

# User-Entity Differential Privacy in Learning Natural Language Models

Phung Lai, NhatHai Phan\*  
New Jersey Institute of Technology, USA  
{tl353, phan}@njit.edu

Tong Sun, Rajiv Jain, Franck Deroncourt,  
Jiuxiang Gu, Nikolaos Barmpalios  
Adobe Systems Inc., USA  
{tsun, rajijain, deronco, jigu, barmpali}@adobe.com

**Abstract**—In this paper, we introduce a novel concept of user-entity differential privacy (UeDP) to provide formal privacy protection simultaneously to both sensitive entities in textual data and data owners in learning natural language models (NLMs). To preserve UeDP, we developed a novel algorithm, called UeDP-Alg, optimizing the trade-off between privacy loss and model utility with a tight sensitivity bound derived from seamlessly combining user and sensitive entity sampling processes. An extensive theoretical analysis and evaluation show that our UeDP-Alg outperforms baseline approaches in model utility under the same privacy budget consumption on several NLM tasks, using benchmark datasets.

**Index Terms**—Differential privacy, natural language models, entities, user identity

## I. INTRODUCTION

Despite remarkable performance in many applications, natural language models (NLMs), such as GPT models [1, 2, 3], are vulnerable to privacy attacks because of such attacks' capacity to memorize unique patterns in training data [4]. Recent data training extraction attacks [5] illustrate that sensitive entities, such as a person's name, email address, phone number, physical address, etc., can be accurately extracted from NLM parameters. These sensitive entities and the language data memorized in NLMs may identify a data owner - explicitly by name or implicitly, e.g., via a rare or unique phrase - and link that data owner to extracted sensitive entities.

Our main goal is to provide a rigorous guarantee that a trained NLM protects the privacy of sensitive entities in the training data and the participation information (membership) of the data owners in learning the model while maintaining high model utility. The simple solution of anonymizing (including removing/de-identifying) sensitive entities is insufficient; since the anonymized entities can be matched with non-anonymized data records in another dataset [6]. Also, the model utility can be notably affected, as shown in our experimental study. While cryptographic approaches can be applied to protect privacy, they introduce computation and resource overhead [7]. Therefore, we proposed to use differential privacy [8], one of the adequate solutions, given its formal protection without undue sacrifice in computation efficiency and model utility.

\*Corresponding author

Differential privacy (DP) provides rigorous privacy protection as a probabilistic term, limiting the knowledge about a data record an ML model can leak while learning features of the whole training set. DP-preserving mechanisms have been investigated and applied in real-world [9, 10, 11], including image processing [12], healthcare data [13], financial records [14], social media [15], and NLMs [16, 17, 18, 19].

However, existing DP protection levels, including sample-level DP [6, 9, 20, 21], user-level DP [16, 22], element-level DP [23], and local (feature-level) DP [17, 18, 24, 25], do not provide the privacy protection level demanded to solve our problem. Given training data: 1) Sample-level DP protects the privacy of a single sample; 2) User-level DP protects privacy of a single data owner, also called a single user, who may contribute one or more data samples; 3) Element-level DP partitions data owners' contribution to the training data into sensitive elements, e.g., a curse word, which will be protected. Element-level DP does not provide privacy protection to data owners; and 4) Local (feature-level) DP protects true values of a data sample from being inferred. Recently, [18] proposed local DP-preserving approaches for text embedding extraction under (word-level) local DP (**Eq. 2**). However, the privacy budget in [18] is accumulated over the dimensions of embedding, resulting in an impractical (loose) privacy guarantee (Appendix D in our supplemental document<sup>1</sup>).

Therefore, there is a demand for a new level of DP to protect privacy simultaneously for both sensitive entities in the training data and the participation information of data owners in learning NLMs. Motivated by this, we structure our paper around the following significant contributions.

- We propose a novel notion of user-entity adjacent databases (**Definition 2**), leading to formal guarantees of user-entity privacy rather than privacy for a single user or a single sensitive entity.

- To preserve UeDP, we introduce a novel algorithm, called **UeDP-Alg**, which leverages the recipe of DP-FEDAVG [16] to protect both sensitive entities and user membership under DP via the moments accountant [9]. Moments accountant was first developed to preserve DP in stochastic gradient descent (SGD) for sample-level privacy. Our federated averaging approach groups multiple SGD updates computed from a two-level random sampling process, including a random sample of

users and a random sample of sensitive entities. That enables large-step model updates and optimizes the trade-off between privacy loss and the model utility through a tight noise scale bound (**Lemma 1** and **Theorem 1**).

• Through theoretical analysis and rigorous experiments conducted on benchmark datasets, we show that our UeDP-Alg outperforms baseline approaches in terms of model utility on fundamental tasks, i.e., next word prediction and text classification, under the same privacy budget consumption. Our code is available<sup>2</sup>.

## II. BACKGROUND

In this section, we revisit NLM tasks, privacy risk, and DP. For the sake of clarity, let us focus on the next word prediction, and we will extend it to text classification in Section VI. A list of sensitive entity categories is summarized in Table I.

a) *Next Word Prediction*: Let  $D$  be a private training data containing  $U$  users (data owners) and a set of sensitive entities  $E$ . Each user  $u \in U$  consists of  $n_u$  sentences. Given a vocabulary  $\mathcal{V}$ , each sentence is a sequence of words, presented as  $x = x_1 x_2 \dots x_{m_u}$ , where  $x_i \in \mathcal{V}$ , ( $i \in [1, m_u]$ ) is a word in  $x$  and  $m_u$  is the length of  $x$ . In next word prediction, the first  $j$  words in  $x$ , i.e.,  $x_1, x_2, \dots, x_j$  ( $\forall j < m_u$ ), are used to predict the next word  $x_{j+1}$ . Here,  $x_{j+1}$  can be considered as a label in the next word prediction task. Perplexity  $PP = 2^{-\sum_{x \in D} p(x) \log_2 p(x)}$  is a measurement of how well a model predicts a sentence and is often used to evaluate language models, where  $p(x)$  is a probability to predict the next word  $x_{j+1}$  in  $x$  [26]. A lower perplexity indicates a better model.

b) *Sensitive Entities and Sentences*: Each sensitive entity  $e \in E$  consists of a word or consecutive words that must be protected. For instance, personal identifiable information (PII) related to an identifiable person, such as person names, locations, and phone numbers, can be considered sensitive entities. If a sentence  $x$  consists of a sensitive entity  $e$ ,  $x$  is considered as a sensitive sentence; otherwise,  $x$  is a non-sensitive sentence.

For instance, in Fig. 1, “David Johnson,” “Maine,” “September 18,” and “Main Hospital” are considered sensitive entities, correspondingly categorized into PII, geopolitical entities (GPE) (i.e., countries, cities, and states), time, and organization names. The first and second sentences consisting of the sensitive entities are considered sensitive sentences. Meanwhile, the third and fourth sentences are non-sensitive since they do not contain any sensitive entities.

c) *Privacy Threat Models*: It is well-known that trained ML model parameters can disclose information about training data [5, 27], especially in NLMs [5, 16]. Given a data sample and model parameters, by using a membership inference attack [28, 29, 30], adversaries can infer whether the training used the sample or not. In NLMs, adversaries can accurately recover individual training examples, such as full names, email addresses, and phone numbers of individuals, using training data extracting attacks [5]. Accessing to these can lead to severe privacy breaches.

TABLE I: Description of sensitive entity categories.

Type	Description
Person	Person, i.e., people, including fictional
Loc	Location, i.e., non-GPE locations, mountain ranges, bodies of water
Org	Organization, i.e., companies, agencies, institutions, etc.
Misc	Miscellaneous, i.e., entities that do not belong to the person, location, and organization in CONLL-2003
GPE	Geopolitical entity, i.e., countries, cities, states
PII	Personal identification information, i.e., name, location, phone, etc.
Date	Absolute or relative dates or periods
NoRP	Nationalities or religious or political groups
Fac	Buildings, airports, highways, bridges, etc.
Product	Objects, vehicles, foods, etc. (Not services.)
Event	Named hurricanes, battles, wars, sports events, etc.
Law	Named documents made into laws
Language	Any named language
Work of art	Titles of books, songs, etc.
Time	Times smaller than a day
Percent	Percentage, including “%”
Money	Monetary values, including unit
Quantity	Measurements, as of weight or distance
Ordinal	“First”, “second”, etc.
Cardinal	Numerals that do not fall under another type

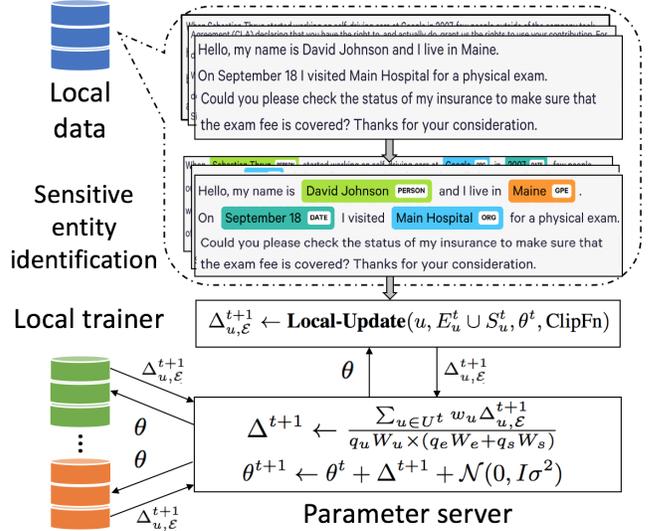


Fig. 1: User-Entity DP. Data from users is processed to identify sensitive entities, before being trained by local trainers. Bounded gradients from local trainers are aggregated at a server with additive noise. Updated model is sent back to local trainers for next rounds.

## III. DIFFERENT LEVELS OF DP

To avoid these privacy risks, DP guarantees restriction on the adversaries in what they can learn from the training data given the model parameters by ensuring similar model outcomes with and without any single training sample. Let us revisit the definition of DP, as follows:

**Definition 1.**  $(\epsilon, \delta)$ -DP [8]. A randomized algorithm  $\mathcal{A}$  fulfills  $(\epsilon, \delta)$ -DP, if for any two adjacent datasets  $D$  and  $D'$  differing by at most one sample, and for all outcomes  $\mathcal{O} \subseteq \text{Range}(\mathcal{A})$ :

$$\Pr[\mathcal{A}(D) = \mathcal{O}] \leq e^\epsilon \Pr[\mathcal{A}(D') = \mathcal{O}] + \delta \quad (1)$$

with a privacy budget  $\epsilon$  and a broken probability  $\delta$ .

<sup>2</sup><https://github.com/PhungLai728/UeDP>

The privacy budget  $\epsilon$  controls the amount by which the distributions induced by  $D$  and  $D'$  may differ. A smaller  $\epsilon$  enforces a stronger privacy guarantee. The broken probability  $\delta$  means the highly unlikely “bad” events, in which an adversary can infer whether a particular data sample belongs to the training data, happen with the probability  $\leq \delta$ .

There are different levels of DP protection in literature categorized into four research lines, including sample-level DP, user-level DP, element-level DP, and local (feature-level) DP. They are different from our goal since we focus on providing simultaneous protections to data owners and sensitive entities in textual data. Let us revisit these DP levels and distinguish them with our goal.

*a) Sample-level DP:* Traditional DP mechanisms [6, 20, 31] ensure DP at the sample-level, in which adjacent datasets  $D$  and  $D'$  are different from at most a single training sample. Sample-level DP does not protect privacy for users. That is different from our goal. We aim at protecting privacy for users and sensitive entities, which are different from data samples.

*b) User-level DP:* To protect privacy for users, who may contribute more than one training sample, rather than a single sample, [16] proposed a user-level DP, in which neighboring databases  $D$  and  $D'$  are defined to be different from all of the samples associated with an arbitrary user in the training set. Several works follow this direction [22, 32]. User-level DP differs from our goal, since it does not provide privacy protection for sensitive entities in the training set.

*c) Element-level DP:* [23] introduce element-level DP, in which users are partitioned based on sensitive elements, which are protected in a way that an adversary cannot infer whether a user has a sensitive element in her/his data, e.g., if a user has ever sent a curse word in his/her messages or not. Similar to sample-level DP, element-level DP is different from our goal, since it does not provide DP protection for users.

*d) Local (feature-level) DP:* [18] proposed a notion of word-level local DP for a sentence’s embedding features, in which two adjacent sentences  $x$  and  $x'$  are different at most one word:

$$Pr[\mathcal{A}(f(x)) = \mathcal{O}] \leq e^\epsilon Pr[\mathcal{A}(f(x')) = \mathcal{O}] \quad (2)$$

where  $f(x)$  extracts embedded features of  $x$  and  $\mathcal{A}$  is a randomized algorithm, such as a Laplace mechanism [6]. In a similar effort, [17] applied a randomized response mechanism [24, 33, 34] on top of binary encoding of embedded features’ real values to achieve local DP feature embedding. The approaches proposed in [17, 18] are different from our goal, since they do not offer either user-level DP or word-level DP.

#### IV. USER-ENTITY DIFFERENTIAL PRIVACY

In this section, we focus on answering the question: “*Could we protect sensitive entities and user membership simultaneously by leveraging existing levels of DP and how?*” Based upon that, we propose our user-entity DP notion.

##### A. Sensitive Entities and User Membership

To protect sensitive entities and user membership, a potential approach is to decouple them into separated protection levels offering by existing DP notions. However, this approach has limitations as discussed next.

Let us consider a sentence consisting of one or more than one sensitive entities. We can leverage sample-level DP to protect the sentence, i.e., each sentence could be a sample, covering all the sensitive entities under DP. If each user has only one sentence, then this approach can also protect the user membership. In practice, one user may contribute many sentences to the training data. To address this issue, we can utilize group privacy [6] resulting in an amplification of the privacy budget proportional to the number of sentences a user may have in the training data.

Instead of group privacy, another potential solution is applying user-level DP on top of the sample-level DP to protect both sentence and user membership. In the sample-level DP, we can clip and inject Gaussian noise into the gradient derived from each sentence [9]. Meanwhile, in the user-level DP, an additional Gaussian noise is injected into the aggregation of gradients, each of which derived from a single user [16]. Although this combination of sample - user levels can cover both sensitive entities and user membership under DP protection, it has disadvantages. First, some sentences are sensitive and other sentences are not. Protecting all (sensitive and non-sensitive) sentences or removing all the sensitive sentences from the training data may cause significant model utility degradation. Second, different sentences may consist of different types and numbers of sensitive entities. Under the same sampling probability for training as in [9] for sample-level DP, these sentences expose different privacy risks to user identity and sensitive entities.

To address these issues, instead of the sentence level, one can work at the word level by extracting embedded features for every words in the training data. Embedded features of sensitive entities are randomized by local DP-preserving mechanisms [34]. The randomized embedded features are aggregated with embedded features of non-sensitive words to train NLMs. Then, user-level DP can be applied to clip gradients derived from each user’s data with adding Gaussian noise into the aggregation of these gradients. However, this approach suffers from a remarkable model utility degradation. Local DP provides rigorous privacy protection but it comes with a cost in terms of utility [35]. Then, adding the user-level DP adversely affects the utility.

The root cause of these limitations is that the combination of sentence-level DP and user-level DP notions does not capture the correlation between sensitive entities and user membership in unifying notion of DP. Meanwhile, working with word-level embedded features under local DP introduces expensive model utility costs. Therefore, there is a demand for a unifying notion of DP and an optimal approach to protect both sensitive entities and user membership in training NLMs.

## B. UeDP Definition

To preserve privacy for both users and sensitive entities in NLMs, we propose a new definition of user-entity adjacent databases, as follows: Two databases  $D$  and  $D'$  are user-entity adjacent if they differ in a single user and a single sensitive entity; that is, one user  $u'$  and one sensitive entity  $e'$  are present in one database (i.e.,  $D'$ ) and are absent in the other (i.e.,  $D$ ). Together with the absence of all sentences from the user  $u'$  in  $D$ , all sentences (across users) consisting of the sensitive entity  $e'$  are also absent in  $D$ . This is because one user can have multiple sentences, and one sensitive entity can exist in multiple sentences for training. The definition of our user-entity adjacent databases is presented as follows:

**Definition 2.** *User-Entity Adjacent Databases.* Two databases  $D$  and  $D'$  are called user-entity adjacent if:  $\|U - U'\|_1 \leq 1$  and  $\|E - E'\|_1 \leq 1$ , where  $U$  and  $E$  are the sets of users and sensitive entities in  $D$ , and  $U'$  and  $E'$  are the sets of users and sensitive entities in  $D'$ .

Given the user-entity adjacent databases, we present our UeDP in the following definition.

**Definition 3.**  $(\epsilon, \delta)$ -UeDP. A randomized algorithm  $\mathcal{A}$  is  $(\epsilon, \delta)$ -UeDP if for all outcomes  $\mathcal{O} \subseteq \text{Range}(\mathcal{A})$  and for all user-entity adjacent databases  $D$  and  $D'$ , we have:

$$\Pr[\mathcal{A}(D) = \mathcal{O}] \leq e^\epsilon \Pr[\mathcal{A}(D') = \mathcal{O}] + \delta \quad (3)$$

with a privacy budget  $\epsilon$  and a broken probability  $\delta$ .

## V. PRESERVING UEDP IN NLMs

UeDP provides rigorous privacy protection to both users and sensitive entities; however, the practicability of UeDP preservation depends on the reliability of sensitive entity detection from the training text data. In practice, misidentifying sensitive entities can introduce extra privacy risks. In addition to addressing this challenge, we focus on bounding the sensitivity of an NLM under UeDP and addressing the trade-off between privacy loss and model utility.

### A. Misidentifying Sensitive Entities

Identifying all the sensitive entities typically requires intensive manual efforts [36]. We are aware of this issue in real-world applications. Fortunately, there are several ways to automatically identify sensitive entities in textual data, such as: 1) Using Named Entity Recognition (NER) [37, 38]; and 2) Using publicly available toolkits for detecting named entities or PII in text, e.g., spaCy [39], Stanza [40], and Microsoft Presidio<sup>3</sup>. These approaches and toolkits are user-friendly and reliable to reduce manual efforts in identifying sensitive entities and information. We found that the results from spaCy cover over 94% of sensitive information identified by Amazon Mechanical Turk (AMT) workers in a diverse set of datasets used in our experiments. More information about identifying sensitive entities is available in Appx. A.

Although effective, the small error rate (i.e.,  $\approx 6\%$ ) from these techniques introduces a certain level of privacy risk, that means, some sensitive entities may be misidentified to be non-sensitive, and vice-versa. Classifying non-sensitive entities to be sensitive entities does not incur any extra privacy risk. Meanwhile, classifying one (or more than one) sensitive entity to be non-sensitive in a sentence introduces two issues, as follows: **(1)** There may be sensitive sentences misidentified to become non-sensitive sentences. In other words, given a set of non-sensitive sentences detected by NER tools, we do not know which sentence is truly non-sensitive; and **(2)** Given a sensitive sentence  $x$ , some sensitive entities in  $x$  may not be identified by NER tools. Preserving UeDP in NLMs by directly using the results of NER tools will expose these misidentifying sensitive sentences and entities unprotected.

### B. Preserving UeDP

To address the problem of sensitive entity misidentification in preserving UeDP, our key idea is:

**(1)** Extending UeDP by considering each sentence, identified to be non-sensitive using NER tools [39, 40], in the private training dataset as a single type of sensitive entity. We denote this extended set of sensitive entities as  $S$ . The private dataset  $D$  now consists of  $U$  users and a (sufficient) set of sensitive entities  $E \cup S$  that will be protected.

**(2)** Upon forming the sufficient set of sensitive entities, we propose a two-step sampling approach to strictly preserve UeDP in NLMs. In our approach, at a training round  $t$ , we sample a set of users from  $U$  and a set of sensitive entities from  $E \cup S$ . We use sentences in the training data of the sampled users consisting of the sampled sensitive entities to train NLMs. In this sampling approach: (i) If a sensitive sentence  $x$  is not sampled for training, i.e., due to the fact that some sensitive entities in  $x$  are not identified by NER tools,  $x$  is not used for training at the round  $t$ ; thus avoiding privacy risks exposed by  $x$ ; and (ii) If the sensitive sentence  $x$  is sampled for training, then the sensitive entities in  $x$ , which are not identified by NER tools, are protected since  $x$  is protected under DP.

By covering all possible cases of sensitive entity misidentification, we strictly preserve UeDP without having additional privacy risks. The pseudo-code of our algorithm is in Alg. 1.

At each iteration  $t$ , we randomly sample  $U^t$  users from  $U$ ,  $E^t$  detected sensitive entities from  $E$ , and  $S^t$  extended sensitive entities from  $S$ , with sampling rates  $q_u$ ,  $q_e$ , and  $q_s$ , respectively (Lines 8 and 10). Then, we use all sensitive sentences in  $E_u^t \cup S_u^t$  consisting of the sensitive entities in  $E^t$  and  $S^t$  belonging to the selected users in  $U^t$  for training. Like [16], we leverage the basic federated learning setting in [41] to compute gradients of model parameters for a particular user, denoted as  $\Delta_{u,\mathcal{E}}^{t+1}$  (Line 11). Here, we clip the per-user gradients so that its  $l_2$ -norm is bounded by a predefined gradient clipping bound  $\beta$  (Lines 20 - 29). Next, a weighted-average estimator  $f_{\mathcal{E}^+}$  is employed to compute the average gradient  $\Delta^{t+1}$  using the clipped gradients  $\Delta_{u,\mathcal{E}}^{t+1}$  gathered from all the selected users (Line 13). Finally, we add random

<sup>3</sup><https://microsoft.github.io/presidio/>

Gaussian noise  $\mathcal{N}(0, I\sigma^2)$  to the model update (Line 15). During the training, the moments accountant  $\mathcal{M}$  is used to compute the  $T$  training steps' privacy budget consumption (Lines 16 - 18).

To tighten the sensitivity bound, our weighted-average estimator  $f_{\mathcal{E}^+}$  (Line 13) is as follows:

$$f_{\mathcal{E}^+}(U^t, E^t) = \frac{\sum_{u \in U^t} w_u \Delta_{u, \mathcal{E}^+}^{t+1}}{q_u W_u (q_e W_e + q_s W_s)} \quad (4)$$

where  $\Delta_{u, \mathcal{E}^+}^{t+1} = \sum_{e \in E_u^t} w_e \Delta_{u, e} + \sum_{s \in S_u^t} w_s \Delta_{u, s}$ , and  $w_u, w_e$ , and  $w_s \in [0, 1]$  are weights associated with a user  $u$ , a detected sensitive entity  $e$ , and an extended sensitive entity  $s$ .

These weights capture the influence of a user and sensitive entities to the model outcome.  $\Delta_{u, e}$  and  $\Delta_{u, s}$  are the parameter gradients computed using the sensitive entities  $e \in E$  and  $s \in S$ . In addition,  $W_u = \sum_{u \in U} w_u$ ,  $W_e = \sum_{e \in E} w_e$ , and  $W_s = \sum_{s \in S} w_s$ .

Since  $\mathbb{E}[\sum_{e \in E_u^t} w_e + \sum_{s \in S_u^t} w_s] = q_e W_e + q_s W_s$ , the estimator  $f_{\mathcal{E}^+}$  is unbiased. The sensitivity of the estimator  $\mathbb{S}(f_{\mathcal{E}^+})$  is computed as:  $\mathbb{S}(f_{\mathcal{E}^+}) = \max_{u', e'} \|f_{\mathcal{E}^+}(\{U^t \cup u', (E^t \cup S^t) \cup e'\}) - f_{\mathcal{E}^+}(\{U^t, E^t \cup S^t\})\|_2$ .  $\mathbb{S}(f_{\mathcal{E}^+})$  is bounded in the following lemma.

**Lemma 1.** *If for all users  $u$  we have  $\|\Delta_{u, \mathcal{E}^+}^{t+1}\|_2 \leq \beta$ , then  $\mathbb{S}(f_{\mathcal{E}^+}) \leq \frac{(q_u |U| + 1) \max(w_u) \beta}{q_u W_u (q_e W_e + q_s W_s)}$ .*

*Proof.* If for all users  $\|\Delta_{u, \mathcal{E}^+}^{t+1}\|_2 \leq \beta$ , then

$$\begin{aligned} \mathbb{S}(f_{\mathcal{E}^+}) &= \frac{\sum_{u \in U^t \cup u'} w_u \left( \sum_{e \in E_u^t} w_e (\sum_{s \text{ consists of } e} \Delta_{u, s}) \right)}{q_u W_u (q_e W_e + q_s W_s)} \\ &+ \frac{\sum_{u \in U^t \cup u'} w_u \left( \sum_{s \in S_u^t} w_s \Delta_{u, s} \right)}{q_u W_u (q_e W_e + q_s W_s)} \\ &+ \frac{\sum_{u \in U^t \cup u'} w_u \left[ w_{e'} (\sum_{s \text{ consists of } e'} \Delta_{u, s}) \right]}{q_u W_u (q_e W_e + q_s W_s)} \\ &- \frac{\sum_{u \in U^t} w_u \left( \sum_{e \in E_u^t} w_e (\sum_{s \text{ consists of } e} \Delta_{u, s}) \right)}{q_u W_u (q_e W_e + q_s W_s)} \\ &- \frac{\sum_{u \in U^t} w_u \left( \sum_{s \in S_u^t} w_s \Delta_{u, s} \right)}{q_u W_u (q_e W_e + q_s W_s)} \\ &\leq \frac{\sum_{u \in U^t \cup u'} [(w_u) \beta]}{q_u W_u (q_e W_e + q_s W_s)} \leq \frac{(q_u |U| + 1) \max(w_u) \beta}{q_u W_u (q_e W_e + q_s W_s)} \quad (5) \end{aligned}$$

Consequently, Lemma 1 holds.  $\square$

Once the sensitivity of the estimator  $f_{\mathcal{E}^+}$  is bounded, we can add Gaussian noise scaled to the sensitivity  $\mathbb{S}(f_{\mathcal{E}^+})$  to obtain a privacy guarantee. By applying Lemma 1, the noise scale  $\sigma$  becomes:

$$\sigma = z \mathbb{S}(f_{\mathcal{E}^+}) = \frac{z(q_u |U| + 1) \max(w_u) \beta}{q_u W_u (q_e W_e + q_s W_s)} \quad (6)$$

The noise scale  $\sigma$  in Eq. 6 is tighter than the noise scale in existing works [16, 22] proportional to the number of sensitive entities used in the training process (i.e.,  $q_e W_e + q_s W_s$ ).

- 1: **Input:** Dataset  $D$ , set of sensitive entities  $E$ , extended set of sensitive entities  $S$ , sampling rates  $q_u, q_e$ , and  $q_s$ , clipping bound  $\beta$ , a hyper-parameter  $z$ , and number of iterations  $T$
- 2: Initialize model  $\theta^0$  and moments accountant  $\mathcal{M}$
- 3:  $w_u \leftarrow \min(\frac{n_u}{\hat{w}_u}, 1)$  for all users  $u$
- 4:  $w_e \leftarrow \min(\frac{n_e}{\hat{w}_e}, 1)$  for all sensitive entities in  $E$
- 5:  $w_s \leftarrow \min(\frac{n_s}{\hat{w}_s}, 1)$  for all extended sensitive entities in  $S$  where  $n_u, n_e$ , and  $n_s$  are the number of sentences in user  $u$ , the number of sentences containing sensitive entities  $e \in E$ , the number of sensitive entities in  $S$ , and  $\hat{w}_u, \hat{w}_e$ , and  $\hat{w}_s$  are per-user sentence cap, per-entity sentence cap, and per-entity entity cap.
- 6:  $W_u \leftarrow \sum_{u \in U} w_u, W_e \leftarrow \sum_{e \in E} w_e, W_s \leftarrow \sum_{s \in S} w_s$
- 7: **for**  $t \in T$  **do**
- 8:  $U^t \leftarrow$  sample users with probability  $q_u$
- 9: **for** each user  $u \in U^t$  **do**
- 10:  $E_u^t \cup S_u^t \leftarrow$  sensitive entities (belonging to the user  $u$ ) consisting of sensitive entities  $E^t$  sampled from  $E$  with probability  $q_e$  and extended sensitive entities  $S^t$  sampled from  $S$  with probability  $q_s$
- 11:  $\Delta_{u, \mathcal{E}^+}^{t+1} \leftarrow$  **Local-Update**( $u, E_u^t \cup S_u^t, \theta^t, \text{ClipFn}$ )
- 12: **end for**
- 13:  $\Delta^{t+1} \leftarrow \frac{\sum_{u \in U^t} w_u \Delta_{u, \mathcal{E}^+}^{t+1}}{q_u W_u (q_e W_e + q_s W_s)}$
- 14:  $\sigma \leftarrow \frac{z(q_u |U| + 1) \max(w_u) \beta}{q_u W_u (q_e W_e + q_s W_s)}$
- 15:  $\theta^{t+1} \leftarrow \theta^t + \Delta^{t+1} + \mathcal{N}(0, I\sigma^2)$
- 16:  $\mathcal{M}.\text{accum\_priv\_spending}(z)$
- 17: **end for**
- 18: **print**  $\mathcal{M}.\text{get\_priv\_spent}()$
- 19: **Output:**  $(\epsilon, \delta)$ -UeDP  $\theta, \mathcal{M}$
- 20: **Local-Update**( $u, E_u^t \cup S_u^t, \theta^0, \text{ClipFn}$ ):
- 21:  $\theta \leftarrow \theta^0$
- 22:  $\mathcal{B} \leftarrow$   $u$ 's data split into size  $B$  batches
- 23: **for** batch  $b \in \mathcal{B}$  **do**
- 24:  $\forall e \in E_u^t : \Delta_{u, e} \leftarrow \sum_{\text{sentence } s \text{ (} \in b \text{) consists of } e} \nabla l(\theta, s)$
- 25:  $\forall s \in S_u^t \cap b : \Delta_{u, s} \leftarrow \nabla l(\theta, s)$
- 26:  $\Delta_{u, \mathcal{E}^+} \leftarrow \sum_{e \in E_u^t} w_e \Delta_{u, e} + \sum_{s \in S_u^t} w_s \Delta_{u, s}$
- 27:  $\theta \leftarrow \theta^0 - \eta \Delta_{u, \mathcal{E}^+}$
- 28: **end for**
- 29: **return**  $\text{ClipFn}(\theta - \theta^0, \beta)$
- 30: **ClipFn**( $\Delta, \beta$ ): **return**  $\pi(\Delta, \beta) \leftarrow \Delta \cdot \min\left(1, \frac{\beta}{\|\Delta\|}\right)$

**Algorithm 1:** UeDP-Alg

Therefore, we can inject less noise into our model under the same privacy budget while improving our model utility.

In extreme cases, that is also true: (1)  $E$  is empty, which means there are no detected sensitive entities. Given a fixed set of training data, while  $E$  is empty,  $S$  becomes larger (i.e., covering the whole dataset), resulting in a larger value of  $W_s$ . Therefore, we obtain a larger value of  $q_s W_s$  (with a pre-defined  $q_s$ ), enabling us to reduce the noise scale under the same UeDP guarantee. That is an advantage compared with the naive approach that only uses detected sensitive entities  $E$  in the training process (i.e., ignoring the term  $q_s W_s$  in Eq. 6). If  $E$  is empty, the naive approach will have no sentences for training; and (2)  $S$  is empty, that is, every sentence in the data consists of at least one detected sensitive entity  $e \in E$ . Similarly, given a fixed set of training data, if  $S$  is empty, then  $E$  and  $W_e$  become larger. It enables us to obtain a larger value of  $q_e W_e$  (with a pre-defined  $q_e$ ), which results in smaller noise

scale while maintaining the high model utility.

*UeDP Guarantee.* Given the bounded sensitivity of the estimator, the moments accountant  $\mathcal{M}$  [9] is used to get a tight bound on the total UeDP privacy consumption of  $T$  steps of the Gaussian mechanism with the noise  $\mathcal{N}(0, I\sigma^2)$  (Line 15).

**Theorem 1.** *For the estimator  $f_{\mathcal{E}^+}$ , the moments accountant of the sampled Gaussian mechanism correctly computes UeDP privacy loss with the scale  $z = \sigma/\mathbb{S}(f_{\mathcal{E}^+})$  for  $T$  training steps.*

*Proof.* At each step, users, detected sensitive entities in  $E$ , and extended sensitive entities in  $S$  are selected randomly with probabilities  $q_u$ ,  $q_e$ , and  $q_s$ , respectively. For  $f_{\mathcal{E}^+}$ , if the  $l_2$ -norm of each user’s gradient update, using the sampled sensitive entities in  $E_u^t \cup S_u^t$ , is bounded by  $\mathbb{S}(f_{\mathcal{E}^+})$ , then the moments accountant can be bounded by that of the sampled Gaussian mechanism with sensitivity 1, the scale  $z = \sigma/\mathbb{S}(f_{\mathcal{E}^+})$ , and sampling probabilities  $q_u$ ,  $q_e$ , and  $q_s$ . Thus, we can apply the composability as in Theorem 2.1 [9] to correctly compute the UeDP privacy loss with the scale  $z = \sigma/\mathbb{S}(f_{\mathcal{E}^+})$  for  $T$  steps.  $\square$

## VI. EXPERIMENTAL RESULTS

We conducted an extensive experiment, both in theory and on benchmark datasets, to shed light on understanding 1) the integrity of sensitive entity identification, 2) the interplay among the UeDP privacy budget  $(\epsilon, \delta)$ , different types of sensitive entities (i.e., organization, location, PII, and miscellaneous entities), and model utility, and 3) whether considering the extended set of sensitive entities  $S$  will improve our model utility under the same UeDP protection.

*a) Baseline Approaches:* We evaluate our **UeDP-Alg** in comparison with both noiseless and privacy-preserving mechanisms (either user level or entity level), including: **(1) User-level DP** [16], which is the state-of-the-art DP-preserving model closely related to our work; **(2) De-Identification** [42], which is considered as a strong baseline to protect privacy for sensitive entities. Although sensitive entities are masked to hide them in the training process, De-Identification does not offer formal privacy protection to either the data owners or sensitive entities; and **(3) A Noiseless** model, which is a language model trained without any privacy-preserving mechanisms. In addition, we consider the naive approach, which is a variation of our algorithm, called **UeDP-Alg  $f_{\mathcal{E}}$** . As a baseline, the estimator  $f_{\mathcal{E}}$  is computed without taking the extended set of sensitive entities  $S$  into account (Appx. B, Supplementary<sup>4</sup>). This is further used to comprehensively evaluate our proposed approach. In our experiment, our algorithms and baselines, i.e., UeDP-Alg, User-level DP, and De-Identification, are applied on the noiseless model in the training process. As in the literature review [18, 23], there are no other appropriate DP-preserving baselines for UeDP protection.

*b) Evaluation Tasks and Metrics:* Our experiment considers two tasks: **(1)** next word prediction and **(2)** text classification. For the next word prediction, we employ the widely used perplexity [43]. The smaller perplexity is, the better

model is. For the text classification, we use the test error rate as in earlier work [44]. Test error rate implies prediction error on a test set, so it is 1 - the test set’s accuracy. The lower the test error rate is, the better model is.

*c) Data and Model Configuration:* For the reproducibility sake, all details about our datasets and data processing are included in Appx. C (Supplementary<sup>4</sup>). We carried out our experiment on three textual datasets, including the CONLL-2003 news dataset [37], AG’s corpus of news articles<sup>4</sup>, and our collected Security and Exchange Commission (SEC) financial contract dataset. The data breakdown is in Table II.

For the next word prediction, we employ a GPT-2 model [2], which is one of the state-of-art text generation models. To make the work easily reproducible, we use a version of the pretrained GPT-2 that has 12-layer, 768-hidden, 12-heads, 117M parameters, and then fine-tune with the aforementioned datasets as our Noiseless GPT-2 model. For the text classification, we fine-tune a Noiseless BERT (i.e., BERT-Base-Uncased<sup>5</sup>) pre-trained model [45] that has 12-layer, 768-hidden, 12-heads, and 110M parameters with an additional softmax function on top of the BERT model. Adam optimizer is used with the learning rate is  $10^{-5}$ . Gradient clipping bound  $\beta = 0.1$  and the scale  $z = 2.5$ . The sampling rates for users, detected sensitive entities, and extended sensitive entities  $q_u$ ,  $q_e$ , and  $q_s$  are 0.05, 0.5, and 1.0.

To test the effectiveness and adaptability of our mechanism across models, we also conducted experiments with an AWD-LSTM model [46], which has a much fewer parameters compared with GPT-2 and BERT. In AWD-LSTM model, we use a three-layer LSTM model with 1,150 units in the hidden layer and an embedding input layer of size 100. Embedding weights are uniformly initialized in the interval  $[-0.1, 0.1]$  with dimension  $d = 100$  and other weights are initialized between  $[-\frac{1}{\sqrt{H}}, \frac{1}{\sqrt{H}}]$ , where  $H$  is the size of all hidden layers. The values used for dropout on the embedding layer, the LSTM hidden-to-hidden matrix, and the final LSTM layer’s output are 0.1, 0.3, and 0.5, respectively. Gradient clipping bound  $\beta = 0.1$  and the scale  $z = 2$ . The sampling rates  $q_u$ ,  $q_e$ , and  $q_s$  are 0.05, 0.5, and 1.0 (note that  $q_s$  is 0.6 in the text classification task). SGD optimizer is used. In the text classification with the AG dataset, a softmax layer is applied on top of the AWD-LSTM with the output dimension is 4, corresponding to four classes in the AG dataset. The same sets of sensitive entity categories are used for all models in the next word prediction and the text classification tasks.

*d) Evaluation Results:* To answer our evaluation questions, we conducted the following experiments: **(1)** examining how the sensitive entities detected by the entity recognition spaCy [39] covers the sensitive information clarified by AMT workers, **(2)** comparing estimators  $f_{\mathcal{E}}$ ,  $f_{\mathcal{E}^+}$ , and User-level DP; **(3)** investigating the interplay between privacy budget and model utility; **(4)** studying the impacts of different sensitive

<sup>4</sup>[http://groups.di.unipi.it/~gulli/AG\\_corpus\\_of\\_news\\_articles.html](http://groups.di.unipi.it/~gulli/AG_corpus_of_news_articles.html)

<sup>5</sup>[https://huggingface.co/transformers/pretrained\\_models.html](https://huggingface.co/transformers/pretrained_models.html)

TABLE II: Breakdown of CONLL-2003, AG, and SEC datasets.

Dataset	$\mathcal{V}$	# of sentences	# of users	# of sensitive sentences				
				Org	Loc	Person	Misc	All
CONLL-2003	8,882	14,040	946	Org	Loc	Person	Misc	All
				5,187	5,433	4,406	3,438	11,176
AG	30,000	112,000	7,536	Org	Loc	GPE	PII	All
				58,177	39,988	18,506	42,683	67,157
SEC	12,651	5,188	1,592	Org	Loc	GPE	PII	All
				1,955	273	60	357	2,166

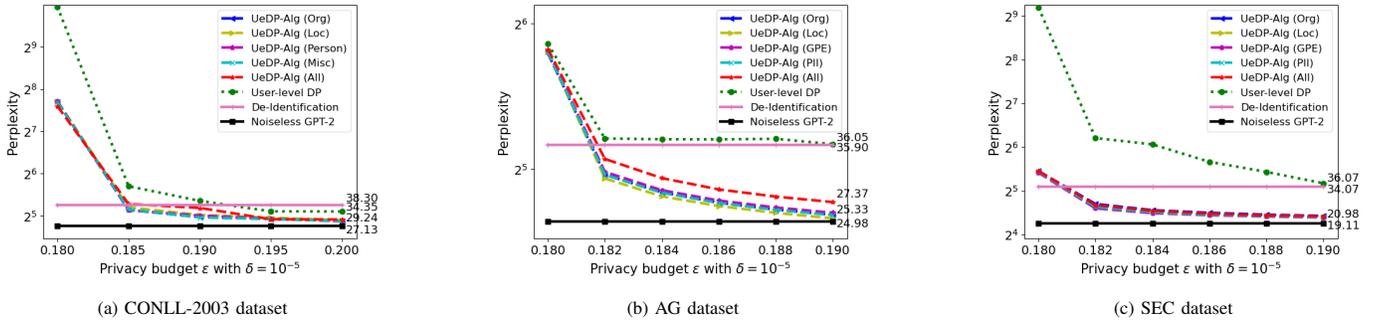


Fig. 2: Next word prediction results using the GPT-2 model. (The lower the better)

entity categories in  $E$  on the privacy budget and model utility; and (5) confirming our results in the text classification task.

Our result is as follows:

- **Integrity of sensitive entities.** Our work utilizes spaCy [39], one of the state-of-the-art large-scale entity recognition systems, to identify sensitive entities for evaluation purposes on datasets that do not have ground-truth sensitive entities, including the AG and SEC datasets. For CONLL-2003, we consider the labels of four sensitive entity types (i.e., location, person, organization, and miscellaneous) from NER as the ground truth. To evaluate the integrity of identified sensitive entities, we conducted a clarification on AMT. We found that the results from spaCy cover over 94% of sensitive information as identified by AMT workers. We recruited master-level AMT workers for a high quality of results, and we provided detailed guidance before AMT workers conducted the task. Each sentence was assigned to 3 workers to mitigate bias and subjective views. Consequently, our experiments using the spaCy identified sensitive entities are solid.

- **Comparing Estimators  $f_{\mathcal{E}}$ ,  $f_{\mathcal{E}+}$ , and User-level DP.** In this analysis, we set  $q_u = 0.05$ ,  $q_e = 0.5$ ,  $q_s = 1$ ,  $z = 2$ , and compute privacy budget  $\epsilon$  at  $\delta = 10^{-5}$  (a typical value of  $\delta$  in DP) as a function of the training steps  $T$ . Fig. 4 shows curves of using different estimators and the User-level DP with all entities in CONLL-2003, AG, and SEC datasets.

Our UeDP-Alg with  $f_{\mathcal{E}+}$  achieves a notably tighter privacy budget compared with  $f_{\mathcal{E}}$  and the User-level DP in all scenarios in CONLL-2003, AG, and SEC datasets. The key reason is that typically detected sensitive entities in  $E$  appear rarely in a dataset compared with extended sensitive entities. Thus, using only sensitive entities in  $E$  identified by the spaCy in training will cause information distortion, which can damage model utility and a loose privacy budget.

User-level DP consumes a much higher privacy budget  $\epsilon$  compared with both of our estimators  $f_{\mathcal{E}+}$  and  $f_{\mathcal{E}}$ . For

instance, at  $T = 50$ , the values of  $\epsilon$  in all entities of  $f_{\mathcal{E}+}$  and  $f_{\mathcal{E}}$ , and the value of  $\epsilon$  of the User-level DP in: (1) the CONLL-2003 dataset are 0.52, 0.62, and 1.18; (2) the AG dataset are 0.50, 0.75, and 1.48; and (3) the SEC dataset are 0.40, 0.71, and 1.40, respectively.

Significantly, the privacy budget ( $\epsilon$ ) gap between User-level DP,  $f_{\mathcal{E}}$ , and  $f_{\mathcal{E}+}$  is proportionally increased to the number of steps  $T$ . That means the more training steps  $T$ , the larger  $\epsilon$  our model can save compared with User-level DP. That is a promising result in the context that our model provides entity DP protection for both users and sensitive entities, compared with only protection for users in User-level DP. We observe a similar phenomenon on different sensitive categories.

- **Privacy Budget ( $\epsilon, \delta$ )-UeDP and Model Utility.** From our theoretical analysis,  $f_{\mathcal{E}+}$  is better than the estimator  $f_{\mathcal{E}}$ . Therefore, for the sake of simplicity, we only consider UeDP-Alg  $f_{\mathcal{E}+}$  instead of showing results from both estimators. From now, UeDP-Alg is used to indicate the use of our estimator  $f_{\mathcal{E}+}$ . Fig. 2 illustrates the perplexity as a function of the privacy budget  $\epsilon$  for an GPT-2 model trained on a variety of sensitive entity categories in UeDP, User-level DP, and De-Identification. The noiseless GPT-2 (for the next word prediction) and BERT (for the text classification) models are considered an upper-bound performance mechanism without offering any privacy protection.

In the CONLL-2003 dataset (Fig. 2a), there are NER labels for person, location, organization, and miscellaneous entities; therefore, we choose these types as sensitive entity categories to protect in UeDP-Alg. UeDP-Alg achieves a better perplexity compared with User-level DP under a tight privacy budget  $\epsilon \in [0.18, 0.20]$ . Also, from  $\epsilon = 0.185$  (a tight privacy protection), our UeDP-Alg achieves a better perplexity than De-Identification. In fact, at  $\epsilon = 0.185$ , our UeDP-Alg achieves 35.09 for person, 35.34 for organization, 35.57 for miscellaneous, and 36.79 for location entities, compared with

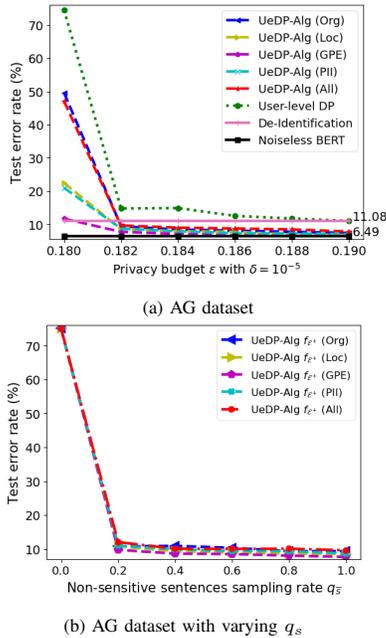


Fig. 3: Text classification results on the AG dataset using the BERT model. With  $q_s = 0.0$ , the test error rate is 75% in all cases. (The lower the better)

52.01 in User-level DP. When spending more privacy budget ( $\epsilon \geq 0.195$ ), both UeDP-Alg and User-level DP converge at a very competitive perplexity level, approaching the upper-bound Noiseless GPT-2. For instance, at  $\epsilon = 0.20$ , there are significant perplexity drops given UeDP-Alg and User-level DP mechanisms, i.e., our UeDP-Alg is 29.24 for person, 29.35 for miscellaneous, 29.58 for organization, and 29.75 for location entities. Meanwhile, the perplexity values of User-level DP, De-Identification, and the noiseless GPT-2 model are 30.15, 38.30, and 27.13.

Results on AG and SEC datasets (Figs. 2b and 2c) further strengthen our observations. In AG and SEC datasets, we applied spaCy to identify different sensitive entity categories, such as GPE, location, organization, and PII (i.e., person and location information). UeDP-Alg achieves better results compared with User-level DP in all considering sensitive entity categories and privacy budgets, and outperforms De-Identification in most cases. That is promising and consistent with our previous analysis. For instance, in the AG dataset, at  $\epsilon = 0.19$ , our UeDP-Alg achieves 25.33 for location, 25.72 for PII, 25.77 for organization, and 26.01 for GPE entities, compared with 36.05 in User-level DP. De-Identification obtains 35.90, and the upper bound result in the noiseless GPT-2 model is 24.98. Similarly, in the SEC dataset (Fig. 2c), at  $\epsilon = 0.19$ , UeDP-Alg achieves perplexity of 20.98 in GPE, 21.12 in PII, 21.22 in location, 21.50 in organization, and 21.33 in all entities, compared with 36.07 in User-level DP, and 34.07 in De-Identification. In AG and SEC datasets, at a tight privacy budget, i.e.,  $\epsilon = 0.19$ , our UeDP-Alg has better perplexity values than the De-Identification, approaching the noiseless GPT-2 model.

• **Sensitive Entity Categories.** In all datasets (Figs. 2 and 9, Appx. E, Supplementary<sup>4</sup>), *the more sensitive sentences to protect, the higher the privacy budget is needed, and the lower performance the model achieves* (i.e., higher perplexity values). For instance, in the SEC dataset, the number of sensitive sentences in each category is as follows: 60 in GPE, 273 in location, 357 in PII, 1,955 in organization, and 2,166 in all entities. After 500 steps, the values of  $\epsilon$  are 0.19 in GPE, 0.24 in location, 0.26 in PII, 0.73 in organization, 0.81 in all entities, and 4.08 in User-level DP (Fig. 4). At  $\epsilon = 0.18$  (Fig. 2c), we obtain perplexity values of 42.63 in GPE, 43.21 in location, 43.30 in PII, 43.70 in organization, 43.77 in all entities, and a 583.06 in User-level DP.

• **Text classification.** Fig. 3a shows that our UeDP-Alg achieves lower test error rates in terms of text classification on the AG dataset than baseline approaches in most cases across different types of sensitive entities under a very tight UeDP protection ( $\epsilon \in [0.18, 0.19]$ ). This is a promising result. When  $\epsilon$  is higher, the test error rates of both UeDP-Alg and User-level DP drop, approaching the noiseless BERT model’s upper-bound result.

• **Extended Sensitive Entities.** To shed light into the impact the extended sensitive entity sampling rate  $q_s$  on model utility under UeDP protection, we varied the value of  $q_s$  from 0 to 1 in all datasets and tasks. Figs. 3b, 5, and 7 show that considering extended sensitive entities (i.e.,  $q_s > 0$ ) significantly improves model utility (i.e., perplexity or test error rate) compared with only considering sensitive entities  $e \in E$  (i.e.,  $q_s = 0$ ). However, different tasks on different datasets may have different optimal values of  $q_s$ . This opens a new research question on how to theoretically approximate the optimal value of  $q_s$ .

Results on the AWD-LSTM model (Figs. 6 and 7) further strengthen our observations. In our experiments, the AWD-LSTM model generally obtains comparable results with the GPT-2 model for next word prediction at a higher privacy budget range (i.e.,  $\epsilon \in [0.5, 3.0]$  in the AWD-LSTM model compared with  $\epsilon \in [0.18, 0.2]$  in the GPT-2 model). This is because the GPT-2 model is pretrained on large-scale datasets, so that it is easily adapted to the idiosyncrasies of a target task (i.e., next word prediction) compared with the AWD-LSTM model trained from scratch.

## VII. CONCLUSION AND FUTURE WORK

In this paper, we developed a novel notion of user-entity DP (UeDP), protecting users’ participation information and sensitive entities in NLMs. By incorporating user and sensitive entity sampling in the training process, we addressed the trade-off between model utility and privacy loss with a tight bound of model sensitivity. Theoretical analysis and rigorous experiments show that UeDP-Alg outperforms baselines in next word prediction and text classification under UeDP protection.

In practice, the list of sensitive entities and users can grow over time. Periodically updating the list of users and sensitive entities may incur extra privacy and computational

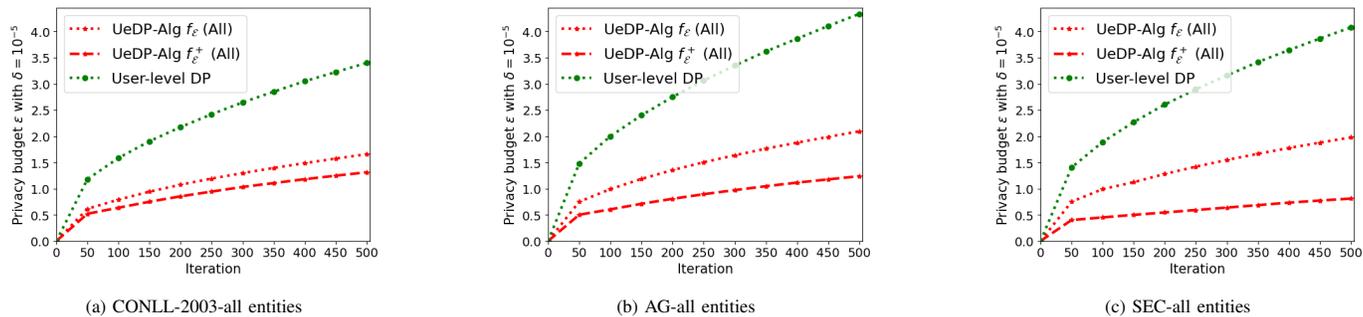


Fig. 4: Privacy budget of UeDP-Alg  $f_{\mathcal{E}}$ , UeDP-Alg  $f_{\mathcal{E}^+}$ , and User-level DP as a function of iterations in CONLL-2003, AG, and SEC datasets. UeDP-Alg  $f_{\mathcal{E}^+}$  achieves a tighter privacy budget compared with UeDP-Alg  $f_{\mathcal{E}}$  and User-level DP.

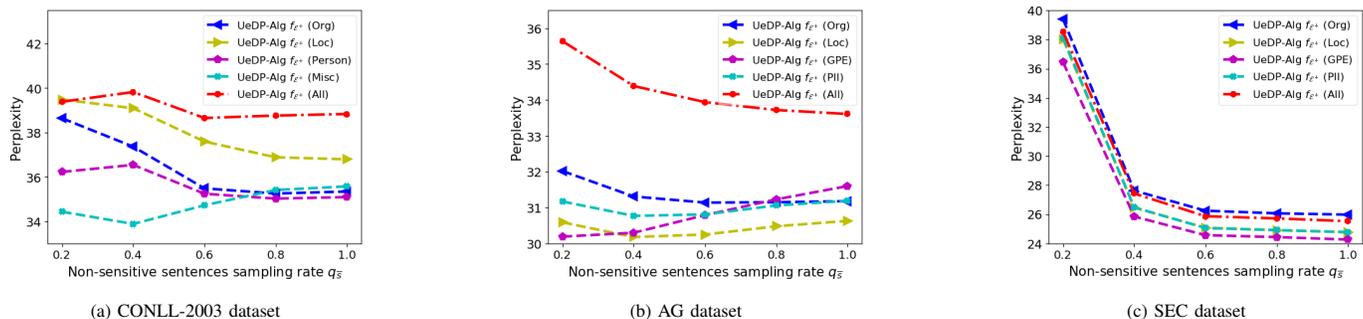


Fig. 5: Next word prediction results using the GPT-2 model with varying extended sensitive entities sampling rate  $q_s$  in training. (The lower the better)

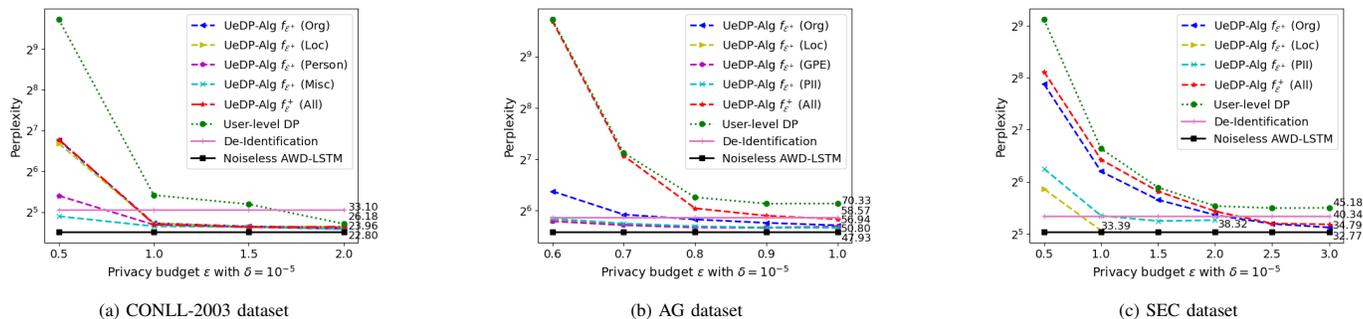


Fig. 6: Next word prediction results using the AWD-LSTM model. (The lower the better)

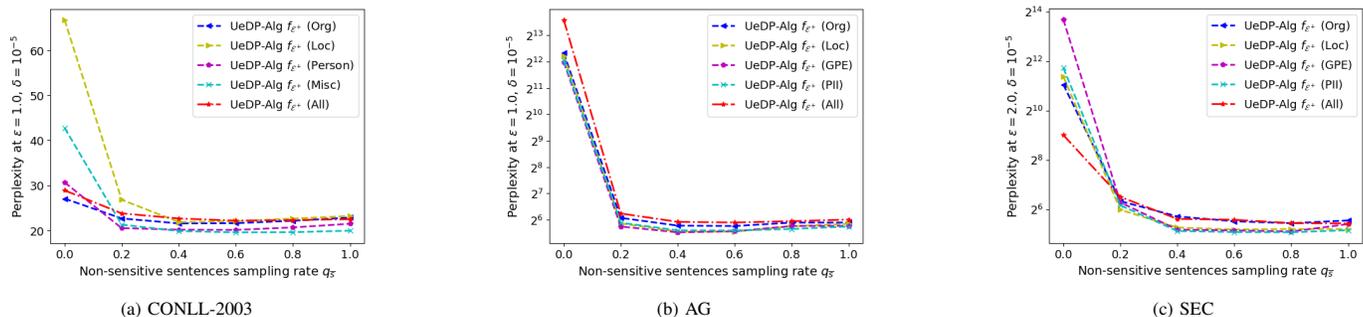


Fig. 7: Next word prediction results using the AWD-LSTM model with varying extended sensitive entities sampling rate  $q_s$  in training. (The lower the better)

cost. Therefore, we will focus on preserving UeDP given a growing list of users and sensitive entities in our future work.

#### ACKNOWLEDGEMENT

This work is partially supported by grants NSF IIS-2041096, NSF CNS-1935928, NSF CNS-1850094, and unrestricted gifts from Adobe System Inc.

#### REFERENCES

- [1] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving language understanding by generative pre-training," 2018.
- [2] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [3] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *NeurIPS*, vol. 33, pp. 1877–1901, 2020.
- [4] N. Carlini, C. Liu, Ú. Erlingsson, J. Kos, and D. Song, "The secret sharer: Evaluating and testing unintended memorization in neural networks," in *USENIX*, 2019, pp. 267–284.
- [5] N. Carlini, F. Tramer, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. Brown, D. Song, U. Erlingsson *et al.*, "Extracting training data from large language models," in *USENIX*, 2021, pp. 2633–2650.
- [6] C. Dwork, A. Roth *et al.*, "The algorithmic foundations of differential privacy," *Fnt-TCS*, vol. 9, pp. 211–407, 2014.
- [7] A. Al Badawi, L. Hoang, C. Mun, K. Laine, and K. Aung, "Privft: Private and fast text classification with homomorphic encryption," *IEEE Access*, vol. 8, pp. 226 544–226 556, 2020.
- [8] C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating noise to sensitivity in private data analysis," in *Theory of Cryptography Conference*, 2006, pp. 265–284.
- [9] M. Abadi, A. Chu, I. Goodfellow, H. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep learning with differential privacy," in *CCS*, 2016, pp. 308–318.
- [10] N. Phan, Y. Wang, X. Wu, and D. Dou, "Differential privacy preservation for deep auto-encoders: an application of human behavior prediction," in *AAAI*, vol. 16, 2016, pp. 1309–1316.
- [11] R. Shokri and V. Shmatikov, "Privacy-preserving deep learning," in *CCS*, 2015, pp. 1310–1321.
- [12] H. Phan, M. Thai, H. Hu, R. Jin, T. Sun, and D. Dou, "Scalable differential privacy with certified robustness in adversarial learning," in *ICML*, 2020, pp. 7683–7694.
- [13] M. Zia, M. Khan, and H. El-Sayed, "Application of differential privacy approach in healthcare data—a case study," in *IIT*, 2020.
- [14] N. Wu, F. Farokhi, D. Smith, and M. A. K., "The value of collaboration in convex machine learning with differential privacy," *IEEE Symposium on Security and Privacy (SP)*, 2020.
- [15] N. Li, W. Qardaji, and D. Su, "On sampling, anonymization, and differential privacy or, k-anonymization meets differential privacy," in *ASIA CCS*, 2012, pp. 32–33.
- [16] H. McMahan, D. Ramage, K. Talwar, and L. Zhang, "Learning differentially private recurrent language models," *ICLR*, 2017.
- [17] L. Lyu, Y. Li, X. He, and T. Xiao, "Towards differentially private text representations," in *ACM SIGIR*, 2020, pp. 1813–1816.
- [18] L. Lyu, X. He, and Y. Li, "Differentially private representation for nlp: Formal guarantee and an empirical study on privacy and fairness," *EMNLP*, 2020.
- [19] E. Bagdasaryan, O. Poursaeed, and V. Shmatikov, "Differential privacy has disparate impact on model accuracy," in *NeurIPS*, 2019, pp. 15 479–15 488.
- [20] A. Roth, "Buying private data at auction: the sensitive surveyor's problem," *ACM SIGecom Exchanges*, vol. 11, no. 1, pp. 1–8, 2012.
- [21] R. Bassily, A. Smith, and A. Thakurta, "Private empirical risk minimization: Efficient algorithms and tight error bounds," in *FOCS*, 2014, pp. 464–473.
- [22] S. Ramaswamy, O. Thakkar, R. Mathews, G. Andrew, H. McMahan, and F. Beaufays, "Training production language models without memorizing user data," *arXiv*, 2020.
- [23] H. Asi, J. Duchi, and O. Javidsbakht, "Element level differential privacy: The right granularity of privacy," *arXiv*, 2019.
- [24] U. Erlingsson, V. Pihur, and A. Korolova, "Rappor: Randomized aggregatable privacy-preserving ordinal response," in *CCS*, 2014, pp. 1054–1067.
- [25] J. Duchi, M. Jordan, and M. Wainwright, "Minimax optimal procedures for locally private estimation," *JASA*, vol. 113, 2018.
- [26] T. Mikolov, A. Deoras, S. Kombrink, L. Burget, and J. Černocký, "Empirical evaluation and combination of advanced language modeling techniques," in *ISCA*, 2011.
- [27] C. Dwork, "Differential privacy: A survey of results," in *TAMC*, 2008, pp. 1–19.
- [28] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, "Membership inference attacks against machine learning models," in *IEEE Symposium on Security and Privacy*, 2017, pp. 3–18.
- [29] A. Salem, Y. Zhang, M. Humbert, P. Berrang, M. Fritz, and M. Backes, "MI-leaks: Model and data independent membership inference attacks and defenses on machine learning models," *NDSS*, 2019.
- [30] S. Yeom, I. Giacomelli, M. Fredrikson, and S. Jha, "Privacy risk in machine learning: Analyzing the connection to overfitting," in *CSF*, 2018, pp. 268–282.
- [31] X. Pan, M. Zhang, S. Ji, and M. Yang, "Privacy risks of general-purpose language models," in *IEEE SP*, 2020, pp. 1314–1331.
- [32] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings *et al.*, "Advances and open problems in federated learning," *Foundations and Trends in Machine Learning*, vol. 14, no. 1–2, pp. 1–210, 2021.
- [33] T. Wang, J. Blocki, N. Li, and S. Jha, "Locally differentially private protocols for frequency estimation," in *USENIX*, 2017, pp. 729–745.
- [34] P. Arachchige, P. Bertok, I. Khalil, D. Liu, S. Camtepe, and M. Atiquzzaman, "Local differential privacy for deep learning," *IoT-J*, vol. 7, no. 7, pp. 5827–5842, 2019.
- [35] S. Wagh, X. He, A. Machanavajjhala, and P. Mittal, "Dp-cryptography: marrying differential privacy and cryptography in emerging applications," *CACM*, vol. 64, pp. 84–93, 2021.
- [36] Z. Yang and Z. Liang, "Automated identification of sensitive data from implicit user specification," *Cybersecurity*, vol. 1, no. 1, pp. 1–15, 2018.
- [37] E. Sang and F. De Meulder, "Introduction to the conll-2003 shared task: Language-independent named entity recognition," *CoNLL*, 2003.
- [38] L. Derczynski, E. Nichols, M. van Erp, and N. Limsopatham, "Results of the WNUT2017 shared task on novel and emerging entity recognition," in *W-NUT*, 2017, pp. 140–147.
- [39] M. Honnibal and I. Montani, "Spacy 2: Natural language understanding with bloom embeddings, convolutional neural networks and incremental parsing," *To appear*, vol. 7, no. 1, pp. 411–420, 2017.
- [40] P. Qi, Y. Zhang, Y. Zhang, J. Bolton, and C. Manning, "Stanza: A python natural language processing toolkit for many human languages," *ACL System Demonstration*, 2020.
- [41] H. McMahan, E. Moore, D. Ramage, and B. y Arcas, "Federated learning of deep networks using model averaging," *arXiv:1602.05629*, 2016.
- [42] F. Dernoncourt, J. Lee, O. Uzuner, and P. Szolovits, "De-identification of patient notes with recurrent neural networks," *JAMIA*, vol. 24, no. 3, pp. 596–606, 2017.
- [43] T. Mikolov, S. Kombrink, L. Burget, J. Černocký, and S. Khudanpur, "Extensions of recurrent neural network language model," in *ICASSP*, 2011, pp. 5528–5531.
- [44] J. Howard and S. Ruder, "Universal language model fine-tuning for text classification," *ACL*, p. 328–339, 2018.
- [45] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *NAACL-HLT*, 2019.
- [46] S. Merity, N. S. Keskar, and R. Socher, "Regularizing and optimizing LSTM language models," *ICLR*, 2018.
- [47] C. Dwork and J. Lei, "Differential privacy and robust statistics," in *STOC*, 2009, pp. 371–380.

### A. Sensitive Entity Recognition and Tool-kits

If a training set does not have sensitive entity indicators, we suggest several ways to identify sensitive entities in textual data, as follows.

**Using Named Entity Recognition (NER) datasets.** NER datasets [37, 38] refer to textual data in which entities in a text are labeled based on several predefined categories. NER typically makes it easy for individuals and systems to identify and understand the subject of the given text quickly. Therefore, extracted entities are critical and should be protected. For instance, in the CONLL-2003 dataset [37], there are four entity types, i.e., location, person, organization, and miscellaneous.

**Using Publicly Available Tool-kits.** For textual datasets that do not have NER labels or sensitive entity indicators, there are publicly available tool-kits for detecting named entities or PII in text, for example, Spacy [39], Stanza [40], and Microsoft Presidio<sup>3</sup>. Spacy and Stanza deploy pre-trained NER models based on statistical learning methods to identify eighteen categories of named entities, including person, nationality or religious groups, facility, etc. (Table I). Microsoft Presidio is another toolbox for PII detectors and NER models based on Spacy and regular expression<sup>6</sup>. For instance, Spacy is used as a sensitive entity identification in Fig. 1 to detect “David Johnson” a person entity, “Main” a GPE entity, “September 18” a date entity, and “Main Hospital” an organization entity.

We present descriptions of different sensitive entity categories in the CONLL-2003, AG, and SEC datasets in Table I. The descriptions are from [37] and spaCy, supporting eighteen different entity types. In the current work, we play with four different types and their combinations. Note that, in UeDP, providing the name of an algorithm and a sensitive entity means we consider that type of entity as sensitive entities in the training process. For instance, in Fig. 4, UeDP-Alg  $f_{\mathcal{E}^+}$  (Org) means we use all organization entities as sensitive entities in the UeDP-Alg algorithm. “All entities” means all types of sensitive entities considered for the dataset are used. For example, “all entities” in the CONLL-2003 dataset means all person, location, organization, and miscellaneous entities are regarded as sensitive entities. Meanwhile, in the AG and SEC datasets, “all entities” means that all organization, location, GPE, and PII entities are considered sensitive entities. More entity types are also presented in Table I so that users can have more choices when identifying sensitive entities.

### B. UeDP without Considering Extended Sensitive Entities

At each iteration  $t$ , we randomly sample  $U^t$  users from  $U$  and  $E^t$  sensitive entities from  $E$ , with sampling rates  $q_u$  and  $q_e$ , respectively. Then, we use all sensitive sentences consisting of the sensitive entities in  $E^t$  belonging to the selected users in  $U^t$  for training. Like [16], we leverage the basic federated learning setting in [41] to compute gradients of model parameters for a particular user, denoted as  $\Delta_{u,\mathcal{E}}^{t+1}$ . Here, we clip the per-user gradients so that its  $l_2$ -norm is

bounded by a predefined gradient clipping bound  $\beta$ . Next, a weighted-average estimator  $f_{\mathcal{E}}$  is employed to compute the average gradient  $\Delta^{t+1}$  using the clipped gradients  $\Delta_{u,\mathcal{E}}^{t+1}$  gathered from all the selected users. Finally, we add random Gaussian noise  $\mathcal{N}(0, I\sigma^2)$  to the model update. During the training, the moments accountant  $\mathcal{M}$  is used to compute the  $T$  training steps’ privacy budget consumption.

In this process, we need to bound the sensitivity of the weighted-average estimator  $f_{\mathcal{E}}$  for per-user gradients  $\Delta_{u,\mathcal{E}}^{t+1}$ . We first consider the following simple estimator, with both sampling rates  $q_u$  for the user-level and  $q_e$  for the sensitive entity-level:

$$f_{\mathcal{E}}(U^t, E^t) = \frac{\sum_{u \in U^t} w_u \Delta_{u,\mathcal{E}}^{t+1}}{q_u W_u q_e W_e} \quad (7)$$

$$s.t. \Delta_{u,\mathcal{E}}^{t+1} = \sum_{e \in E_u^t} w_e \left( \sum_{s \text{ consists of } e} \Delta_{u,s} \right)$$

where  $w_u$  and  $w_e \in [0, 1]$  are weights associated with a user  $u$  and with a sensitive entity  $e$ . These weights capture the influence of a user and a sensitive entity to the model outcome.  $\Delta_{u,s}$  is the parameter gradients computed using a sensitive sentence  $s$  consisting of the sensitive entity  $e$ . In addition,  $W_u = \sum_u w_u$  and  $W_e = \sum_e w_e$ .

The estimator  $f_{\mathcal{E}}$  is unbiased to the sampling process; since  $\mathbb{E}[\sum_{u \in U^t} w_u] = q_u W_u$  and  $\mathbb{E}[\sum_{e \in E_u^t} w_e] = q_e W_e$ . The sensitivity of the estimator  $f_{\mathcal{E}}$  can be computed as:  $\mathbb{S}(f_{\mathcal{E}}) = \max_{u', e'} \|f_{\mathcal{E}}(\{U^t \cup u', E^t \cup e'\}) - f_{\mathcal{E}}(\{U^t, E^t\})\|_2$ , where the added user  $u'$  can have arbitrary data and  $e'$  is an arbitrary sensitive entity.

Given that  $\Delta_{u,\mathcal{E}}^{t+1}$  is  $l_2(\beta)$ -norm bounded, where  $\beta$  is the radius of the norm ball by replacing  $\Delta_{u,\mathcal{E}}^{t+1}$  with  $\Delta_{u,\mathcal{E}}^{t+1} \cdot \min\left(1, \frac{\beta}{\|\Delta_{u,\mathcal{E}}^{t+1}\|}\right)$ , the sensitivity of  $\mathbb{S}(f_{\mathcal{E}})$  is also bounded.

**Lemma 2.** *If for all users  $u$  we have  $\|\Delta_{u,\mathcal{E}}^{t+1}\|_2 \leq \beta$ , then  $\mathbb{S}(f_{\mathcal{E}}) \leq \frac{(q_u|U|+1)\max(w_u)\beta}{q_u W_u \times q_e W_e}$ .*

*Proof.* If for all users  $u$  we have  $\|\Delta_{u,\mathcal{E}}^{t+1}\|_2 \leq \beta$ , then  $\mathbb{S}(f_{\mathcal{E}})$

$$= \frac{\sum_{u \in U^t \cup u'} w_u [(\sum_{e \in E^t} w_e (\sum_{s \in S_{ue}^t} \Delta_{u,s}))]}{(q_u W_u \times q_e W_e)}$$

$$+ \frac{\sum_{u \in U^t \cup u'} w_u [w_{e'} (\sum_{s \in S_{ue'}^t} \Delta_{u,e'})]}{(q_u W_u \times q_e W_e)}$$

$$- \frac{\sum_{u \in U^t} w_u [\sum_{e \in E^t} w_e (\sum_{s \in S_{ue}^t} \Delta_{u,s})]}{(q_u W_u \times q_e W_e)}$$

$$\leq \frac{\sum_{u \in U^t \cup u'} w_u \beta}{q_u W_u \times q_e W_e} \leq \frac{(q_u|U|+1)\max(w_u)\beta}{q_u W_u \times q_e W_e} \quad (8)$$

Consequently, Lemma 2 holds.  $\square$

By applying Lemma 2, given a hyper-parameter  $z$ , the noise scale  $\sigma$  for the estimator  $f_{\mathcal{E}}$  is:

$$\sigma = z\mathbb{S}(f_{\mathcal{E}}) = \frac{z(q_u|U|+1)\max(w_u)\beta}{q_u W_u \times q_e W_e} \quad (9)$$

We show that this approach achieves  $(\epsilon, \delta)$ -UeDP, by applying the moments accountant  $\mathcal{M}$  to bound the total privacy

<sup>6</sup><https://github.com/google/re2/>

loss of  $T$  steps of the Gaussian mechanism with the noise  $\mathcal{N}(0, I\sigma^2)$  in Theorem 1. However, this mechanism only uses sensitive entities detected by automatic toolkits to train the model ignoring a large number of extended sensitive entities. As a result, it introduces a loose sensitivity bound (Lemma 2) and affects our model utility.

### C. Datasets and Data Processing

CONLL-2003 consists of Reuters news stories published between August 1996 and August 1997. CONLL-2003 is an NER dataset, where there are labels for four different types of named entities, including location, organization, person, and miscellaneous entities. These types of named entities are considered sensitive entities. In the CONLL-2003 dataset, there is no obvious user information; hence, we consider each document as a user consisting of multiple sentences in the next word prediction task.

AG dataset is a collection of news articles gathered from more than 2,000 news sources by ComeToMyHead academic news search engine<sup>7</sup>. It is categorized into four classes: world, sport, business, and science/technology. Similar to the CONLL-2003 dataset, there is no user information in AG. To imitate a user indicator, we randomly divide news into different users based on Gaussian distribution. There are no named entities; thus, we apply pre-trained Spacy to find named entities and PII in the dataset. We choose different types of these named entities to be sensitive entities: organization, GPE (i.e., countries, cities, and states), location, and PII entities.

Our SEC dataset consists contract clauses collected from contracts submitted in SEC filings<sup>8</sup>. Since the contracts can be associated with a company ID, we use the ID as a user indicator. Similar to the AG dataset, we consider organization, GPE, location, and PII entities as sensitive entities to protect.

In addition to the next word prediction, we conducted text classification on the AG dataset to further strengthen our observations. For text classification, the number of labels is not sufficient in the SEC dataset, and the labels do not exist in the CONLL-2003 dataset. Therefore, we do not utilize CONLL-2003 and SEC datasets for text classification in this study.

For data preprocessing, we changed all words to lower-case and removed punctuation marks. Fig. 8 shows the distribution of the number of users and sentences in the CONLL-2003, AG, and SEC datasets. In the CONLL-2003 dataset, there is no obvious user information; hence, we consider each document as a user consisting of multiple sentences. Like the CONLL-2003 dataset, in the AG dataset, there is no user information. Therefore, to imitate a user indicator, we randomly divide news into different users. The number of sentences per user follows a Gaussian distribution  $\mathcal{N}(15, 2^2)$ , i.e., there are 15 sentences per user on average, and the standard deviation is 2 sentences. In the SEC dataset, since the contracts can be associated with a company ID, we use the ID as a user indicator. The document related to the ID is considered to be that user’s data.

<sup>7</sup><http://newsengine.di.unipi.it/>

<sup>8</sup><https://www.sec.gov/edgar.shtml>

### D. Revisiting Word-level LDP Analysis in [18]

This section aims at revisiting privacy protection in [18] and describes a privacy accumulation issue over the embedding dimension. Then, we revise Theorems 1 and 2 in [18] and compare them with our approaches.

In [18], the authors aim at preserving the privacy of the extracted test representation from users while maintaining the good performance of the classifier, which is trained at a server by the data collected from users. To achieve the goal, they consider a word-level DP, that is, two inputs  $x$  and  $x'$  are adjacent if they differ by at most 1 word. Additionally, they introduce a DP noise layer  $r$  after a predefined feature extractor  $f(x)$ . To train a robust classifier at the server, they add the same level of noise as the test phase in the training process and optimize the classifier by minimizing the loss function as follows:

$$\mathcal{L}(x, y) = \mathcal{X}(C(f(x) + r), y) \quad (10)$$

where  $C$  is the classifier,  $y$  is the true label, and  $\mathcal{X}$  is the cross entropy loss function.

The Laplace noise layer  $r$  is injected into the embedding  $f(x)$  in which its coordinates  $r = \{r_1, r_2, \dots, r_k\}$  are random variables drawn from the Laplace distribution defined by  $Lap(b)$  with  $b = \frac{\Delta_f}{\epsilon}$ ,  $\epsilon$  is the privacy budget, and  $\Delta_f$  is the sensitivity of the extracted representation. Here,  $k$  is the dimension of  $f(x)$ .

Algorithm 1 describes how to derive DP-preserving representation from the feature extractor  $f$ . Note that  $x_s$  in the Algorithm 1 is a sentence (equivalent to  $x$  in our notation), which is considered to be sensitive and needs to be protected.

**Revisiting Theorems 1 and 2 in [18].** In the paper, the authors consider adjacent sentences differing by one word. Changing one word in  $x$  may change the entire embedding vector  $f(x)$ . Each element of  $f(x)$  is normalized into the range  $[0, 1]$  (Line 5, Algorithm 1), hence each element sensitivity of  $f(x)$  is  $\Delta_f = 1$ , the noise is  $Lap(\Delta_f/\epsilon)$ . Therefore, each element of the embedding  $f(x)$  consumes a privacy budget  $\epsilon$ . Since the  $k$  elements of the embedding are derived from a single sensitive input  $x$ , applying the LDP mechanism  $\mathcal{A}(\cdot)$ , i.e.,  $Lap(b)$ ,  $k$  times will consume the privacy budget  $k \times \epsilon$ . This follows the composition property in DP. Note that the  $k$  elements cannot be treated by using the parallel property in DP [47], since all of them are derived from a single (data) input  $x$ , NOT from  $k$  different inputs ( $k$  different data samples). Consequently, the privacy guarantees in Theorems 1 and 2 of [18] is  $k\epsilon$ -DP, instead of  $\epsilon$ -DP as reported.

In their experimental results, e.g., Table 2 of [18], the approach could achieve almost the same (and even better) model utility with noiseless model given the extremely low  $\epsilon = 0.05$  using BERT embeddings. As our analysis, the privacy budget in Theorems 1 and 2 is  $k\epsilon$ , instead of  $\epsilon$ . Therefore, the proper privacy budget is at least  $0.05 \times 768 = 38.4$ . Similar results were reported through out the all in experiments. With this high value of the privacy budget, the word-level DP in [18] provides loose privacy protection.

---

**Algorithm 1** Differentially Private Neural Representation (DPNR) [18]

---

- 1: **Input:** Each sensitive input  $x_s \in \mathbb{R}^d$ , feature extractor  $f$
  - 2: **Parameters:** Dropout vector  $I_n \in \{0, 1\}^d$
  - 3: Word dropout:  $\tilde{x}_s \leftarrow x_s \odot I_n$ , where  $\odot$  performs a word-wise multiplication.
  - 4: Extraction:  $x_r \leftarrow f(\tilde{x}_s)$
  - 5: Normalization:  $x_r \leftarrow x_r - \min(x_r) / (\max(x_r) - \min(x_r))$
  
  - 6: Perturbation:  $\hat{x}_r \leftarrow x_r + r$ ,  $r_i \sim \text{Lap}(b)$
  - 7: **Output:** Perturbed representation  $\hat{x}_r$ .
- 

**Revisiting Element-level DP in [18].** During our discussion with the authors of [18], the authors mentioned that their approach preserves a new notion of  $(\epsilon, 0)$ -element-level DP, i.e., two embeddings differ from one element, instead of a word-level DP. However, for the element-DP to hold, all the elements in the embedding  $f(x)$  must be independent from each other, that is, changing one element will not result in changing any other element. If changing one element results in changing all the remaining elements, then element-DP will be suffered from the dimension of the embedding by following group privacy. In the current approach, changing one element means there is a change in the input data  $x$  to occur. Equivalently, using BERT, any change in the input data  $x$  will result in changing the whole embedding (all elements). Therefore, the condition of two neighboring embeddings only differing in only one element does NOT hold in theory and practice. Consequently, the introduced element-level DP does NOT hold at the level of  $(\epsilon, 0)$ -DP.

**Our revising Theorems 1 and 2 in [18].** Based upon our analysis, we introduce revised versions of the Theorems 1 and 2 in [18], as follows.

**Theorem 1. Revised Theorem 1 in [18].** Let the entries of the noise vector  $r$  be drawn from  $\text{Lap}(b)$  with  $b = \frac{\Delta_f}{\epsilon}$ . The Algorithm 1 is  $k\epsilon$ -word-level DP, where  $k$  is dimension of the embedding  $f(x)$ .

*Proof.* Each element of the embedding  $f$  is bounded in  $[0, 1]$ , so  $\Delta_f = 1$  for each element. By adding random noise variables drawn from the Laplace  $\text{Lap}(b)$  with  $b = \frac{\Delta_f}{\epsilon}$  into each element of  $f$ , each element consumes  $\epsilon/k$  privacy budget. Since the  $k$  elements of the embedding are derived from a single sensitive input  $x$ , applying the mechanism  $\text{Lap}(b)$   $k$  times on the  $k$  elements will consume the privacy budget  $k\epsilon$ . Therefore, the Algorithm 1 is  $k\epsilon$ -word-level DP.  $\square$

**Theorem 2.** Given an input  $x \in D$ , suppose  $\mathcal{A}(x) = f(x) + r$  is  $k\epsilon$ -word-level DP, let  $I_n$  with dropout rate  $\mu$  be applied to

$x$ :  $\tilde{x} = x \odot I_n$ , then  $\mathcal{A}(\tilde{x})$  is  $\epsilon'$ -word level-DP, where  $\epsilon' = \ln[(1 - \mu) \exp(k\epsilon) + \mu]$ .

*Proof.* Suppose there are two adjacent inputs  $x_1$  and  $x_2$  that differ only in the  $i$ -th coordinate (word), say  $x_{1i} = v$ ,  $x_{2i} \neq v$ . For arbitrary binary vector  $I_n$ , after dropout,  $\tilde{x}_1 = x_1 \odot I_n$ ,  $\tilde{x}_2 = x_2 \odot I_n$ , there are two possible cases, i.e.,  $I_{ni} = 0$  and  $I_{ni} = 1$ .

If  $I_{ni} = 0$ : Since  $x_1$  and  $x_2$  differ only in  $i$ -th coordinate, after dropout  $\tilde{x}_{1i} = \tilde{x}_{2i} = 0$ , hence  $x_1 \odot I_n = x_2 \odot I_n$ . Then  $\Pr\{\mathcal{A}(x_1 \odot I_n) = S\} = \Pr\{\mathcal{A}(x_2 \odot I_n) = S\}$ .

If  $I_{ni} = 1$ : Since  $x_1$  and  $x_2$  differ only in  $i$ -th coordinate, after dropout  $\tilde{x}_{1i} = v$ , and  $\tilde{x}_{2i} \neq v$ . Since  $\mathcal{A}(x)$  is  $k\epsilon$ -word level-DP, then  $\Pr\{\mathcal{A}(x_1 \odot I_n) = S\} \leq \exp(k\epsilon) \Pr\{\mathcal{A}(x_2 \odot I_n) = S\}$ .

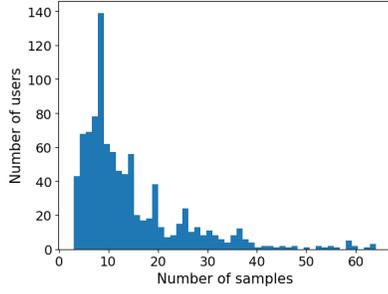
Combining these two cases, and  $\Pr[I_{ni} = 0] = \mu$ , we have:

$$\begin{aligned} & \Pr\{\mathcal{A}(x_1 \odot I_n) = S\} \\ &= \mu \Pr\{\mathcal{A}(x_1 \odot I_n) = S\} + (1 - \mu) \Pr\{\mathcal{A}(x_1 \odot I_n) = S\} \\ &\leq \mu \Pr\{\mathcal{A}(x_2 \odot I_n) = S\} \\ &+ (1 - \mu) \exp(k\epsilon) \Pr\{\mathcal{A}(x_2 \odot I_n) = S\} \\ &= [(1 - \mu) \exp(k\epsilon) + \mu] \Pr\{\mathcal{A}(x_2 \odot I_n) = S\} \\ &= \exp\left(\ln[(1 - \mu) \exp(k\epsilon) + \mu]\right) \Pr\{\mathcal{A}(x_2 \odot I_n) = S\} \end{aligned} \tag{11}$$

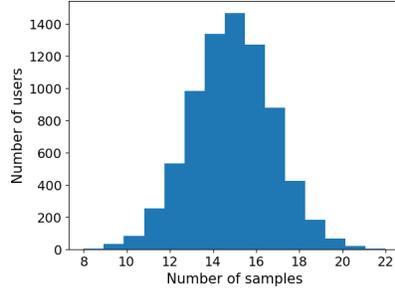
Therefore, after dropout, the privacy budget is  $\epsilon' = \ln[(1 - \mu) \exp(k\epsilon) + \mu]$ .  $\square$

**Comparison with UeDP.** Apart from the privacy accumulation over the embedding dimension, in [18], during training the model, the Laplace or Gaussian noise is drawn at every training iteration. Therefore, the model accesses the raw data at every iteration. As a result, the privacy budget at the training phase is accumulated over the number of training iterations, which can be a large number causing an exploded privacy budget in training. [18] focuses on protecting privacy at the inference time and use the noise in the training phase to obtain a more robust model without considering training data privacy. This is different from our goal to protect users and sensitive entities of training data, which is a more challenging task. Our UeDP-preserving model can be deployed to the end-users for a direct use in the inference phase, without demanding that the end-users send their data embedding to our server; therefore offering a more rigorous privacy protection and better usability. In addition to this, our approach offers more rigorous DP budget bounds compared with the DPNR algorithm in [18], since DPNR consumes large DP budgets that is proportional to the commonly large dimension of the embedding  $k$ .

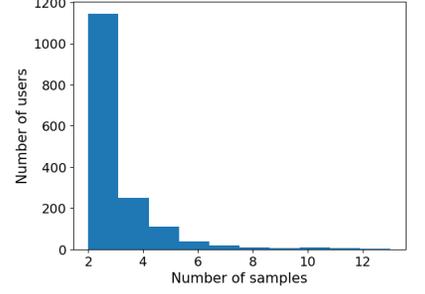
### E. Supplemental Experimental Results



(a) CONLL-2003 dataset

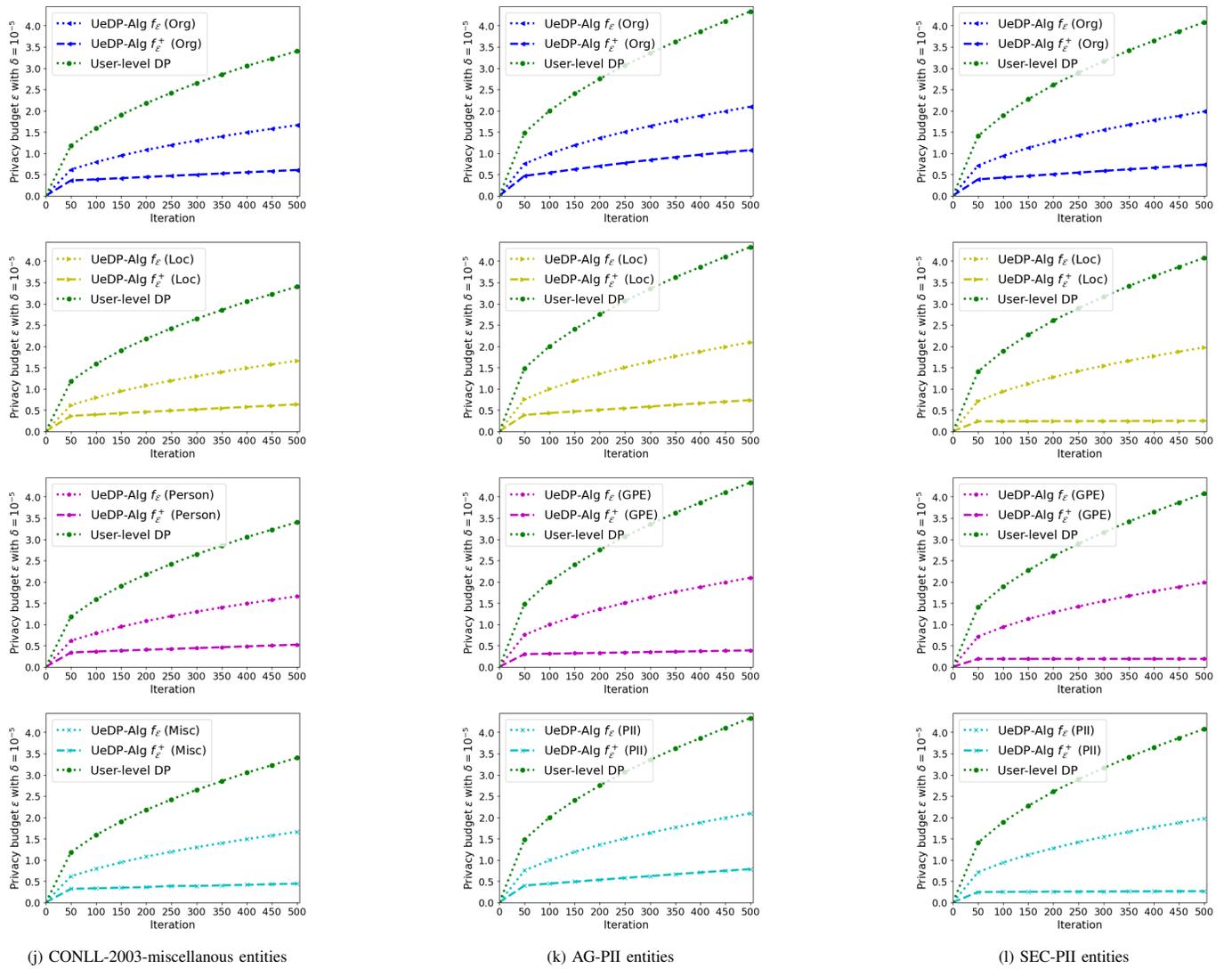


(b) AG dataset



(c) SEC dataset

Fig. 8: Distribution of users and sentences.



(j) CONLL-2003-miscellaneous entities

(k) AG-PII entities

(l) SEC-PII entities

Fig. 9: Privacy budget of UeDP-Alg  $f_{\mathcal{E}}$ , UeDP-Alg  $f_{\mathcal{E}^+}$ , and User-level DP as a function of iterations in CONLL-2003, AG, and SEC datasets. UeDP-Alg  $f_{\mathcal{E}^+}$  achieves a tighter privacy budget compared with UeDP-Alg  $f_{\mathcal{E}}$  and User-level DP.