7):5827–5842.	668
	Iyadh Ben Cheikh Larbi, Aljoscha Burchardt, and Roland Roller. 2023. Clinical text anonymization, its influence on downstream NLP tasks and the risk of re-identification . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop</i> , pages 105–111, Dubrovnik, Croatia. Association for Computational Linguistics.	669
	Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar	670

776		Richard Socher, Alex Perelygin, Jean Wu, Jason	833
777		Chuang, Christopher D. Manning, Andrew Ng, and	834
778		Christopher Potts. 2013. Recursive deep models for	835
779		semantic compositionality over a sentiment treebank .	836
		In <i>Proceedings of the 2013 Conference on Empirical</i>	837
780	Shervin Minaee, Nal Kalchbrenner, Erik Cambria, Nar-	<i>ical Methods in Natural Language Processing</i> , pages	838
781	jes Nikzad, Meysam Chenaghlu, and Jianfeng Gao.	1631–1642, Seattle, Washington, USA. Association	839
782	2021. Deep learning–based text classification: a com-	for Computational Linguistics.	840
783	prehensive review. <i>ACM computing surveys (CSUR)</i> ,		
784	54(3):1–40.		
785	Ahmadreza Mosallanezhad, Ghazaleh Beigi, and Huan	William Frederick Spitz. 1995. <i>High-order compact fi-</i>	841
786	Liu. 2019. Deep reinforcement learning-based text	<i>nite difference schemes for computational mechanics</i> .	842
787	anonymization against private-attribute inference . In	The University of Texas at Austin.	843
788	<i>Proceedings of the 2019 Conference on Empirical</i>		
789	<i>Methods in Natural Language Processing and the</i>	Alex Wang, Amanpreet Singh, Julian Michael, Felix	844
790	<i>9th International Joint Conference on Natural Lan-</i>	Hill, Omer Levy, and Samuel R. Bowman. 2019.	845
791	<i>guage Processing (EMNLP-IJCNLP)</i> , pages 2360–	GLUE: A multi-task benchmark and analysis plat-	846
792	2369, Hong Kong, China. Association for Computa-	form for natural language understanding. In <i>Proceed-</i>	847
793	tional Linguistics.	<i>ings of ICLR</i> .	848
794	OpenAI. 2021. Chatgpt . OpenAI Website. Accessed	Xin Zhou, Jinzhu Lu, Tao Gui, Ruotian Ma, Zichu Fei,	849
795	on 2023.	Yuran Wang, Yong Ding, Yibo Cheung, Qi Zhang,	850
796	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	and Xuanjing Huang. 2022. TextFusion: Privacy-	851
797	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	preserving pre-trained model inference via token fu-	852
798	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	sion . In <i>Proceedings of the 2022 Conference on</i>	853
799	2022. Training language models to follow instruc-	<i>Empirical Methods in Natural Language Processing</i> ,	854
800	tions with human feedback. <i>Advances in Neural</i>	pages 8360–8371, Abu Dhabi, United Arab Emirates.	855
801	<i>Information Processing Systems</i> , 35:27730–27744.	Association for Computational Linguistics.	856
802	Richard Plant, Dimitra Gkatzia, and Valerio Giuffrida.	Xin Zhou, Yi Lu, Ruotian Ma, Tao Gui, Yuran Wang,	857
803	2021. CAPE: Context-aware private embeddings	Yong Ding, Yibo Zhang, Qi Zhang, and Xuanjing	858
804	for private language learning . In <i>Proceedings of the</i>	Huang. 2023. TextObfuscator: Making pre-trained	859
805	<i>2021 Conference on Empirical Methods in Natural</i>	language model a privacy protector via obfuscating	860
806	<i>Language Processing</i> , pages 7970–7978, Online and	word representations . In <i>Findings of the Associa-</i>	861
807	Punta Cana, Dominican Republic. Association for	<i>tion for Computational Linguistics: ACL 2023</i> , pages	862
808	Computational Linguistics.	5459–5473, Toronto, Canada. Association for Com-	863
809	Hanieh Poostchi, Ehsan Zare Borzeshi, and Massimo	putational Linguistics.	864
810	Piccardi. 2018. BiLSTM-CRF for Persian named-		
811	entity recognition ArmanPersonNERCorpus: the first		
812	entity-annotated Persian dataset . In <i>Proceedings of</i>		
813	<i>the Eleventh International Conference on Language</i>		
814	<i>Resources and Evaluation (LREC 2018)</i> , Miyazaki,		
815	Japan. European Language Resources Association		
816	(ELRA).		
817	Chen Qu, Weize Kong, Liu Yang, Mingyang Zhang,		
818	Michael Bendersky, and Marc Najork. 2021. Natural		
819	language understanding with privacy-preserving bert.		
820	In <i>Proceedings of the 30th ACM International Con-</i>		
821	<i>ference on Information & Knowledge Management</i> ,		
822	pages 1488–1497.		
823	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,		
824	Dario Amodei, Ilya Sutskever, et al. 2019. Language		
825	models are unsupervised multitask learners. <i>OpenAI</i>		
826	<i>blog</i> , 1(8):9.		
827	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine		
828	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,		
829	Wei Li, and Peter J Liu. 2020. Exploring the limits		
830	of transfer learning with a unified text-to-text trans-		
831	former. <i>The Journal of Machine Learning Research</i> ,		
832	21(1):5485–5551.		