

Differentially Private Natural Language Models: Recent Advances and Future Directions

Anonymous EACL submission

Abstract

Recent developments in deep learning have led to great success in various natural language processing (NLP) tasks. However, these applications may involve data that contain sensitive information. Therefore, how to achieve good performance while also protecting the privacy of sensitive data is a crucial challenge in NLP. To preserve privacy, Differential Privacy (DP), which can prevent reconstruction attacks and protect against potential side knowledge, is becoming a de facto technique for private data analysis. In recent years, NLP in DP models (DP-NLP) has been studied from different perspectives, which deserves a comprehensive review. In this paper, we provide the first systematic review of recent advances in DP deep learning models in NLP. In particular, we first discuss some differences and additional challenges of DP-NLP compared with the standard DP deep learning. Then we investigate some existing work on DP-NLP and present its recent developments from two aspects: gradient perturbation based methods and embedding vector perturbation based methods. We also discuss some challenges and future directions.

1 Introduction

The recent advances in deep neural networks have led to significant success in various tasks in Natural Language Processing (NLP), such as sentiment analysis, question answering, information retrieval, and text generation. However, such applications always involve data that contains sensitive information. For example, a model of aid typing on a model keyboard is trained from language data which might contain sensitive information such as passwords, text messages, and search queries. Moreover, language data can also identify a speaker explicitly by name or implicitly, for example via a rare or unique phrase. Thus, one often encountered challenge in NLP is how to handle this sensitive information. To overcome the challenge, privacy-

preserving NLP has been intensively studied in recent years. One of the commonly used approaches is based on text anonymization (Pilán et al., 2022), which identifies sensitive attributes and then replaces these sensitive words with some other values. Another approach is injecting additional words into the original text without detecting sensitive entities in order to achieve text redaction (Sánchez and Batet, 2016). However, removing personally identifiable information or injecting additional words is often unsatisfactory, as it has been shown that an adversary can still infer an individual’s membership in the dataset with high probability via the summary statistics on the datasets (Narayanan and Shmatikov, 2008). Moreover, recent studies claim that deep neural networks for NLP tasks often tend to memorize their training data, which makes them vulnerable to leaking information about training data (Shokri et al., 2017; Carlini et al., 2021, 2019). One way that takes into account the limitations of existing approaches by preventing individual re-identification and protecting against any potential data reconstruction and side-knowledge attacks is designing Differentially Private (DP) algorithms. DP (Dwork et al., 2006) provides provable protection against identification and is resilient to arbitrary auxiliary information that might be available to attackers. Thanks to its formal guarantees, DP has become a de facto standard tool for private statistical data analysis.

Although there are numerous studies on DP machine learning and DP deep learning such as (Abadi et al., 2016; Bu et al., 2019; Yu et al., 2019), most of them mainly focus on either the continuous tabular data or image data and less attention has been paid to adapting variants of DP algorithms to the context of NLP and the text domain. On the other side, while there are several surveys on DP and its applications such as (Ji et al., 2014; Dankar and Emam, 2013; Xiong et al., 2020; Wang et al., 2020a; Desfontaines and Pejó, 2020), all of them

do not study its applications to the NLP domain. Recently, [Klymenko et al. \(2022\)](#) gave a brief introduction to applications of DP in NLP, but the reviewed work is not exhaustive and it lacks a technical and systematic view of DP-NLP. Thus, to fill in this gap, in this paper, we provide the first technical overview of the recent developments and challenges of DP in language models.

Specifically, we give a survey on the most recent 65¹ papers on deep learning based approaches for NLP tasks under DP constraints. First, we show some specificities of DP-NLP compared with the general deep learning with DP. Then we discuss current results from two perspectives via the ways of adding randomness to ensure DP: the first one is gradient perturbation based methods which includes DP-SGD and DP-Adam; the second one is embedding vector perturbation based methods which includes DP auto-encoder. For each type of approach, we also consider its applications to different NLP tasks. Finally, we present some potential challenges and future directions.

Due to space limits, in [Appendix A](#) we give a preliminary introduction to DP to readers who are unfamiliar with DP.

2 Specificities of NLP with DP

We first discuss some specificities for DP-NLP compared with the standard DP deep learning. Generally speaking, there are two aspects, one is privacy notations and another is privacy levels.

2.1 Variants of DP Notions in NLP

Recall that DP ensures data analysts or adversaries will get almost the same information if we change any single data sample in the training data, i.e., it treats all records as sensitive. However, such an assumption is quite stringent. On the one side, unlike image data, for text data it is more common that only several instead of all attributes need to be protected. For example, for the sentence "My cell phone number is 1234567890", only the last token with the actual cell phone number needs to be protected. On the other side, canonical DP requires that the log of the ratio between the distribution probabilities is always upper bounded by the privacy parameter ϵ for any pair of neighboring data. However, such a requirement is also quite restrictive. For example, for the sentence "I will arrive at 2:00 pm", we want the adversary not to distin-

¹Note that we did not cover all related works, see the Limitations and Future Directions sections for the works that are not included in this paper.

guish it from the sentence "I will arrive at 4:00 pm". However, DP also can ensure the adversary cannot distinguish it from the sentence "I will arrive at 10:00 pm", which is meaningless. Thus, for language data, besides the canonical DP, it is also reasonable to study its relaxations for some specific scenarios. Actually, this is quite different from the existing work on DP deep learning, which mainly focuses on standard DP definitions. In the following, we will discuss some commonly used relaxations of DP for language models.

SDP. As we mentioned above, in some scenarios, the sensitive information in text data is sparse and we only need to protect some sensitive attributes instead of the whole sentence. Based on this, [Shi et al. \(2021\)](#) propose a new privacy notion namely selective differential privacy (SDP), to provide privacy guarantees on the sensitive portion of the data to improve model utility. From the definition aspect, the main difference between SDP and DP is the definition of neighboring datasets. Informally, in SDP, two datasets are adjacent if they differ in at least one sensitive attribute. However, it is hard to define such neighboring datasets directly as there are some correlations between sensitive and non-sensitive attributes, indicating that we can still infer information on sensitive attributes ([Kifer and Machanavajjhala, 2011](#)). To address the issue, [Shi et al. \(2021\)](#) leverage the Pufferfish framework in ([Kifer and Machanavajjhala, 2014](#)).

Metric DP. To relax the requirement that the log probability ratio is uniformly bounded by ϵ for all neighboring data pairs, [Feyisetan et al. \(2020\)](#) first adopt the Metric DP (or d_χ -privacy) to the problem of private embedding, which is proposed by ([Chatzikokolakis et al., 2013](#)) for location data originally. In particular, a Metric DP mechanism could report a token in a privacy-preserving manner while giving higher probability to tokens that are close to the current token, and negligible probability to tokens in a completely different part of the vocabulary, where we will use some distance function d to measure the distance between two tokens.

Definition 1. For a data domain (vocabulary) \mathcal{X} , a randomized algorithm $\mathcal{A} : \mathcal{X} \mapsto \mathcal{R}$ is called (ϵ, δ) -Metric DP with distance function d if for any $S, S' \in \mathcal{X}^l$ and $T \subseteq \mathcal{R}$ we have

$$\Pr[\mathcal{A}(S) \in T] \leq e^{d(S,S')\epsilon} \Pr[\mathcal{A}(S') \in T] + \delta.$$

From the above definition, we can see the probability ratio of observing any particular output y

181 given two possible inputs S and S' is bounded by
182 $e^{\epsilon d(S', S)}$ instead of e^ϵ in DP. Motivated by Metric
183 DP and local DP, (Feyisetan et al., 2020) provides
184 the Local Metric DP (LMDP) and uses it for pri-
185 vate word embeddings (see Section 4 for details).
186 Motivated by Utility-optimized LDP (ULDP) (Mu-
187 rakami and Kawamoto, 2019) rather than LDP, re-
188 cently Yue et al. (2021) propose Utility-optimized
189 Metric LDP (UMLDP). It exploits the fact that
190 different inputs have different sensitivity levels to
191 achieve higher utility. By assuming the input space
192 such as the set of tokens is split into sensitive and
193 non-sensitive parts, UMLDP achieves a privacy
194 guarantee equivalent to LDP for sensitive inputs.

195 2.2 Variants Levels of Privacy in NLP

196 When we consider using DP, the first question is
197 what kind of information we aim to protect. In the
198 previous studies on DP deep learning, we always
199 wanted to protect the whole data sample. However,
200 in the NLP domain, such one data sample could be
201 either a word, a sentence a paragraph, etc. If we
202 ignore the concrete privacy level and directly apply
203 the previous DP methods, we may have mediocre
204 results. Thus, unlike the sample level privacy in DP
205 deep learning, researchers in NLP consider differ-
206 ent levels of privacy. Especially, they focus on the
207 word level and sentence level, which aims to pro-
208 tect each word and sentence respectively (Meehan
209 et al., 2022; Feyisetan et al., 2019).

210 In the federated learning setting, there is a cen-
211 tral server and several users each of them has a
212 local dataset, the sample level of DP may be insuf-
213 ficient. For example, in language modeling each
214 user may contribute many thousands of words to
215 the training data and each typed word makes its
216 own contribution to the RNN’s training objective.
217 In this case, just protecting each word is unsatis-
218 factory and it is still possible to re-identify users.
219 Thus, besides the sample level, we also have the
220 user level of privacy, which aims to protect users’
221 histories.

222 After discussing some specificities of DP-NLP.
223 In the following we categorize its recent studies
224 into two classes based on their methods to ensure
225 DP: gradient perturbation based methods and em-
226 bedding vector perturbation based methods. See
227 Tab. 1 in Appendix for an overview.

228 3 Gradient Perturbation Based Methods

229 Generally speaking, a gradient perturbation method
230 is based on adding noises to gradients of the loss

231 during training the network to ensure DP. As the
232 baseline and canonical algorithm for this type of ap-
233 proach, Differentially Private Stochastic Gradient
234 Descent (DP-SGD) (Abadi et al., 2016) is a DP ver-
235 sion of SGD. Its main idea is to use the noisy and
236 clipped subsampled gradient g^t to approximate the
237 whole gradient $\nabla L(\theta^t, D)$. In fact, besides SGD,
238 we can use this idea for any optimizer, such as
239 Adam (Kingma and Ba, 2015), whose private ver-
240 sion DP-Adam is proposed and applied in BERT by
241 (Anil et al., 2021). In the past few years, there has
242 been a long list of work on DP-SGD from differ-
243 ent perspectives, such as the subsampling strategy,
244 faster clipping procedures, private clipping param-
245 eter tuning, and the selection of batch size. In the
246 following, we will only discuss the previous work
247 on using DP-SGD-based methods for variants of
248 NLP tasks. See Appendix B for an introduction to
249 DP-SGD.

250 3.1 DP Pre-trained Models

251 Recent developments in NLP have led to successful
252 applications in large-scale language models with
253 the appearance of transformer (Devlin et al., 2019).
254 It combines the contextual information into lan-
255 guage models with a more powerful ability of rep-
256 resentation. These models are called pre-trained
257 models, which train word embedding in large cor-
258 pora targeting various tasks and gain the knowledge
259 for downstream tasks (Peters et al., 2018). In this
260 section, we review some papers that focus on pre-
261 trained NLP models under DP constraints.

262 The workflow of BERT (Devlin et al., 2019) is
263 pre-training the unlabeled text using some large cor-
264 pora first. Then, the downstream tasks first initial-
265 ize the model using the same parameters and fine-
266 tune the parameters according to different tasks.
267 Despite the benefits of powerful representation abil-
268 ity given by the pre-training process, it also has
269 privacy issues since the model would memorize
270 sensitive information such as words or phrases.

271 In order to solve this privacy leakage issue, there
272 are several studies on how to train BERT privately.
273 Hoory et al. (2021) successfully trained a differ-
274 entially private BERT model by modifying the
275 WordPiece algorithm to satisfy DP, and conducted
276 experiments on the problem of entity extraction
277 tasks from medical text. They construct a tailored
278 domain-specific DP-based trained vocabulary de-
279 signed to generate a new domain-specific vocabu-
280 lary while maintaining user privacy and then use
281 the original DP-SGD in the training process. For

the DP vocabulary part, they first construct a word histogram by dividing the text into a sequence of N -word tuples and then add Gaussian noise to the histogram to ensure (ϵ, δ) -DP. Finally, they clip the histogram with some threshold. For the training phase, they use the original DP-SGD to meet privacy guarantees. Besides, they also use the parallel training trick to make the training faster. Very recently [Yin and Habernal \(2022\)](#) apply DP-BERT to the legal NLP domain. While DP-BERT can achieve good performance with privacy guarantees in language tasks. There are still two problems: a large gap between non-private accuracy and private accuracy, and computation inefficiency of clipping every sample gradient in DP-SGD. In order to mitigate these issues, [Anil et al. \(2021\)](#) later privatizes the Adam optimizer to improve the performance. Instead of adding noise and clipping every entry in every batch in DP-SGD, it selects a pre-defined number of samples randomly and sums the clipped gradients of these selected samples, then it updates average gradients with Gaussian noise adding the sum in each batch. Besides, it also uses an increasing batch size schedule instead of a fixed one. It finds that large batch size can improve accuracy and the increasing batch size schedule can improve training efficiency. ([Senge et al., 2022](#)) recently studied five different typical NLP tasks with varying complexity using modern neural models based on BERT and XtremeDistil architectures. They showed that to achieve adequate performance, each task and privacy regime requires special treatment.

Besides BERT, [Ponomareva et al. \(2022\)](#) privately pre-train T5 ([Raffel et al., 2020](#)) via their proposed private tokenizer called DP-SentencePiece and DP-SGD. They show that DP-T5 does not suffer a large drop in pre-training utility, nor in training speed, and can still be fine-tuned to high accuracy on downstream tasks

3.2 DP Fine-tuning

Besides training pre-trained models using DP algorithms, another direction is how to fine-tune pre-trained models privately. Here the main difference is that we assume the pre-trained models such as BERT have been trained with some public data and our goal is to privately fine-tune targeting specific downstream tasks that involve sensitive data. It is noted that in this section we also include some related work on training shallow neural networks in DP such as RNN or LSTM such as ([Li et al., 2022](#); [Amid et al., 2022](#)) as these methods can be directly applied to DP fine-tuning.

In this topic, the first direction is to investigate different tasks in the DP model and to compare its performance compared to the non-private one for studying the utility-privacy tradeoff. [Yue et al. \(2022\)](#) consider the task of synthetic text generation and show that simply fine-tuning a pre-trained GPT2 with the vanilla DP-SGD enables the model to generate useful synthetic text. [Mireshghallah et al. \(2022\)](#) recently extended to generating latent semantic parses in the DP model and then generating utterances based on the parses. [Carranza et al. \(2023\)](#) use DP-SGD to fine-tune a publicly pre-trained LLM on a query generation task. The resulting model can generate private synthetic queries representative of the original queries which can be freely shared for downstream non-private recommendation training procedures. Very recently, [Lee and Sogaard \(2023\)](#) adopted the DP-SGD to the meeting summarization task and showed that DP can improve performance when evaluated on unseen meeting types. [Aziz et al. \(2022\)](#) use GPT-2 and DP-SGD based methods to generate synthetic EHR data which can de-identify sensitive information for clinical text. [Wunderlich et al. \(2021\)](#) study the hierarchical text classification task and they use DP-SGD to Bag of Words (BoW), CNNs and Transformer-based architectures. They find that Transformer-based models achieve better performance than CNN-based models in large datasets while CNN-based models are superior to Transformer-based models in small datasets.

The second direction is to reduce the huge memory cost of storing individual gradients, and decrease the added noise suffering notorious dimensional dependence in DP-SGD. Specifically, the studies in this direction always propose a general method for DP-SGD and then perform the method for different NLP tasks. [Yu et al. \(2021\)](#) propose a variant of DP-SGD called the Reparametrized Gradient Perturbation (RGP) method. The framework of RGP parametrizes each weight matrix with two low-rank carrier matrices and a residual weight matrix, which will be used to approximate the original one. Such a way can reduce the memory cost for computing individual gradient matrices and can maintain the optimization process via forward/backward signals. Later, based on RGP, [Yu et al. \(2022\)](#) show that advanced parameter-efficient methods such as ([Houlsby et al., 2019](#); [Karimi Mahabadi et al., 2021](#)) can lead to simpler and significantly improved algorithms for private

385 fine-tuning. Instead of DP-SGD, [Du and Mi \(2021\)](#)
386 propose a DP version of Forward-Propagation.
387 Specifically, it clips representations followed by
388 noise addition in the forward propagation stage.

389 Besides adapting the optimization method in
390 vanilla DP-SGD, there are also some works on
391 modifying the clipping operation or the fine-tuning
392 method directly to save the memory cost. [Li et al.](#)
393 [\(2021\)](#) propose a memory-saving technique that
394 allows clipping in DP-SGD for fine-tuning to run
395 without instantiating per-example gradients for any
396 linear layer in the model. The technique enables
397 private training Transformers with almost the same
398 memory cost as non-private training at a modest
399 run-time overhead. [Dupuy et al. \(2021\)](#) propose
400 another variant of DP-SGD via micro-batch com-
401 putations per GPU and noise decay and apply it
402 to fine-tuning models. Specifically, they scale gra-
403 dients in each micro-batch and set a decreasing
404 noise multiplier with epoch. Then, they add scaled
405 Gaussian noise to gradients. In this way, they can
406 make the training more faster and adapt it for GPU
407 training. [Bu et al. \(2023\)](#) develop a novel Book-
408 Keeping (BK) technique that implements existing
409 DP optimizers, with a substantial improvement on
410 the computational cost while also keeping almost
411 the same accuracy as DP-SGD. [Gupta et al. \(2023\)](#)
412 propose a novel language transformer finetuning
413 strategy that introduces task-specific parameters in
414 multiple transformer layers. They show that the
415 method of combining RGP and their novel strat-
416 egy is more suitable to low-resource applications.
417 [Bu et al. \(2022\)](#) privatize the bias-term fine-tuning
418 (BiTFiT) and show that DP-BiTFiT matches the
419 state-of-the-art accuracy for DP algorithms and the
420 efficiency of the standard BiTFiT ([Zaken et al.,](#)
421 [2022](#)). [Igamberdiev and Habernal \(2021\)](#) apply
422 DP-Adam in Graph Convolutional Networks to per-
423 form the private fine-tuning for text classification.
424 Specifically, they first split the graph into discon-
425 nected sub-graphs and then add noise to gradients.

426 Rather than reducing the memory cost, there are
427 some papers considering developing variants of
428 DP-SGD method to improve the performance. For
429 example, [Xia et al. \(2023\)](#) propose a per-sample
430 adaptive clipping algorithm, which is a new per-
431 spective and orthogonal to dynamic adaptive noise
432 and coordinate clipping methods. [Behnia et al.](#)
433 [\(2022\)](#) use the Edgeworth accountant ([Wang et al.,](#)
434 [2022](#)) to compute the amount of noise that is re-
435 quired to be added to the gradients in SGD to guar-

436 antee a certain privacy budget, which is lower than
437 the original DP-SGD. [Li et al. \(2022\)](#); [Amid et al.](#)
438 [\(2022\)](#) propose new private optimization methods
439 under the setting where there are some public and
440 non-sensitive data.

441 The last direction is to relax the definition of
442 DP and propose new DP-SGD variants. [Shi et al.](#)
443 [\(2021\)](#) tailor DP-SGD to SDP. Their method SDP-
444 SGD first splits the text into the sensitive and non-
445 sensitive parts, and apply normal SGD to the non-
446 sensitive part while applying DP-SGD to the sensi-
447 tive part respectively. Later, [Shi et al. \(2022\)](#) extend
448 to large language models and propose a method
449 namely Just Fine-tune Twice to private fine-tuning
450 with the guarantee of SDP.

451 3.3 Federated Learning Setting

452 In the previous parts, we reviewed the related work
453 on DP pre-trained models and DP fine-tuning mod-
454 els. Note that all the previous work only considers
455 the central DP setting where all the training data
456 samples are already collected before training, in-
457 dicating that these methods cannot be applied to
458 the federated learning (FL) setting. Compared to
459 central DP, there are fewer studies on DP Federated
460 Learning for NLP. [McMahan et al. \(2018\)](#) apply
461 DP-SGD in the FedAvg algorithm to protect user-
462 level privacy for LSTM and RNN architectures in
463 the federated learning setting. Specifically, they
464 first sample users with some probability, and then
465 add Gaussian noise to model updates of the sam-
466 pled users on the server side. Based on this, [Ra-](#)
467 [maswamy et al. \(2020\)](#) develop the first consumer-
468 scale next-word prediction model.

469 Rather than adopting DP-SGD, [Kairouz et al.](#)
470 [\(2021\)](#) provide a new paradigm for DP-FL by using
471 the Follow-The-Regularized-Leader (FTRL) algo-
472 rithm, which achieves state-of-the-art performance,
473 which is recently improved by [Choquette-Choo](#)
474 [et al. \(2022\)](#); [Koloskova et al. \(2023\)](#); [Denisov et al.](#)
475 [\(2022\)](#); [Agarwal et al. \(2021\)](#).

476 It is notable that all the previous studies only
477 consider shallow neural networks such as RNN
478 and LSTM and do not consider the large language
479 model. Until very recently, there have been some
480 papers studying DP-FL fine-tuning. For example,
481 [Wang et al. \(2023\)](#) consider the cross-device setting
482 and use DP-FTRL to privately fine-tune. Moreover,
483 they propose a distribution matching algorithm that
484 leverages both private on-device LMs and public
485 LLMs to select public records close to private data
486 distribution. [Xu et al. \(2023\)](#) deploy DP-FL ver-
487 sions of Gboard Language Models ([Hard et al.,](#)

2018) via DP-FTRL and quantile-based clip estimation method in Andrew et al. (2021).

4 Embedding Vector Perturbation Based Methods

Generally speaking, this type of approach considers to privatize the embedding vector for each token. Specifically, in this framework, the text data is first transformed into a vector (text representation) via some word embedding method such as Word2Vec (Mikolov et al., 2013) and BERT. Then we use some DP mechanism to privatize each representation and train NLP models based on these privatized text representations. Due to the post-processing property of DP, we can see the main strength of this approach is any further training on these private embeddings also preserves the DP property, while gradient perturbation based methods heavily rely on the network structure. We can see that the main step of this method is to design the best private text representation. Note that since we need to privatize each embedding representation separately, the whole algorithm could be considered as an LDP algorithm and thus it can also be used in the LDP setting. It is also notable that different studies may consider different notions and levels of privacy. In fact, most of the existing work considers the word level of privacy.

4.1 Vanilla DP

The most direct approach is to design private embedding mechanisms that satisfy the standard DP. Lyu et al. (2020b) first study this problem and they propose a framework. Specifically, firstly for each word the embedding module of such framework outputs a 1-dimensional real representation with length r , then it privatizes the vector via a variant of the Unary Encoding mechanism in (Wang et al., 2017). In order to remove the dependence of dimensionality in the Unary Encoding mechanism, they propose an Optimized Multiple Encoding, which embeds vectors with a certain fixed size. Their post-processing procedure was then improved by (Plant et al., 2021). In (Plant et al., 2021), it first gets the final layer representation of the pre-trained model for each token, then normalizes it with sequence and adds Laplacian noise, and finally trains this classifier with adversarial training. To further improve the fairness for the downstream tasks on private embedding, later Lyu et al. (2020a) propose to dropout perturbed embeddings to amplify privacy and a robust training algorithm that incorporates the noisy training representation in the training pro-

cess to derive a robust target model, which also reduces model discrimination in most cases.

Krishna et al. (2021); Habernal (2021); Alnasser et al. (2021) also study privatizing word embeddings. However, instead of using the Unary Encoding mechanism or dropout, Krishna et al. (2021); Alnasser et al. (2021) propose ADePT which is an auto-encoder-based DP algorithm. Let \mathbf{u} be the input, an auto-encoder model consists of an encoder that returns a vector representation $\mathbf{r} = \text{Enc}(\mathbf{u})$ for the input \mathbf{u} , which is then passed into the decoder to construct an output $\mathbf{v} = \text{Dec}(\mathbf{r})$. In (Krishna et al., 2021), it first normalized the word embedded vector by some parameter C i.e., $w = \text{Enc}(\mathbf{u}) \min\{1, \frac{C}{\|\text{Enc}(\mathbf{u})\|_2}\}$, then it adds Laplacian noise to the normalized vector w and get \mathbf{r} . Unfortunately, Habernal (2021) points out that ADePT is not differentially private by thorough theoretical proof. The problem of ADePT lies in the sensitivity calculation and could be remedied by adding calibrated noise or tighter bounded clipping norm. Later, Igamberdiev et al. (2022) provide the source code of DP Auto-Encoder methods to improve reproducibility. Recently, Maheshwari et al. (2022) propose a method that combines differential privacy and adversarial training techniques to solve the privacy-fairness-accuracy trade-off in local DP. In their framework, first, the input text will be fed into encoders, then it will be normalized and privatized by using the Laplacian mechanism. Next, it will be fed into a normal classifier and adversarial training separately to combine a loss that contains normal classification loss and adversarial loss. They find that the model can improve privacy and fairness simultaneously. To further improve the performance, (Bollegala et al., 2023) propose a Neighbourhood-Aware Differential Privacy (NADP) mechanism considering the neighborhood of a word in a pretrained static word embedding space to determine the minimal amount of noise required to guarantee a specified privacy level.

Besides the work on word-level privacy we mentioned above, recently there have been some works studying sentence-level and token-level private embeddings. Meehan et al. (2022) propose a method namely DeepCandidate to achieve sentence-level privacy. They first put public and private sentences into a sentence encoder to get sentence embeddings. Then, they use a method namely DeepCandidate to choose the candidate sentence embeddings that are near to private embeddings. Finally, they use some

DP mechanism to sample from the candidate embeddings for each private embedding. This method somehow solves the challenge of the sentence-level privacy problem by taking advantage of clustering in differential privacy. (Du et al., 2023b) consider sentence-level privacy for private fine-tuning and propose DP-Forward fine-tuning, which perturbs the forwardpass embeddings of every user’s (labeled) sequence. However, it is notable that they consider a variant of LDP called sequence local DP. Chen et al. (2023) propose a novel Customized Text (CusText) sanitization mechanism that provides more advanced privacy protection at the token level.

4.2 Metric DP

In Metric DP for text data, each sample of the input can be represented as a string x with at most l words, thus the data universe will be W^l where W is a dictionary. Also we assume that there is a word embedding model $\phi : W \mapsto \mathbb{R}^n$ and its associated distance $d(x, x') = \sum_{i=1}^l \|\phi(w_i) - \phi(w'_i)\|_2$, where $x = w_1 w_2 \dots w_l$ and $x' = w'_1 w'_2 \dots w'_l$ are two samples. Thus, the goal is to design a mechanism for each $\phi(w_i)$ with the guarantee of Metric DP. Since we aim to randomize each $\phi(w_i)$ for each sample. The whole algorithm is also suitable for local metric DP with word-level privacy.

Feyisetan et al. (2020) first study this problem. Generally speaking, their mechanism consists of two steps. The first step is perturbation, we add some noise N to text vector $\phi(w_i)$ to ensure ϵ -LDP, where N has then density probability function $p_N(z) \propto \exp(-\epsilon\|z\|_2)$. The main issue of this approach is that after the perturbation, $\hat{\phi}_i$ may be inconsistent with the word embedding. That is, there may not exist a word u such that $u = \hat{\phi}_i$. Thus, to address this issue, we need to project the perturbed vector into the embedding space. That is the second step. Feyisetan et al. (2020) show that the algorithm is ϵ -local Metric DP.

Note that the method was later improved from different aspects. For example, Xu et al. (2020) reconsider the problem setting and they observe that the distance used in (Feyisetan et al., 2020) is the Euclidean norm $d(x, x') = \sum_{i=1}^l \|\phi(w_i) - \phi(w'_i)\|_2$, which cannot describe the similarity between two words in the embedding space. To address the issue, they propose to use the Mahalanobis Norm and modify the algorithm by using the Mahalanobis mechanism, which can improve performance. To further improve the utility in the projection step, Xu et al. (2021b) further

propose the Vickrey mechanism in case the first nearest neighbors are the original input or some rare words need large-scale noise to perturb and hard to find the corresponding words. In order to solve this problem, they use a hyperparameter in their algorithm to adjust the selection of the first and second nearest neighbors (words). To further allow a smaller range of nearby words to be considered than the multivariate Laplace mechanism, (Xu et al., 2021a; Carvalho et al., 2021b) propose an improved perturbation method via the Truncated Gumbel Noise. To further address the high dimensional issue, Feyisetan and Kasiviswanathan (2021) uses the random projection for the original text representation to a lower dimensional space and then projects back to the original space after adding random noise to preserve DP. Besides, Feyisetan et al. (2019) define the hyperbolic embeddings and use the Metropolis-Hastings (MH) algorithm to sample from hyperbolic distribution. However, it is remarkable that if we consider the LDP setting, then all the previous methods need to send real numbers to the server, which has a high communication cost. To address the issue, Carvalho et al. (2021a) propose to use the binary randomized response mechanism by using binary embedding vectors. Recently, Tang et al. (2020) consider the case where different words may have different levels of privacy. They first divide the word into two types, and then add corresponding noise according to different levels of privacy. Imola et al. (2022) recently proposed an optimal Meric DP mechanism for finite vocabulary, they then provided an algorithm that could quickly calculate the mechanism. Finally, they applied it to private word embedding. Instead of developing new private mechanisms, there are also some studies on improving the embedding process. The previous metric DP mechanisms are expected to fall short of finding substitutes for words with ambiguous meanings. Address these ambiguous words, Arnold et al. (2023a) provide a sense embedding and incorporate a sense disambiguation step prior to noise injection. Arnold et al. (2023b) account the common semantic context issue that appeared in the previous private embedding mechanisms. They incorporate grammatical categories into the privatization step in the form of a constraint to the candidate selection and show that selecting a substitution with matching grammatical properties amplifies the performance in downstream tasks. Qu et al. (2021) recently points out that (Lyu et al.,

2020a) does not address privacy issues in the training phase since the server needs users’ raw data to fine-tuning. Moreover, its method has a high computational cost due to the heavy encoder workload on the user side. Thus, Qu et al. (2021) improve it and consider the federated setting where users send their privatized samples via some local metric DP mechanism to the server and the server conducts privacy-constrained fine-tuning methods. Moreover, besides the text-to-text privatization given in (Feyisetan et al., 2020) and the sequence private representation proposed by Lyu et al. (2020a), Qu et al. (2021) proposed new token-level privatization and text-to-text privatization methods. In the token representation privatization method, they add random noise using metric DP to token embedding and send it to the server. They add noise to the embedded token and output the closest neighbor token in the embedding space.

Instead of the local Metric DP, Yue et al. (2021) consider UMLDP and propose SANTEXT and SANTEXT+ algorithms for text sanitization tasks. Specifically, they divide all the text into a sensitive token set \mathcal{V}_S and a remaining token set \mathcal{V}_N . Then \mathcal{V}_S and \mathcal{V}_N will use a privacy budget of ϵ and ϵ_0 respectively via the composition theorem in LDP. After deriving token vectors, SANTEXT samples new tokens via local Metric DP with Euclidean distance. Compared with SANTEXT, SANTEXT+ samples new tokens when the original tokens are in sensitive set \mathcal{V}_S . They apply it to BERT pre-training and fine-tuning models.

While there are many studies on the benefits of private embedding with word-level privacy. There are also some shortcomings to such notion of privacy, as mentioned by (Mattern et al., 2022) recently. For example, in the previous private word embedding methods we need to assume the length of the string for each sample is the same. Moreover, since we consider word level of privacy, the total privacy budget will grow linearly with the length of the sample. To mitigate some shortcomings, Mattern et al. (2022) propose an alternative text anonymization method based on fine-tuning of large language models for paraphrasing. To ensure DP, they adopt the exponential mechanism to sample from the softmax distribution. They apply their method in fine-tuning models with GPT-2.

Recently Du et al. (2023a) studied sentence-level private embedding in local metric DP. Borrowing the wisdom of normalizing sentence embedding for

robustness, they impose a consistency constraint on their sanitization. They propose two instantiations from the Euclidean and angular distances. The first one utilizes the Purkayastha mechanism (Weggenmann and Kerschbaum, 2021) and the other is upgraded from the generalized planar Laplace mechanism with post-processing.

5 Challenges and Future Directions

Large-scale Training. Dealing with large-scale text data and training large models like GPT-3 are tough tasks in deep learning with DP. Due to the high dimensionality of embedding vectors, even adding small noise can have a significant influence on the training speed and performance of models. It is more severe for DP-SGD-based methods, which need high memory cost and their per-example clipping procedure is time-consuming. These methods will be inefficient when they are applied to large language models. Thus, how to reduce the memory cost and accelerate the training of DP-SGD become core concerns in gradient perturbation-based methods. Although there is some work in this direction, there is still a gap in accuracy between private and non-private models and these methods still need catastrophic cost of memory compared with the non-private ones. Moreover, it is well known that we need a heavy workload on hyperparameter-tuning for large-scale models in the non-private case. From the privacy view, each try-on hyperparameter-tuning will cost an additional privacy budget, which makes our final private model cost a large privacy budget. Thus, how to efficiently and privately tune the hyperparameters in large models is challenging.

Private Inference. It is notable that in this paper we mainly discussed how to privately train and release a language model without leaking information about training data. However, in some scenarios (such as Machine Learning as a Service) we only want to use the model for inference instead of releasing the model. Thus, for these scenarios, we only need to perform inference tasks based on our trained model while we do not want to leak information of training data. From the DP side, such private inference corresponds to the DP prediction algorithm, which is proposed by (Dwork and Feldman, 2018). Compared with private training, DP inference for text data is still far from well-understood and there is only few studies on it (Ginart et al., 2022; Weggenmann et al., 2022a; Majmudar et al., 2022; Zhou et al., 2023; Li et al., 2023).

796 Limitations

797 First, in this paper, we mainly focused on the deep
798 learning-based models for NLP tasks in the differ-
799 ential privacy model. Actually, there are also some
800 studies on classical statistical models or approaches
801 for NLP in DP, such as topic modeling (Park et al.,
802 2016; Zhao et al., 2021; Huang and Chen, 2021)
803 and n-gram extraction (Kim et al., 2021). Secondly,
804 due to the space limit, we did not discuss all the
805 related work for DP-SGD and we only focused on
806 the work that uses DP-SGD to NLP-related tasks.
807 Thirdly, while we tried our best to discuss all the
808 existing work on deep learning-based methods for
809 DP-NLP, we have to say we may missed some re-
810 lated work. Moreover, since we aim to classify all
811 the current work into two categories based on their
812 methods of adding randomness, there is still some
813 work that does not belong to these two classes, such
814 as (Bo et al., 2021; Weggenmann et al., 2022b; Tian
815 et al., 2022; Duan et al., 2023; Tang et al., 2023;
816 Wu et al., 2023). To make our paper be consist-
817 ent, we did not mention these work here. Fourthly,
818 although DP can provide rigorous guarantees on
819 privacy-preserving, it also has been shown that DP
820 machine learning models can cause fairness issues.
821 For example, they always have a disparate impact
822 on model accuracy (Bagdasaryan et al., 2019). Fi-
823 nally, it is notable that in this paper we did not
824 discuss the narrow assumptions made by differen-
825 tial privacy, and the broadness of natural language
826 and of privacy as a social norm. More details can
827 be found in (Brown et al., 2022).

828 References

829 Martín Abadi, Andy Chu, Ian J. Goodfellow, H. Bren-
830 dan McMahan, Ilya Mironov, Kunal Talwar, and
831 Li Zhang. 2016. [Deep learning with differential pri-
832 vacy](#). In *Proceedings of the 2016 ACM SIGSAC
833 Conference on Computer and Communications Se-
834 curity, Vienna, Austria, October 24-28, 2016*, pages
835 308–318. ACM.

836 Naman Agarwal, Peter Kairouz, and Ziyu Liu. 2021.
837 The skellam mechanism for differentially private fed-
838 erated learning. *Advances in Neural Information
839 Processing Systems*, 34:5052–5064.

840 Walaa Alnasser, Ghazaleh Beigi, and Huan Liu. 2021.
841 [Privacy preserving text representation learning using
842 BERT](#). In *Social, Cultural, and Behavioral Modeling
843 - 14th International Conference, SBP-BRiMS 2021,
844 Virtual Event, July 6-9, 2021, Proceedings*, volume
845 12720 of *Lecture Notes in Computer Science*, pages
846 91–100. Springer.

Ehsan Amid, Arun Ganesh, Rajiv Mathews, Swa-
847 roop Ramaswamy, Shuang Song, Thomas Steinke,
848 Vinith M Suriyakumar, Om Thakkar, and Abhradeep
849 Thakurta. 2022. [Public data-assisted mirror descent
850 for private model training](#). In *International Confer-
851 ence on Machine Learning*, pages 517–535. PMLR.
852

Galen Andrew, Om Thakkar, Brendan McMahan, and
853 Swaroop Ramaswamy. 2021. [Differentially private
854 learning with adaptive clipping](#). *Advances in Neural
855 Information Processing Systems*, 34:17455–17466.
856

Rohan Anil, Badih Ghazi, Vineet Gupta, Ravi Kumar,
857 and Pasin Manurangsi. 2021. [Large-scale differen-
858 tially private BERT](#). *CoRR*, abs/2108.01624.
859

Stefan Arnold, Dilara Yesilbas, and Sven Weinzierl.
860 2023a. [Driving context into text-to-text privatization](#).
861 *CoRR*, abs/2306.01457.
862

Stefan Arnold, Dilara Yesilbas, and Sven Weinzierl.
863 2023b. [Guiding text-to-text privatization by syntax](#).
864 *CoRR*, abs/2306.01471.
865

Shahab Asoodeh, Jiachun Liao, Flávio P. Calmon,
866 Oliver Kosut, and Lalitha Sankar. 2021. [Three
867 variants of differential privacy: Lossless conversion
868 and applications](#). *IEEE J. Sel. Areas Inf. Theory*,
869 2(1):208–222.
870

Md Momin Al Aziz, Tanbir Ahmed, Tasnia Faequa,
871 Xiaoqian Jiang, Yiyu Yao, and Noman Mohammed.
872 2022. [Differentially private medical texts genera-
873 tion using generative neural networks](#). *ACM Trans.
874 Comput. Heal.*, 3(1):5:1–5:27.
875

Eugene Bagdasaryan, Omid Poursaeed, and Vitaly
876 Shmatikov. 2019. [Differential privacy has disparate
877 impact on model accuracy](#). In *Advances in Neural
878 Information Processing Systems 32: Annual Confer-
879 ence on Neural Information Processing Systems 2019,
880 NeurIPS 2019, December 8-14, 2019, Vancouver, BC,
881 Canada*, pages 15453–15462.
882

Borja Balle, Gilles Barthe, and Marco Gaboardi. 2018.
883 [Privacy amplification by subsampling: Tight anal-
884 yses via couplings and divergences](#). In *Advances
885 in Neural Information Processing Systems 31: An-
886 nual Conference on Neural Information Processing
887 Systems 2018, NeurIPS 2018, December 3-8, 2018,
888 Montréal, Canada*, pages 6280–6290.
889

Borja Balle, Gilles Barthe, Marco Gaboardi, and Joseph
890 Geumlek. 2019. [Privacy amplification by mixing
891 and diffusion mechanisms](#). In *Advances in Neural
892 Information Processing Systems 32: Annual Confer-
893 ence on Neural Information Processing Systems 2019,
894 NeurIPS 2019, December 8-14, 2019, Vancouver, BC,
895 Canada*, pages 13277–13287.
896

Borja Balle, Peter Kairouz, Brendan McMahan, Om Di-
897 pakbhai Thakkar, and Abhradeep Thakurta. 2020.
898 [Privacy amplification via random check-ins](#). In *Ad-
899 vances in Neural Information Processing Systems 33:
900 Annual Conference on Neural Information Process-
901 ing Systems 2020, NeurIPS 2020, December 6-12,
902 2020, virtual*.
903

904	Rouzbeh Behnia, Mohammadreza Reza Ebrahimi, Jason Pacheco, and Balaji Padmanabhan. 2022. Ew-tune: A framework for privately fine-tuning large language models with differential privacy. In <i>2022 IEEE International Conference on Data Mining Workshops (ICDMW)</i> , pages 560–566. IEEE.	<i>USENIX Security 2019, Santa Clara, CA, USA, August 14-16, 2019</i> , pages 267–284. USENIX Association.	959 960 961
910	Haohan Bo, Steven H. H. Ding, Benjamin C. M. Fung, and Farkhund Iqbal. 2021. ER-AE: Differentially private text generation for authorship anonymization. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3997–4007, Online. Association for Computational Linguistics.	Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. 2021. Extracting training data from large language models . In <i>30th USENIX Security Symposium, USENIX Security 2021, August 11-13, 2021</i> , pages 2633–2650. USENIX Association.	962 963 964 965 966 967 968 969
918	Danushka Bollegala, Shuichi Otake, Tomoya Machide, and Ken-ichi Kawarabayashi. 2023. A neighbourhood-aware differential privacy mechanism for static word embeddings . <i>CoRR</i> , abs/2309.10551.	Aldo Gael Carranza, Rezsza Farahani, Natalia Ponomareva, Alex Kurakin, Matthew Jagielski, and Milad Nasr. 2023. Privacy-preserving recommender systems with synthetic query generation using differentially private large language models. <i>arXiv preprint arXiv:2305.05973</i> .	970 971 972 973 974 975
923	Hannah Brown, Katherine Lee, Fatemehsadat Miresghallah, Reza Shokri, and Florian Tramèr. 2022. What does it mean for a language model to preserve privacy? In <i>FACt '22: 2022 ACM Conference on Fairness, Accountability, and Transparency, Seoul, Republic of Korea, June 21 - 24, 2022</i> , pages 2280–2292. ACM.	Ricardo Silva Carvalho, Theodore Vasiloudis, and Oluwaseyi Feyisetan. 2021a. BRR: preserving privacy of text data efficiently on device . <i>CoRR</i> , abs/2107.07923.	976 977 978 979
930	Zhiqi Bu, Jinshuo Dong, Qi Long, and Weijie J. Su. 2019. Deep learning with gaussian differential privacy . <i>CoRR</i> , abs/1911.11607.	Ricardo Silva Carvalho, Theodore Vasiloudis, and Oluwaseyi Feyisetan. 2021b. TEM: high utility metric differential privacy on text . <i>CoRR</i> , abs/2107.07928.	980 981 982 983
933	Zhiqi Bu, Yu-Xiang Wang, Sheng Zha, and George Karypis. 2022. Differentially private bias-term only fine-tuning of foundation models. <i>arXiv preprint arXiv:2210.00036</i> .	Konstantinos Chatzikokolakis, Miguel E. Andrés, Nicolás Emilio Bordenabe, and Catuscia Palamidessi. 2013. Broadening the scope of differential privacy using metrics . In <i>Privacy Enhancing Technologies - 13th International Symposium, PETS 2013, Bloomington, IN, USA, July 10-12, 2013. Proceedings</i> , volume 7981 of <i>Lecture Notes in Computer Science</i> , pages 82–102. Springer.	984 985 986 987 988 989 990 991
937	Zhiqi Bu, Yu-Xiang Wang, Sheng Zha, and George Karypis. 2023. Differentially private optimization on large model at small cost. In <i>International Conference on Machine Learning</i> , pages 3192–3218. PMLR.	Sai Chen, Fengran Mo, Yanhao Wang, Cen Chen, Jian-Yun Nie, Chengyu Wang, and Jamie Cui. 2023. A customized text sanitization mechanism with differential privacy . In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 5747–5758. Association for Computational Linguistics.	992 993 994 995 996 997 998
942	Mark Bun, Cynthia Dwork, Guy N. Rothblum, and Thomas Steinke. 2018. Composable and versatile privacy via truncated CDP . In <i>Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018, Los Angeles, CA, USA, June 25-29, 2018</i> , pages 74–86. ACM.	Albert Cheu, Adam D. Smith, Jonathan R. Ullman, David Zeber, and Maxim Zhilyaev. 2019. Distributed differential privacy via shuffling . In <i>Advances in Cryptology - EUROCRYPT 2019 - 38th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Darmstadt, Germany, May 19-23, 2019, Proceedings, Part I</i> , volume 11476 of <i>Lecture Notes in Computer Science</i> , pages 375–403. Springer.	999 1000 1001 1002 1003 1004 1005 1006 1007
948	Mark Bun and Thomas Steinke. 2016. Concentrated differential privacy: Simplifications, extensions, and lower bounds . In <i>Theory of Cryptography - 14th International Conference, TCC 2016-B, Beijing, China, October 31 - November 3, 2016, Proceedings, Part I</i> , volume 9985 of <i>Lecture Notes in Computer Science</i> , pages 635–658.	Christopher A Choquette-Choo, H Brendan McMahan, Keith Rush, and Abhradeep Thakurta. 2022. Multi-epoch matrix factorization mechanisms for private machine learning. <i>arXiv preprint arXiv:2211.06530</i> .	1008 1009 1010 1011
955	Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. 2019. The secret sharer: Evaluating and testing unintended memorization in neural networks . In <i>28th USENIX Security Symposium</i> ,	Edwige Cyffers and Aurélien Bellet. 2022. Privacy amplification by decentralization . In <i>International Conference on Artificial Intelligence and Statistics</i> ,	1012 1013 1014

1015	<i>AISTATS 2022, 28-30 March 2022, Virtual Event</i> ,	Cynthia Dwork, Frank McSherry, Kobbi Nissim, and	1071
1016	volume 151 of <i>Proceedings of Machine Learning</i>	Adam D. Smith. 2006. Calibrating noise to sensitiv-	1072
1017	<i>Research</i> , pages 5334–5353. PMLR.	ity in private data analysis . In <i>Theory of Cryptogra-</i>	1073
1018	Fida Kamal Dankar and Khaled El Emam. 2013. Prac-	<i>phy, Third Theory of Cryptography Conference, TCC</i>	1074
1019	ticing differential privacy in health care: A review .	2006, New York, NY, USA, March 4-7, 2006, <i>Pro-</i>	1075
1020	<i>Trans. Data Priv.</i> , 6(1):35–67.	<i>ceedings</i> , volume 3876 of <i>Lecture Notes in Computer</i>	1076
1021	Sergey Denisov, H Brendan McMahan, John Rush,	<i>Science</i> , pages 265–284. Springer.	1077
1022	Adam Smith, and Abhradeep Guha Thakurta. 2022.	Cynthia Dwork and Aaron Roth. 2014. The algorithmic	1078
1023	Improved differential privacy for sgd via optimal pri-	foundations of differential privacy . <i>Founda-</i>	1079
1024	private linear operators on adaptive streams. <i>Advances</i>	<i>tions and Trends® in Theoretical Computer Science</i> ,	1080
1025	<i>in Neural Information Processing Systems</i> , 35:5910–	9(3–4):211–407.	1081
1026	5924.	Cynthia Dwork and Guy N. Rothblum. 2016. Concen-	1082
1027	Damien Desfontaines and Balázs Pejő. 2020. Sok: Dif-	trated differential privacy . <i>CoRR</i> , abs/1603.01887.	1083
1028	ferential privacies . <i>Proceedings on Privacy Enhanc-</i>	Cynthia Dwork, Guy N. Rothblum, and Salil P. Vadhan.	1084
1029	<i>ing Technologies</i> , 2020(2):288–313.	2010. Boosting and differential privacy . In <i>51th An-</i>	1085
1030	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	<i>annual IEEE Symposium on Foundations of Computer</i>	1086
1031	Kristina Toutanova. 2019. BERT: pre-training of	<i>Science, FOCS 2010, October 23-26, 2010, Las Ve-</i>	1087
1032	deep bidirectional transformers for language under-	<i>gas, Nevada, USA</i> , pages 51–60. IEEE Computer	1088
1033	standing . In <i>Proceedings of the 2019 Conference of</i>	Society.	1089
1034	<i>the North American Chapter of the Association for</i>	Oluwaseyi Feyisetan, Borja Balle, Thomas Drake, and	1090
1035	<i>Computational Linguistics: Human Language Tech-</i>	Tom Diethe. 2020. Privacy- and utility-preserving	1091
1036	<i>nologies, NAACL-HLT 2019, Minneapolis, MN, USA,</i>	textual analysis via calibrated multivariate perturba-	1092
1037	<i>June 2-7, 2019, Volume 1 (Long and Short Papers),</i>	tions . In <i>WSDM '20: The Thirteenth ACM Interna-</i>	1093
1038	pages 4171–4186. Association for Computational	<i>tional Conference on Web Search and Data Mining,</i>	1094
1039	Linguistics.	<i>Houston, TX, USA, February 3-7, 2020</i> , pages 178–	1095
1040	Jinshuo Dong, Aaron Roth, and Weijie J. Su. 2022.	186. ACM.	1096
1041	Gaussian differential privacy . <i>Journal of the Royal</i>	Oluwaseyi Feyisetan, Tom Diethe, and Thomas Drake.	1097
1042	<i>Statistical Society: Series B (Statistical Methodol-</i>	2019. Leveraging hierarchical representations for	1098
1043	<i>ogy)</i> , 84(1):3–37.	preserving privacy and utility in text . In <i>2019 IEEE</i>	1099
1044	Jian Du and Haitao Mi. 2021. Dp-fp: Differentially	<i>International Conference on Data Mining, ICDM</i>	1100
1045	private forward propagation for large models . <i>arXiv</i>	2019, Beijing, China, November 8-11, 2019, pages	1101
1046	<i>preprint arXiv:2112.14430</i> .	210–219. IEEE.	1102
1047	Minxin Du, Xiang Yue, Sherman S. M. Chow, and Huan	Oluwaseyi Feyisetan and Shiva Kasiviswanathan. 2021.	1103
1048	Sun. 2023a. Sanitizing sentence embeddings (and	Private release of text embedding vectors. In <i>Pro-</i>	1104
1049	labels) for local differential privacy . In <i>Proceedings</i>	<i>ceedings of the First Workshop on Trustworthy Natu-</i>	1105
1050	<i>of the ACM Web Conference 2023, WWW '23</i> , page	<i>ral Language Processing</i> , pages 15–27.	1106
1051	2349–2359, New York, NY, USA. Association for	Antonio Ginart, Laurens van der Maaten, James Zou,	1107
1052	Computing Machinery.	and Chuan Guo. 2022. Submix: Practical private	1108
1053	Minxin Du, Xiang Yue, Sherman S. M. Chow, Tianhao	prediction for large-scale language models . <i>CoRR</i> ,	1109
1054	Wang, Chenyu Huang, and Huan Sun. 2023b. Dp-	abs/2201.00971.	1110
1055	forward: Fine-tuning and inference on language mod-	Sivakanth Gopi, Yin Tat Lee, and Lukas Wutschitz.	1111
1056	els with differential privacy in forward pass . <i>CoRR</i> ,	2021. Numerical composition of differential privacy .	1112
1057	abs/2309.06746.	In <i>Advances in Neural Information Processing Sys-</i>	1113
1058	Haonan Duan, Adam Dziedzic, Nicolas Papernot, and	<i>tems 34: Annual Conference on Neural Information</i>	1114
1059	Franziska Boenisch. 2023. Flocks of stochastic par-	<i>Processing Systems 2021, NeurIPS 2021, December</i>	1115
1060	rots: Differentially private prompt learning for large	<i>6-14, 2021, virtual</i> , pages 11631–11642.	1116
1061	language models . <i>CoRR</i> , abs/2305.15594.	Umang Gupta, Aram Galstyan, and Greg Ver Steeg.	1117
1062	Christophe Dupuy, Radhika Arava, Rahul Gupta, and	2023. Jointly reparametrized multi-layer adapta-	1118
1063	Anna Rumshisky. 2021. An efficient DP-SGD	tion for efficient and private tuning . <i>arXiv preprint</i>	1119
1064	mechanism for large scale NLP models . <i>CoRR</i> ,	<i>arXiv:2305.19264</i> .	1120
1065	abs/2107.14586.	Ivan Habernal. 2021. When differential privacy meets	1121
1066	Cynthia Dwork and Vitaly Feldman. 2018. Privacy-	NLP: The devil is in the detail . In <i>Proceedings of the</i>	1122
1067	preserving prediction . In <i>Conference On Learning</i>	<i>2021 Conference on Empirical Methods in Natural</i>	1123
1068	<i>Theory, COLT 2018, Stockholm, Sweden, 6-9 July</i>	<i>Language Processing</i> , pages 1522–1528, Online and	1124
1069	2018, volume 75 of <i>Proceedings of Machine Learn-</i>	Punta Cana, Dominican Republic. Association for	1125
1070	<i>ing Research</i> , pages 1693–1702. PMLR.	Computational Linguistics.	1126

1127	Andrew Hard, Kanishka Rao, Rajiv Mathews, Swaroop Ramaswamy, Françoise Beaufays, Sean Augenstein, Hubert Eichner, Chloé Kiddon, and Daniel Ramage. 2018. Federated learning for mobile keyboard prediction. <i>arXiv preprint arXiv:1811.03604</i> .	1184
1128		1185
1129		1186
1130		1187
1131		1188
1132	Shlomo Hoory, Amir Feder, Avichai Tendler, Sofia Erell, Alon Peled-Cohen, Itay Laish, Hootan Nakhost, Uri Stemmer, Ayelet Benjamini, Avinatan Hassidim, and Yossi Matias. 2021. Learning and evaluating a differentially private pre-trained language model . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 1178–1189, Punta Cana, Dominican Republic. Association for Computational Linguistics.	1189
1133		1190
1134		1191
1135		1192
1136		1193
1137		1194
1138		1195
1139		1196
1140		1197
1141	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In <i>International Conference on Machine Learning</i> , pages 2790–2799. PMLR.	1198
1142		1199
1143		
1144		1200
1145		1201
1146		1202
1147	Tao Huang and Hong Chen. 2021. Improving privacy guarantee and efficiency of latent dirichlet allocation model training under differential privacy . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021</i> , pages 143–152. Association for Computational Linguistics.	1203
1148		
1149		1204
1150		1205
1151		1206
1152		1207
1153		1208
1154		1209
1155	Timour Igamberdiev, Thomas Arnold, and Ivan Habernal. 2022. Dp-rewrite: Towards reproducibility and transparency in differentially private text rewriting . In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , page (to appear), Gyeongju, Republic of Korea. International Committee on Computational Linguistics.	1210
1156		1211
1157		1212
1158		1213
1159		1214
1160		
1161	Timour Igamberdiev and Ivan Habernal. 2021. Privacy-preserving graph convolutional networks for text classification . <i>CoRR</i> , abs/2102.09604.	1215
1162		1216
1163		1217
1164	Jacob Imola and Kamalika Chaudhuri. 2021. Privacy amplification via bernoulli sampling . <i>CoRR</i> , abs/2105.10594.	1218
1165		1219
1166		1220
1167	Jacob Imola, Shiva Prasad Kasiviswanathan, Stephen White, Abhinav Aggarwal, and Nathanael Teissier. 2022. Balancing utility and scalability in metric differential privacy . In <i>Uncertainty in Artificial Intelligence, Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence, UAI 2022, 1-5 August 2022, Eindhoven, The Netherlands</i> , volume 180 of <i>Proceedings of Machine Learning Research</i> , pages 885–894. PMLR.	1221
1168		1222
1169		1223
1170		1224
1171		1225
1172		1226
1173		1227
1174		1228
1175		1229
1176	Zhanglong Ji, Zachary Chase Lipton, and Charles Elkan. 2014. Differential privacy and machine learning: a survey and review . <i>CoRR</i> , abs/1412.7584.	1230
1177		1231
1178		1232
1179	Peter Kairouz, Brendan McMahan, Shuang Song, Om Thakkar, Abhradeep Thakurta, and Zheng Xu. 2021. Practical and private (deep) learning without sampling or shuffling. In <i>International Conference on Machine Learning</i> , pages 5213–5225. PMLR.	1233
1180		1234
1181		1235
1182		
1183		1236
	Peter Kairouz, Sewoong Oh, and Pramod Viswanath. 2015. The composition theorem for differential privacy . In <i>Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015</i> , volume 37 of <i>JMLR Workshop and Conference Proceedings</i> , pages 1376–1385. JMLR.org.	1237
		1238
		1239
	Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. 2021. Compacter: Efficient low-rank hypercomplex adapter layers. <i>Advances in Neural Information Processing Systems</i> , 34:1022–1035.	
	Daniel Kifer and Ashwin Machanavajjhala. 2011. No free lunch in data privacy . In <i>Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2011, Athens, Greece, June 12-16, 2011</i> , pages 193–204. ACM.	
	Daniel Kifer and Ashwin Machanavajjhala. 2014. Pufferfish: A framework for mathematical privacy definitions . <i>ACM Trans. Database Syst.</i> , 39(1):3:1–3:36.	
	Kunho Kim, Sivakanth Gopi, Janardhan Kulkarni, and Sergey Yekhanin. 2021. Differentially private n-gram extraction . In <i>Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual</i> , pages 5102–5111.	
	Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization . In <i>3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings</i> .	
	Oleksandra Klymenko, Stephen Meisenbacher, and Florian Matthes. 2022. Differential privacy in natural language processing the story so far . In <i>Proceedings of the Fourth Workshop on Privacy in Natural Language Processing</i> , pages 1–11, Seattle, United States. Association for Computational Linguistics.	
	Anastasia Koloskova, Ryan McKenna, Zachary Charles, Keith Rush, and Brendan McMahan. 2023. Convergence of gradient descent with linearly correlated noise and applications to differentially private learning. <i>arXiv preprint arXiv:2302.01463</i> .	
	Satyapriya Krishna, Rahul Gupta, and Christophe Dupuy. 2021. ADePT: Auto-encoder based differentially private text transformation . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 2435–2439, Online. Association for Computational Linguistics.	
	Seolhwa Lee and Anders Søgaard. 2023. Private meeting summarization without performance loss. <i>arXiv preprint arXiv:2305.15894</i> .	
	Tian Li, Manzil Zaheer, Sashank Reddi, and Virginia Smith. 2022. Private adaptive optimization with side information. In <i>International Conference on Machine Learning</i> , pages 13086–13105. PMLR.	

1240	Xuechen Li, Florian Tramer, Percy Liang, and Tatsunori Hashimoto. 2021. Large language models can be strong differentially private learners. In <i>International Conference on Learning Representations</i> .	Fatemehsadat Mireshghallah, Richard Shin, Yu Su, Tatsunori Hashimoto, and Jason Eisner. 2022. Privacy-preserving domain adaptation of semantic parsers. <i>arXiv preprint arXiv:2212.10520</i> .	1294
1241			1295
1242			1296
1243			1297
1244	Yansong Li, Zhixing Tan, and Yang Liu. 2023. Privacy-preserving prompt tuning for large language model services. <i>CoRR</i> , abs/2305.06212.	Ilya Mironov. 2017. Rényi differential privacy. In <i>30th IEEE Computer Security Foundations Symposium, CSF 2017, Santa Barbara, CA, USA, August 21-25, 2017</i> , pages 263–275. IEEE Computer Society.	1298
1245			1299
1246			1300
1247	Lingjuan Lyu, Xuanli He, and Yitong Li. 2020a. Differentially private representation for NLP: Formal guarantee and an empirical study on privacy and fairness. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2355–2365, Online. Association for Computational Linguistics.	Ilya Mironov, Kunal Talwar, and Li Zhang. 2019. Rényi differential privacy of the sampled gaussian mechanism. <i>CoRR</i> , abs/1908.10530.	1302
1248			1303
1249			1304
1250			1305
1251			1306
1252			1307
1253	Lingjuan Lyu, Yitong Li, Xuanli He, and Tong Xiao. 2020b. Towards differentially private text representations. In <i>Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020</i> , pages 1813–1816. ACM.	Takao Murakami and Yusuke Kawamoto. 2019. Utility-optimized local differential privacy mechanisms for distribution estimation. In <i>28th USENIX Security Symposium, USENIX Security 2019, Santa Clara, CA, USA, August 14-16, 2019</i> , pages 1877–1894. USENIX Association.	1308
1254			1309
1255			1310
1256			1311
1257			1312
1258			1313
1259	Gaurav Maheshwari, Pascal Denis, Mikaela Keller, and Aurélien Bellet. 2022. Fair nlp models with differentially private text encoders. <i>arXiv preprint arXiv:2205.06135</i> .	Jack Murtagh and Salil P. Vadhan. 2016. The complexity of computing the optimal composition of differential privacy. In <i>Theory of Cryptography - 13th International Conference, TCC 2016-A, Tel Aviv, Israel, January 10-13, 2016, Proceedings, Part I</i> , volume 9562 of <i>Lecture Notes in Computer Science</i> , pages 157–175. Springer.	1314
1260			1315
1261			1316
1262			1317
1263	Jimit Majmudar, Christophe Dupuy, Charith Peris, Sami Smaili, Rahul Gupta, and Richard S. Zemel. 2022. Differentially private decoding in large language models. <i>CoRR</i> , abs/2205.13621.	Arvind Narayanan and Vitaly Shmatikov. 2008. Robust de-anonymization of large sparse datasets. In <i>2008 IEEE Symposium on Security and Privacy (S&P 2008), 18-21 May 2008, Oakland, California, USA</i> , pages 111–125. IEEE Computer Society.	1318
1264			1319
1265			1320
1266			1321
1267	Justus Mattern, Benjamin Weggenmann, and Florian Kerschbaum. 2022. The limits of word level differential privacy. <i>arXiv preprint arXiv:2205.02130</i> .	Mijung Park, James R. Foulds, Kamalika Chaudhuri, and Max Welling. 2016. Private topic modeling. <i>CoRR</i> , abs/1609.04120.	1322
1268			1323
1269			1324
1270	H. Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. 2018. Learning differentially private recurrent language models. In <i>6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings</i> . OpenReview.net.	Manas A. Pathak, Shantanu Rane, and Bhiksha Raj. 2010. Multiparty differential privacy via aggregation of locally trained classifiers. In <i>Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010. Proceedings of a meeting held 6-9 December 2010, Vancouver, British Columbia, Canada</i> , pages 1876–1884. Curran Associates, Inc.	1325
1271			1326
1272			1327
1273			1328
1274			1329
1275			1330
1276	Casey Meehan, Khalil Mrini, and Kamalika Chaudhuri. 2022. Sentence-level privacy for document embeddings. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3367–3380, Dublin, Ireland. Association for Computational Linguistics.	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.	1331
1277			1332
1278			1333
1279			1334
1280			1335
1281			1336
1282	Sebastian Meiser and Esfandiar Mohammadi. 2018. Tight on budget?: Tight bounds for r-fold approximate differential privacy. In <i>Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, CCS 2018, Toronto, ON, Canada, October 15-19, 2018</i> , pages 247–264. ACM.	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.	1337
1283			1338
1284			1339
1285			1340
1286			1341
1287			1342
1288	Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. In <i>1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings</i> .	Ildikó Pilán, Pierre Lison, Lilja Øvrelid, Anthi Papadopoulou, David Sánchez, and Montserrat Batet. 2022. The text anonymization benchmark (TAB): A dedicated corpus and evaluation framework for text anonymization. <i>CoRR</i> , abs/2202.00443.	1343
1289			1344
1290			1345
1291			1346
1292			1347
1293			

1348	Richard Plant, Dimitra Gkatzia, and Valerio Giuffrida.	Jingye Tang, Tianqing Zhu, Ping Xiong, Yu Wang, and	1402
1349	2021. CAPE: Context-aware private embeddings	Wei Ren. 2020. Privacy and utility trade-off for tex-	1403
1350	for private language learning . In <i>Proceedings of the</i>	tual analysis via calibrated multivariate perturbations .	1404
1351	<i>2021 Conference on Empirical Methods in Natural</i>	In <i>Network and System Security - 14th International</i>	1405
1352	<i>Language Processing</i> , pages 7970–7978, Online and	<i>Conference, NSS 2020, Melbourne, VIC, Australia,</i>	1406
1353	Punta Cana, Dominican Republic. Association for	<i>November 25-27, 2020, Proceedings</i> , volume 12570	1407
1354	Computational Linguistics.	of <i>Lecture Notes in Computer Science</i> , pages 342–	1408
		353. Springer.	1409
1355	Natalia Ponomareva, Jasmijn Bastings, and Sergei Vas-	Xinyu Tang, Richard Shin, Huseyin A. Inan, Andre	1410
1356	silvitskii. 2022. Training text-to-text transformers	Manoel, Fatemehsadat Miresghallah, Zinan Lin,	1411
1357	with privacy guarantees. In <i>Findings of the Associa-</i>	Sivakanth Gopi, Janardhan Kulkarni, and Robert Sim.	1412
1358	<i>tion for Computational Linguistics: ACL 2022</i> , pages	2023. Privacy-preserving in-context learning with	1413
1359	2182–2193.	differentially private few-shot generation . <i>CoRR</i> ,	1414
		abs/2309.11765.	1415
1360	Chen Qu, Weize Kong, Liu Yang, Mingyang Zhang,	Zhiliang Tian, Yingxiu Zhao, Ziyue Huang, Yu-Xiang	1416
1361	Michael Bendersky, and Marc Najork. 2021. Natu-	Wang, Nevin L. Zhang, and He He. 2022. Seqpate:	1417
1362	ral language understanding with privacy-preserving	Differentially private text generation via knowledge	1418
1363	BERT . In <i>CIKM '21: The 30th ACM International</i>	distillation . In <i>NeurIPS</i> .	1419
1364	<i>Conference on Information and Knowledge Manage-</i>		
1365	<i>ment, Virtual Event, Queensland, Australia, Novem-</i>	Boxin Wang, Yibo Jacky Zhang, Yuan Cao, Bo Li,	1420
1366	<i>ber 1 - 5, 2021</i> , pages 1488–1497. ACM.	H Brendan McMahan, Sewoong Oh, Zheng Xu, and	1421
		Manzil Zaheer. 2023. Can public large language	1422
1367	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	models help private cross-device federated learning?	1423
1368	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	<i>arXiv preprint arXiv:2305.12132</i> .	1424
1369	Wei Li, and Peter J Liu. 2020. Exploring the limits		
1370	of transfer learning with a unified text-to-text trans-	Hua Wang, Sheng Gao, Huanyu Zhang, Milan Shen,	1425
1371	former. <i>The Journal of Machine Learning Research</i> ,	and Weijie J Su. 2022. Analytical composition of dif-	1426
1372	21(1):5485–5551.	ferential privacy via the edgeworth accountant. <i>arXiv</i>	1427
		<i>preprint arXiv:2206.04236</i> .	1428
1373	Swaroop Ramaswamy, Om Thakkar, Rajiv Mathews,	Teng Wang, Xuefeng Zhang, Jingyu Feng, and Xinyu	1429
1374	Galen Andrew, H Brendan McMahan, and Françoise	Yang. 2020a. A comprehensive survey on local dif-	1430
1375	Beaufays. 2020. Training production language mod-	ferential privacy toward data statistics and analysis .	1431
1376	els without memorizing user data. <i>arXiv preprint</i>	<i>Sensors</i> , 20(24).	1432
1377	<i>arXiv:2009.10031</i> .		
		Tianhao Wang, Jeremiah Blocki, Ninghui Li, and	1433
1378	David Sánchez and Montserrat Batet. 2016. C-sanitized:	Somesh Jha. 2017. Locally differentially private pro-	1434
1379	A privacy model for document redaction and saniti-	tocols for frequency estimation . In <i>26th USENIX</i>	1435
1380	zation . <i>J. Assoc. Inf. Sci. Technol.</i> , 67(1):148–163.	<i>Security Symposium, USENIX Security 2017, Van-</i>	1436
		<i>couver, BC, Canada, August 16-18, 2017</i> , pages 729–	1437
1381	Manuel Senge, Timour Igamberdiev, and Ivan Habern-	745. USENIX Association.	1438
1382	al. 2022. One size does not fit all: Investigating		
1383	strategies for differentially-private learning across	Yu-Xiang Wang, Borja Balle, and Shiva Prasad Ka-	1439
1384	nlp tasks. In <i>Proceedings of the 2022 Conference on</i>	saviswanathan. 2020b. Subsampled rényi differential	1440
1385	<i>Empirical Methods in Natural Language Processing</i> ,	privacy and analytical moments accountant . <i>J. Priv.</i>	1441
1386	pages 7340–7353.	<i>Confidentiality</i> , 10(2).	1442
		Benjamin Weggenmann and Florian Kerschbaum. 2021.	1443
1387	Weiyang Shi, Aiqi Cui, Evan Li, Ruoxi Jia, and Zhou	Differential privacy for directional data . In <i>CCS '21:</i>	1444
1388	Yu. 2021. Selective differential privacy for language	<i>2021 ACM SIGSAC Conference on Computer and</i>	1445
1389	modeling . <i>CoRR</i> , abs/2108.12944.	<i>Communications Security, Virtual Event, Republic of</i>	1446
		<i>Korea, November 15 - 19, 2021</i> , pages 1205–1222.	1447
1390	Weiyang Shi, Ryan Shea, Si Chen, Chiyuan Zhang, Ruoxi	ACM.	1448
1391	Jia, and Zhou Yu. 2022. Just fine-tune twice: Selec-	Benjamin Weggenmann, Valentin Rublack, Michael An-	1449
1392	tive differential privacy for large language models.	drejczuk, Justus Mattern, and Florian Kerschbaum.	1450
1393	In <i>Proceedings of the 2022 Conference on Empiri-</i>	2022a. Dp-vae: Human-readable text anonymization	1451
1394	<i>cal Methods in Natural Language Processing</i> , pages	for online reviews with differentially private varia-	1452
1395	6327–6340.	tional autoencoders . In <i>Proceedings of the ACM Web</i>	1453
		<i>Conference 2022</i> , pages 721–731.	1454
1396	Reza Shokri, Marco Stronati, Congzheng Song, and Vi-	Benjamin Weggenmann, Valentin Rublack, Michael An-	1455
1397	taly Shmatikov. 2017. Membership inference attacks	drejczuk, Justus Mattern, and Florian Kerschbaum.	1456
1398	against machine learning models . In <i>2017 IEEE Sym-</i>		
1399	<i>posium on Security and Privacy, SP 2017, San Jose,</i>		
1400	<i>CA, USA, May 22-26, 2017</i> , pages 3–18. IEEE Com-		
1401	puter Society.		

1457	2022b. DP-VAE: human-readable text anonymization for online reviews with differentially private variational autoencoders . In <i>WWW '22: The ACM Web Conference 2022, Virtual Event, Lyon, France, April 25 - 29, 2022</i> , pages 721–731. ACM.	<i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.	1512 1513 1514
1462	Tong Wu, Ashwinee Panda, Jiachen T. Wang, and Prateek Mittal. 2023. Privacy-preserving in-context learning for large language models .	Da Yu, Huishuai Zhang, Wei Chen, Jian Yin, and Tie-Yan Liu. 2021. Large scale private learning via low-rank reparametrization . In <i>Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event</i> , volume 139 of <i>Proceedings of Machine Learning Research</i> , pages 12208–12218. PMLR.	1515 1516 1517 1518 1519 1520 1521
1466	Dominik Wunderlich, Daniel Bernau, Francesco Aldà, Javier Parra-Arnau, and Thorsten Strufe. 2021. On the privacy-utility trade-off in differentially private hierarchical text classification . <i>CoRR</i> , abs/2103.02895.	Lei Yu, Ling Liu, Calton Pu, Mehmet Emre Gursoy, and Stacey Truex. 2019. Differentially private model publishing for deep learning . In <i>2019 IEEE Symposium on Security and Privacy, SP 2019, San Francisco, CA, USA, May 19-23, 2019</i> , pages 332–349. IEEE.	1522 1523 1524 1525 1526
1470	Tianyu Xia, Shuheng Shen, Su Yao, Xinyi Fu, Ke Xu, Xiaolong Xu, and Xing Fu. 2023. Differentially private learning with per-sample adaptive clipping . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 37, pages 10444–10452.	Xiang Yue, Minxin Du, Tianhao Wang, Yaliang Li, Huan Sun, and Sherman S. M. Chow. 2021. Differential privacy for text analytics via natural text sanitization . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 3853–3866, Online. Association for Computational Linguistics.	1527 1528 1529 1530 1531 1532 1533
1474	Xingxing Xiong, Shubo Liu, Dan Li, Zhaohui Cai, and Xiaoguang Niu. 2020. A comprehensive survey on local differential privacy . <i>Secur. Commun. Networks</i> , 2020:8829523:1–8829523:29.	Xiang Yue, Huseyin A Inan, Xuechen Li, Girish Kumar, Julia McAnallen, Huan Sun, David Levitan, and Robert Sim. 2022. Synthetic text generation with differential privacy: A simple and practical recipe . <i>arXiv preprint arXiv:2210.14348</i> .	1534 1535 1536 1537 1538
1478	Nan Xu, Oluwaseyi Feyisetan, Abhinav Aggarwal, Zekun Xu, and Nathanael Teissier. 2021a. Density-aware differentially private textual perturbations using truncated gumbel noise . In <i>Proceedings of the Thirty-Fourth International Florida Artificial Intelligence Research Society Conference, North Miami Beach, Florida, USA, May 17-19, 2021</i> .	Elad Ben Zaken, Yoav Goldberg, and Shauli Ravfogel. 2022. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 1–9.	1539 1540 1541 1542 1543 1544
1485	Zekun Xu, Abhinav Aggarwal, Oluwaseyi Feyisetan, and Nathanael Teissier. 2020. A differentially private text perturbation method using a regularized mahalanobis metric . <i>CoRR</i> , abs/2010.11947.	Fangyuan Zhao, Xuebin Ren, Shusen Yang, Qing Han, Peng Zhao, and Xinyu Yang. 2021. Latent dirichlet allocation model training with differential privacy . <i>IEEE Trans. Inf. Forensics Secur.</i> , 16:1290–1305.	1545 1546 1547 1548
1489	Zekun Xu, Abhinav Aggarwal, Oluwaseyi Feyisetan, and Nathanael Teissier. 2021b. On a utilitarian approach to privacy preserving text generation . <i>arXiv preprint arXiv:2104.11838</i> .	Qinqing Zheng, Jinshuo Dong, Qi Long, and Weijie J. Su. 2020. Sharp composition bounds for gaussian differential privacy via edgeworth expansion . In <i>Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event</i> , volume 119 of <i>Proceedings of Machine Learning Research</i> , pages 11420–11435. PMLR.	1549 1550 1551 1552 1553 1554 1555
1491	Zekun Xu, Abhinav Aggarwal, Oluwaseyi Feyisetan, and Nathanael Teissier. 2021c. On a utilitarian approach to privacy preserving text generation . <i>CoRR</i> , abs/2104.11838.	Xin Zhou, Yi Lu, Ruotian Ma, Tao Gui, Yuran Wang, Yong Ding, Yibo Zhang, Qi Zhang, and Xuanjing Huang. 2023. Textobfuscator: Making pre-trained language model a privacy protector via obfuscating word representations . In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 5459–5473. Association for Computational Linguistics.	1556 1557 1558 1559 1560 1561 1562 1563
1497	Zheng Xu, Yanxiang Zhang, Galen Andrew, Christopher A Choquette-Choo, Peter Kairouz, H Brendan McMahan, Jesse Rosenstock, and Yuanbo Zhang. 2023. Federated learning of gboard language models with differential privacy . <i>arXiv preprint arXiv:2305.18465</i> .	Yuqing Zhu and Yu-Xiang Wang. 2019. Poisson subsampled rényi differential privacy . In <i>Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA</i> , volume 97 of <i>Proceedings of Machine Learning Research</i> , pages 7634–7642. PMLR.	1564 1565 1566 1567 1568 1569
1503	Ying Yin and Ivan Habernal. 2022. Privacy-preserving models for legal natural language processing . In <i>Proceedings of the Natural Legal Language Processing Workshop 2022</i> , pages 172–183.		
1507	Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A. Inan, Gautam Kamath, Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, Sergey Yekhanin, and Huishuai Zhang. 2022. Differentially private fine-tuning of language models . In		

A Differential Privacy Preliminaries

Differential Privacy (DP) is a data post-processing technique, which guarantees data privacy by confusing the attacker. To be more specific, suppose there is one dataset noted as S , and we can get another dataset S' by changing or deleting one data record in this dataset. Denote the output distribution when S is the input as P_1 , and the output distribution when S' is the input as P_2 , if P_1 and P_2 are almost the same then we cannot distinguish these two distributions, i.e., we cannot infer whether the deleted or replaced data sample based on the output we observed. The formal details are given by [Dwork et al. \(2006\)](#). Note that in the definition of DP, adjacency is a key notion. One of the commonly used adjacency definitions is that two datasets S and S' are adjacent (denoted as $S \sim S'$) if S' can be obtained by modifying one record in S .

Definition 2. Given a domain of dataset \mathcal{X} . A randomized algorithm $\mathcal{A} : \mathcal{X} \mapsto \mathcal{R}$ is (ϵ, δ) -differentially private (DP) if for all adjacent datasets S, S' with each sample is in \mathcal{X} and for all $T \subseteq \mathcal{R}$, the following holds

$$\Pr(\mathcal{A}(S) \in T) \leq \exp(\epsilon) \Pr(\mathcal{A}(S') \in T) + \delta.$$

When $\delta = 0$, we call the algorithm \mathcal{A} is ϵ -DP.

Illustration: For example, let \mathcal{X} be a collection of labeled product reviews, each belonging to a single individual, and let \mathcal{R} be parameters of a classifier trained on \mathcal{X} . If the classifier’s training procedure \mathcal{A} satisfies the DP definition above, an attacker’s ability to find out whether a particular individual was present in the training data or not is limited by ϵ and δ .

In the definition of DP, there are two parameters ϵ and δ . Specifically, ϵ measures the closeness between the output distribution when the input is S and the output distribution when the input is S' , smaller ϵ indicates the two distributions are more indistinguishable, i.e., the algorithm \mathcal{A} will be more private. In practice we set $\epsilon = 0.1 - 0.5$ as high privacy regime. Informally, δ could be thought as the probability that ratio between the two distributions is not bounded by e^ϵ . Thus, it is preferable to set δ as small as possible. In practice we always set δ as a value from $\frac{1}{n^{1.1}}$ to $\frac{1}{n^2}$, where n is the number of samples in the dataset S . It is notable that besides ϵ and (ϵ, δ) -DP, there are also other definitions DP such as Rényi DP ([Mironov,](#)

[2017](#)), Concentrated DP ([Bun and Steinke, 2016;](#) [Dwork and Rothblum, 2016](#)), Gaussian DP ([Dong et al., 2022](#)) and Truncated CDP ([Bun et al., 2018](#)). However, all of them can be transformed to the original definition of DP. Thus, in this survey we mainly focus on Definition 2.

There are several important properties of DP, see ([Dwork and Roth, 2014](#)) for details. Here we only introduce those which are commonly used in NLP tasks. The first one is post-processing which means that any post-processing on the output of an (ϵ, δ) -DP algorithm will remain (ϵ, δ) -DP. Equivalently, if an algorithm is DP, then any side information available to the adversary cannot increase the risk of privacy leakage.

Proposition 1. Let $\mathcal{A} : \mathcal{X} \mapsto \mathbb{R}$ be (ϵ, δ) -DP, and let $f : \mathcal{R} \mapsto \mathcal{R}'$ be a (randomized) algorithm. Then $f \circ \mathcal{A} : \mathcal{X} \mapsto \mathbb{R}'$ is (ϵ, δ) -DP.

Example: Continuing with our scenario of training a review classifier under DP, let us imagine we take the model from the previous example, which was trained under (ϵ, δ) -DP, and perform a domain adaptation by fine-tuning on a different dataset, this time without any privacy. The resulting model still remains (ϵ, δ) -DP with respect to the original data, that is privacy cannot be weakened by any post-processing.

The second property is the composition property. Generally speaking, the composition property guarantees that the composition of several DP mechanisms is still DP.

Proposition 2 (Basic Composition Theorem). Let $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_k$ be k be a sequence of randomized algorithms, where $\mathcal{A}_1 : \mathcal{X} \mapsto \mathcal{R}_1$ and $\mathcal{A}_i : \mathcal{R}_1 \times \dots \times \mathcal{R}_{i-1} \times \mathcal{X} \mapsto \mathcal{R}_i$ for $i = 2, \dots, k$. Suppose that for each $i \in [k]$, $\mathcal{A}_i(a_1, \dots, a_{i-1}, \cdot)$ is (ϵ_i, δ_i) -DP. Then the algorithm $\mathcal{A} : \mathcal{X} \mapsto \mathcal{R}_1 \times \dots \times \mathcal{R}_k$ that runs the algorithms \mathcal{A}_i in sequence is (ϵ, δ) -DP with $\epsilon = \sum_{i=1}^k \epsilon_i$ and $\delta = \sum_{i=1}^k \delta_i$.

The basic composition allows us to design complex algorithms by putting together smaller pieces. We can view the overall privacy parameter ϵ as a budget to be divided among these pieces. We will thus often refer to (ϵ, δ) as the “privacy budget”: each algorithm we run leaks some information, and consumes some of our budget. Differential privacy allows us to view information leakage as a resource to be managed. For example, if we fix the privacy budget (ϵ, δ) , then making each \mathcal{A}_i be $(\frac{\epsilon}{k}, \frac{\delta}{k})$ -DP is sufficient to ensure the composition is (ϵ, δ) -DP.

Method Type	Publications	Scenarios	Definition	Model Architecture	DP Level	Downstream Tasks	
Gradient Perturbation Based Methods	Hoory et al. (2021) Anil et al. (2021) Yin and Habernal (2022) Senge et al. (2022) Ponomareva et al. (2022)	Pre-trained	DP	BERT	Sample-level	Entity-extraction	
	BERT			Sample-level	—		
	BERT			Sample-level	Classification, QA		
	BERT, XtremeDistil			Sample-level	Classification, NER, POS, QA		
	T5			Sample-level	NLU		
	Yu et al. (2022) Yu et al. (2021) Dupuy et al. (2021) Li et al. (2021) Lee and Søgaard (2023) Xia et al. (2023) Behnia et al. (2022) Bu et al. (2023) Gupta et al. (2023) Du and Mi (2021) Bu et al. (2022) Yue et al. (2022) Mireshghallah et al. (2022) Carranza et al. (2023) Igamberdiev and Habernal (2021) Aziz et al. (2022) Wunderlich et al. (2021) Li et al. (2022) Amid et al. (2022) Shi et al. (2021) Shi et al. (2022)	Fine-tuning	DP	RoBERT, GPT-2	Sample-level	NLG, NLU	
	BERT			Sample-level	Classification, NLU		
	BERT, BiLSTM			Sample-level	Classification, NER		
	GPT-2, (Ro)BERT			Sample-level	Classification, NLG		
	GPT-2, DialoGPT			Sample-level	Meeting Summarization		
	GPT-2, (Ro)BERT			Sample-level	Classification		
	(Ro)BERT			Sample-level	NLU		
	GPT-2, (Ro)BERT			Sample-level	Classification		
	(Ro)BERT			Sample-level	GLU		
	GPT-2, (Ro)BERT			Sample-level	Classification, NLG		
	(Ro)BERT			Sample-level	Classification, NLG		
	GPT-2			Sample-level	Synthetic Text Generation		
	GPT-2			Sample-level	Synthetic Text Generation		
	T5			Sample-level	Query Generation		
	GPT-2			Sample-level	Classification		
	GPT-2	Sample-level	Synthetic Text Generation				
	BERT, CNN	Sample-level	Classification				
	LSTM	Sample-level	Classification				
	LSTM	Sample-level	Classification				
	RNN	Sample-level	NLG, Dialog System				
	GPT-2, (Ro)BERT	Sample-level	NLG, NLU				
	Federated Learning	McMahan et al. (2018) Ramawamy et al. (2020) Kairouz et al. (2021) Choquette-Choo et al. (2022) Koloskova et al. (2023) Denisov et al. (2022) Agarwal et al. (2021) Wang et al. (2023) Xu et al. (2023)	Federated Learning	LDP	LSTM, RNN	User-level	Prediction, Classification
LSTM		User-level			Prediction, Classification		
LSTM		User-level, Sample-level			Prediction, Classification		
LSTM		User-level, Sample-level			Prediction		
LSTM		User-level, Sample-level			Prediction		
LSTM		User-level, Sample-level			Prediction		
LSTM		User-level, Sample-level			Prediction		
LaMDA		User-level			Prediction		
Gboard		User-level			Prediction		
Embedding Vector Perturbation Based Methods		Lyu et al. (2020b) Lyu et al. (2020a) Plant et al. (2021) Krishna et al. (2021) Habernal (2021) Alnasser et al. (2021) Igamberdiev et al. (2022) Maheshwari et al. (2022) Bollegala et al. (2023) Chen et al. (2023) Du et al. (2023b)	Private Embedding	LDP	BERT	Word-level	Classification
		BERT			Word-level	Classification	
		BERT			Word-level	Classification	
		Auto-Encoder			Word-level	Classification	
		Auto-Encoder			Word-level	Classification	
		Auto-Encoder			Word-level	Classification	
	Auto-Encoder	Word-level	Classification				
	GloVe	Word-level	Classification				
	GloVe, BERT	Token-level	Classification				
	BERT	Sentence-level	Classification, QA				
	Meehan et al. (2022)	Private Embedding	DP	SBERT	Sentence-level	Classification	
	Feyisetan et al. (2020) Xu et al. (2020) Xu et al. (2021c) Xu et al. (2021a) Carvalho et al. (2021b) Feyisetan and Kasiviswanathan (2021) Feyisetan et al. (2019) Carvalho et al. (2021a) Tang et al. (2020) Imola et al. (2022) Arnold et al. (2023a) Arnold et al. (2023b) Qu et al. (2021) Du et al. (2023a)	Private Embedding	LMDP	GloVe, BiLSTM	Word-level	Classification, QA	
	GloVe			Word-level	Classification		
	GloVe, FastText			Word-level	Classification		
	GloVe, CNN			Word-level	Classification		
GloVe	Word-level			Classification			
GloVe, FastText	Word-level			Classification			
GloVe	Word-level			Classification, Prediction			
GloVe, FastText	Word-level			Classification			
GloVe	Word-level			Classification			
GloVe, FastText	Word-level			Classification			
GloVe	Word-level			Classification			
GloVe	Word-level			Classification			
BERT, BiLSTM	Token-level	Classification, NLU					
BERT	Sentence-level	Classification, QA					
Yue et al. (2021)	Private Embedding	UMLDP	BERT, GloVe	Word-level	Classification, QA		

Table 1: An overview of studies for DP-NLP.

Example: In most of the NLP tasks we need to train a model by using variants of optimization methods, such as SGD or Adam. In general, these optimizers include several iterations to update the model, which could be thought as a composition algorithm and each iteration could be thought as an algorithm. Thus, it is sufficient to design DP algorithm for each iteration and we can use the composition theorem to calculate the budget of the whole process.

Beside the basic composition property, there are also several advanced composition theorem for (ϵ, δ) -DP, which could provide tighter privacy guarantees than the basic one. For example, consider each $\mathcal{A}_i, i \in [k]$ is (ϵ, δ) -DP. Then the basic composition theorem implies their composition is $(k\epsilon, k\delta)$ -DP. However, this is not tight as we can use the advanced composition theorem to show their composition could be improved to $(O(\sqrt{k}\epsilon), O(k\delta))$ -DP (Dwork et al., 2010). We refer to reference (Kairouz et al., 2015; Murtagh and Vadhan, 2016; Meiser and Mohammadi, 2018) for details.

The third property is the privacy amplification via subsampling. Intuitively, every differentially private algorithm has a much lower privacy parameter ϵ when it is run on a secret sample than when it is run on a sample whose identities are known to the attacker. And there a secret sample can be obtained by subsampling as it introduces additional randomness.

Proposition 3. Let A be an (ϵ, δ) -DP algorithm. Now we construct the algorithm B as follows: On input $D = \{x_1, \dots, x_n\}$, first we construct a new sub-sampled dataset D_S where each $x_i \in D_s$ with probability q . Then we run algorithm A on the dataset D_S . Then $B(D) = A(D_S)$ is $(\tilde{\epsilon}, \tilde{\delta})$ -DP, where $\tilde{\epsilon} = \ln(1 + (e^\epsilon - 1)q)$ and $\tilde{\delta} = q\delta$.

Example: The subsampling property can be used to private version of the stochastic optimization method. As in these methods, a common strategy is to use subsampled gradient to estimate the whole gradient.

It is notable that, besides subsampling, some other procedures could also amplify privacy such as random check-in (Balle et al., 2020), mixing (Balle et al., 2019) and decentralization (Cyffers and Bellet, 2022). And for different subsampling method, the privacy amplification guarantee is also different (Imola and Chaudhuri, 2021; Zhu and Wang, 2019; Balle et al., 2018).

In the following, we will introduce some mechanisms commonly used in NLP tasks to achieve DP.

We first give the definition of a (numeric) query. The query is simply something we want to learn from the dataset. Formally, a query could be any function f applied to a dataset S and outputting a real valued vector, formally $f : \mathcal{X} \mapsto \mathbb{R}^d$. For example, numeric queries might return the sum of the gradient of the loss on all samples, number of females in the database, or a textual summary of medical records of all persons in the database represented as a dense vector. Given a dataset S , a common paradigm for approximating $f(S)$ differentially privately is via adding some randomized noise. And Laplacian noise and Gaussian noise are the most commonly used ones, which correspond to the Laplacian and Gaussian mechanism respectively.

Definition 3 (Laplacian Mechanism). Given a query $f : \mathcal{X} \mapsto \mathbb{R}^d$, the Laplacian Mechanism is defined as: $\mathcal{M}_L(S, f, \epsilon) = q(S) + (Y_1, Y_2, \dots, Y_d)$, where Y_i is i.i.d. drawn from a Laplacian Distribution $\text{Lap}(\frac{\Delta_1(f)}{\epsilon})$, where $\Delta_1(f)$ is the ℓ_1 -sensitivity of the function f , i.e., $\Delta_1(f) = \sup_{S' \sim S'} \|f(S) - f(S')\|_1$. For a parameter λ , the Laplacian distribution has the density function $\text{Lap}(\lambda)(x) = \frac{1}{2\lambda} \exp(-\frac{x}{\lambda})$. Laplacian Mechanism preserves ϵ -DP.

Definition 4 (Gaussian Mechanism). Given a query $f : \mathcal{X} \mapsto \mathbb{R}^d$, the Gaussian mechanism is defined as $\mathcal{M}_F(S, f, \epsilon, \delta) = q(S) + \xi$ where $\xi \sim \mathcal{N}(0, \frac{2\Delta_2^2(f) \log(1.25/\delta)}{\epsilon^2} \mathbb{I}_d)$, where $\Delta_2(f)$ is the ℓ_2 -sensitivity of the function f , i.e., $\Delta_2(f) = \sup_{S \sim S'} \|f(S) - f(S')\|_2$. Gaussian mechanism preserves (ϵ, δ) -DP when $0 < \epsilon \leq 1$.

From the previous two mechanisms we can see that to privately release $f(S)$ it is sufficient to calculate the ℓ_1 -norm or ℓ_2 -norm sensitivity first and add random noise. Moreover, as $\Delta_2(f) \leq \Delta_1(f)$, Gaussian mechanism will has lower error than the Laplacian mechanism, while we relax the definition from ϵ -DP to (ϵ, δ) -DP.

Instead of answering $f(S)$ privately, we also always meet the selection problem, i.e., we want to output the best candidate among several candidates based on some score of the dataset. Exponential mechanism is the one that can output a nearly best candidate privately.

Definition 5 (Exponential Mechanism). The Ex-

ponential Mechanism allows differentially private computation over arbitrary domains and range \mathcal{R} , parameterized by a score function $u(S, r)$ which maps a pair of input data set S and candidate result $r \in \mathcal{R}$ to a real valued score. With the score function u and privacy budget ϵ , the mechanism yields an output with exponential bias in favor of high scoring outputs. Let $\mathcal{M}(S, u, \mathcal{R})$ denote the exponential mechanism, and Δ be the sensitivity of u in the range \mathcal{R} , *i.e.*, $\Delta = \max_{r \in \mathcal{R}} \max_{D \sim D'} |u(D, r) - u(D', r)|$. Then if $\mathcal{M}(S, u, \mathcal{R})$ selects and outputs an element $r \in \mathcal{R}$ with probability proportional to $\exp(\frac{\epsilon u(S, r)}{2\Delta u})$, it preserves ϵ -DP.

In the original definition of DP, we assume that data are managed by a trusted centralized entity which is responsible for collecting them and for deciding which differentially private data analysis to perform and to release. A classical use case for this model is the one of census data. Compared with the above model (which is called central model), there is another model namely local DP model, where each individual manages his/her proper data and discloses them to a server through some differentially private mechanisms. The server collects the (now private) data of each individual and combines them into a resulting data analysis. A classical use case for this model is the one aiming at collecting statistics from user devices like in the case of Google’s Chrome browser. Formally it is defined as follows.

Definition 6. For a data domain \mathcal{X} , a randomized algorithm $\mathcal{A} : \mathcal{X} \mapsto \mathcal{R}$ is called (ϵ, δ) -local DP (LDP) if for any $s, s' \in \mathcal{X}$ and $T \subseteq \mathcal{R}$ we have

$$\Pr[\mathcal{A}(s) \in T] \leq e^\epsilon \Pr[\mathcal{A}(s') \in T] + \delta.$$

Compared with Definition 2 we can see that here the main difference is the inequality hold for all elements $s, s' \in \mathcal{X}$ instead of all adjacent pairs of dataset. In this case, each individual could ensure that their own disclosures are DP via the randomizer \mathcal{A} . In some sense, the trust barrier is moved closer to the user. While this has a benefit of providing a stronger privacy guarantee, it also comes at a cost in terms of accuracy.

It is notable that besides the central DP and local DP model, there are also other intermediate models such as shuffle model (Cheu et al., 2019) and multi-party setting (Pathak et al., 2010). However, as they are seldom studied in NLP, we will not cover these protocols in this survey.

B An Introduction to DP-SGD

Given a training data with n samples $D = \{x_i\}_{i=1}^n$, a loss function (such as cross-entropy loss) is defined to train the model, which takes the parameter $\theta \in \mathbb{R}^d$ of neural network and samples and outputs a real value:

$$L(\theta, D) = \sum_{i=1}^n \ell(\theta, x_i). \quad (1)$$

The goal is to find the weights of the network that minimizes $L(\theta, D)$, *i.e.*, $\theta^* = \arg \min_{\theta} L(\theta, D)$. With additional constraint on DP, now we aim to design an $(\epsilon, \delta)/\epsilon$ -DP algorithm \mathcal{A} to make the private estimated parameter θ_{priv} close to θ^* .

Example: In Language Modeling (LM), we have a corpus $D = \{x_1, \dots, x_n\}$ where each text sequence x_i consists of multiple tokens $x_i = (x_{i1}, \dots, x_{im_i})$ with x_{ij} as the j -th token of x_i . The goal of LM is to train a neural network (e.g., RNN) parameterized by θ to learn the probability of the sequence $p_{\theta}(x)$, which can be represented as the following objective function

$$- \sum_{i=1}^n \sum_{j=1}^{m_i} \log p_{\theta}(x_{ij} | x_{i1}, \dots, x_{i(j-1)}). \quad (2)$$

We first review the DP-SGD method (Abadi et al., 2016). In the non-private case, to minimize the objective function (1), the most fundamental method is SGD, *i.e.*, in the t -th iteration we update the model as follows:

$$\theta^{t+1} = \theta^t - \eta \frac{1}{|B|} \sum_{x \in B} \nabla \ell(\theta^t, x), \quad (3)$$

where B is a subsampled batch of random examples, η is the learning rate and θ^t is the current parameter. DP-SGD modifies the SGD-based methods by adding Gaussian noise to perturb the (stochastic) gradient in each iteration of the training, *i.e.* during the t -th iteration DP-SGD will compute a noisy gradient as follows:

$$g^t = \frac{1}{|B|} \left(\sum_{x_i \in B} \hat{g}_i^t + \mathcal{N}(0, \sigma^2 C^2 I_d) \right), \quad (4)$$

σ is noise multiplier, \hat{g}_i^t is some vector computed from $\nabla \ell(\theta^t, x_i)$ and g^t is the (noisy) gradient used to update the model. The main reason here we use \hat{g}_i^t instead of the original gradient vector is that we wish to make the term $\sum \hat{g}_i^t$ has bounded

1854 ℓ_2 -sensitivity so that we can use the Gaussian
1855 mechanism to ensure DP. The most commonly
1856 used approach to get a \hat{g}_i^t is clipping the gradient:
1857 $\hat{g}_i^t = \nabla \ell(\theta^t, x_i) \min\{1, \frac{C}{\|\nabla \ell(\theta^t, x_i)\|_2}\}$ *i.e.*, each gra-
1858 dient vector is clipped by a hyper-parameter $C > 0$.
1859 Since the ℓ_2 -sensitivity of $\sum \hat{g}_i^k$ is bounded by C ,
1860 after the clipping, we can add Gaussian noise to
1861 ensure DP. As there are several iterations and in
1862 each iteration, we use some subsampling strategy,
1863 we can use the composition theorem and privacy
1864 amplification to compute the total privacy cost of
1865 DP-SGD. Equivalently, given a fixed privacy bud-
1866 get (ϵ, δ) , number of iterations and subsampling
1867 strategy, one can get the minimal noise multiplier σ
1868 to ensure DP, see (Asoodeh et al., 2021; Gopi et al.,
1869 2021; Mironov et al., 2019; Wang et al., 2020b;
1870 Zheng et al., 2020; Zhu and Wang, 2019) for de-
1871 tails.