Private Language Models via Truncated Laplacian Mechanism

Anonymous ACL submission

Abstract

Recently it has been shown that deep learning 002 models for NLP tasks are prone to attacks that can even reconstruct the verbatim training texts. To prevent privacy leakage, researchers have investigated word-level perturbations, relying 006 on the formal guarantees of differential privacy (DP) in the embedding space. However, many existing approaches either achieve unsatisfactory performance in the high privacy regime when using the Laplacian or Gaussian mechanism, or resort to weaker relaxations of DP that are inferior to the canonical DP in terms of privacy strength. This raises the question of whether a new method for private word embedding can be designed to overcome these limitations.

007

011

017

026

027

033

037

041

In this paper, we propose a novel private embedding method called the high dimensional truncated Laplacian mechanism. Specifically, we introduce a non-trivial extension of the truncated Laplacian mechanism, which was previously only investigated in one-dimensional space cases. Theoretically, we show that our method has a lower variance compared to the previous private word embedding methods. To further validate its effectiveness, we conduct comprehensive experiments on private embedding and downstream tasks using three datasets. Remarkably, even in the high privacy regime, our approach only incurs a slight decrease in utility compared to the non-private scenario.

1 Introduction

The recent developments of deep learning have led to significant success in various tasks in Natural Language Processing (NLP), from next word prediction in mobile keyboards (Ramaswamy et al., 2019), to critical tasks like predicting patient health conditions from clinical records (Yao et al., 2019). However, such applications may always involve user-generated textual data as the training dataset, which contains sensitive information. To address

privacy concerns, text anonymization (Anandan et al., 2012; Pilán et al., 2022) has been commonly used, which involves identifying sensitive attributes and replacing them with alternative values. Nevertheless, such heuristic approaches become ineffective as deep neural networks often tend to memorize training data, making them susceptible to information leakage about the training data (Shokri et al., 2017; Carlini et al., 2021, 2019). One way that takes into account the limitations of existing approaches is designing Differentially Private (DP) algorithms. DP (Dwork et al., 2006a) is resilient to arbitrary side information that might be available to attackers and has become a de facto method for private data analysis.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

Recently, there has been significant research focusing on differentially private (DP) versions of word embedding from various perspectives (Yue et al., 2021; Feyisetan et al., 2019; Krishna et al., 2021; Feyisetan et al., 2020; Xu et al., 2021a,b; Carvalho et al., 2021b,a; Habernal, 2021, 2022). However, there are still some shortcomings in these approaches. On the one hand, several works consider adding Laplacian or Gaussian noise to the embedding space to ensure DP (Habernal, 2021; Krishna et al., 2021; Habernal, 2022). However, these mechanisms suffer from high noise levels, resulting in low utility, especially in the high privacy regime when the privacy parameter (ϵ) is small. Moreover, these mechanisms can even alter the semantics of sentences (see Fig.1). On the other hand, there is a growing body of work that focuses on a relaxation of the canonical definition of DP, known as metric DP, which can achieve better performance. However, as a relaxed notion of DP, Metric DP cannot provide the same level of strong privacy guarantees as the canonical DP (Mattern et al., 2022). This raises the question of whether we can develop improved private word embedding mechanisms that go beyond the limitations of Laplacian or Gaussian mechanisms within the framework of canonical DP.

Compariso	Comparison of Private Embedding											
Original :	Oh and we came on a Saturday night around	11:30	for	context.	(→Privacy Leakage)							
Trlaplace:	Oh and we came on a Saturday night around	9:30pm	for	<unk></unk>	• (\rightarrow Private and Fluent)							
Laplace:	Oh and we came on a Saturday night around	around	for	<unk></unk>	(→Semantic Problem)							
Gaussian:	Oh and we came on a Saturday night around	11:30 f	for	<unk></unk>	(→Privacy Leakage)							

Figure 1: An example of (private) text re-write for different mechanisms with $\epsilon = 0.1$.

In this paper, we provide an affirmative answer to the previous question by proposing a novel private mechanism for word embedding. Our approach is inspired by the superior performance of the truncated Laplacian mechanism in onedimensional space (Geng et al., 2020). However, it remains unclear whether this superiority can extend to high dimensional cases, as directly ex-090 tending the one-dimensional truncated Laplacian 091 mechanism is challenging. To bridge this gap, we develop a high dimensional truncated Laplacian mechanism(TrLaplace), which is a non-trivial extension of the one-dimensional case. Theoretically, we show that compared with Laplacian and Gaussian mechanisms for private word embedding, 097 TrLaplace-based private embedding has a lower variance. Moreover, we also conduct intensive experiments on both private embedding and down-100 stream tasks to show our approach significantly 101 102 outperforms the previous methods in the high privacy regime, and it will not drop much accuracy 103 and utility compared with the non-private case. 104

> Due to space limitations, more details and experiments are included in Appendix.

2 Background

105

106

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

Differential Privacy is a data post-processing technique designed to ensure data privacy by adding confusion to potential attackers. Specifically, suppose there is one dataset noted as \mathcal{D} , and we change or delete one data record in this dataset which we call \mathcal{D}' . If the output distributions of \mathcal{D} and \mathcal{D}' are close enough, then we cannot distinguish these two distributions, i.e., we cannot infer whether the deleted or replaced data sample is really in this dataset. The formal details are given by (Dwork et al., 2006b).

In this work, we adopt a similar setting to previous research on private word embedding (Feyisetan et al., 2020; Xu et al., 2021a; Krishna et al., 2021). We consider a scenario where a user inputs a word w from a discrete fixed vocabulary W. Our goal is to preserve the user's privacy with respect to her/his word. To achieve this goal, we aim to design an algorithm that accepts w as input and whose distribution of output is close to the case where $w' \in W$ is the input, with $w' \neq w$ is any other word. From the attacker's perspective, based on the output, he cannot distinguish whether the user's input word is w or w' as their output distributions are almost the same. Formally, we have the following definition. 124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

Definition 1 Given a discrete vocabulary \mathcal{W} , a randomized algorithm $\mathcal{A} : \mathcal{W} \mapsto \mathcal{R}$ is word-level (ϵ, δ) -differentially private (DP) if for all pair of words $w, w' \in \mathcal{W}$ and for all $T \subseteq \mathcal{R}$ we have $\mathbb{P}(\mathcal{A}(w) \in T) \leq e^{\epsilon} \mathbb{P}(\mathcal{A}(w') \in T) + \delta$. When $\delta = 0$, we call the algorithm \mathcal{A} is ϵ -DP.

In this paper, we assume the user holds a sentence $s = w_1 w_2 \cdots w_n$ with *n* words. And we aim to design an (ϵ, δ) -DP algorithm, which is private w.r.t. each word w_i .

3 Private Embedding via Truncated Laplacian Mechanism

In this section, we will provide details of our method. Generally speaking, for each token w_i , to achieve DP, our approach consists of three steps. First, each token w_i is mapped to an *d*-dimensional pre-trained word embedding $\phi(w_i)$. And we perform a clipping step to get a clipped embedding:

CLIPEmb
$$(w_i) = \phi(w_i) \min\{1, \frac{C}{\|\phi(w_i)\|_2}\},$$
 (1)

where the threshold C > 0 is a hyper-parameter. In the second step, we add some random noise to the clipped embedding vector to make it satisfies DP. Finally, we will perform the projection step by finding the nearest word \hat{w}_i to the perturbed and clipped embedding vector within the embedding space:

$$\hat{w}_i = \arg\min_{w \in \mathcal{W}} \|\phi(w) - \text{CLIPEmb}(w_i) - \eta\|_2,$$
(2)

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

178

179

181

182

183

184

185

186

188

190

191

193

where η is the randomized noise we add in the second step. See Algorithm 1 for details. It is notable

|--|

Input: String $s = w_1 w_2 \dots w_n$, clipping threshold *C*, privacy parameter $\epsilon > 0$.

Output: String $\hat{s} = \hat{w}_1 \hat{w}_2 \dots \hat{w}_n$.

1: for all $i \in \{1, ..., n\}$ do

- 2: Sample η from the truncated Laplacian distribution in Theorem 3.
- 3: Obtain the perturbed clipped embedding $\mathbf{r}_i = \text{CLIPEmb}(w_i) + \eta.$

4: Let
$$\hat{w}_i = \operatorname{Proj}(\mathbf{r_i})$$
 as in (2).

5: end for

6: **return** $\hat{s} = \hat{w}_1 \hat{w}_2 \dots \hat{w}_n$.

that the goal of clipping is to make the ℓ_2 -norm of embedding vector be bounded so that we can adding noise to ensure DP, such as the Laplacian mechanism or Gaussian mechanism (Dwork and Roth, 2014).

Theorem 1 (Laplacian Mechanism) Suppose

CLIPEmb(\mathbf{w}) $\in \mathbb{R}^d$ denote the clipped embedding vector with threshold C. Then the mechanism $\mathcal{A}_{lap}(w) = \text{CLIPEmb}(w) + \eta_1$ is ϵ -DP, where $\eta_1 = (\eta_{1,1}, \cdots, \eta_{1,d})$ and $\eta_{i,j}$ is drawn from a Laplacian Distribution $Lap(\frac{\Delta_1(f)}{\epsilon})$ with $\Delta_1 = 2\sqrt{dC}$. For a parameter λ , the Laplacian distribution has the density function $Lap(\lambda)(x) = \frac{1}{2\lambda} \exp(-\frac{x}{\lambda})$.

Theorem 2 (Gaussian Mechanism) Suppose

CLIPEmb(**w**) $\in \mathbb{R}^d$ denote the clipped embedding vector with threshold C. Then the mechanism $\mathcal{A}_{lap}(w) = \text{CLIPEmb}(w) + \eta_2 \text{ is } (\epsilon, \delta) \text{-}DP$ when $\epsilon \leq 1$, where $\eta_2 \sim \mathcal{N}(0, \frac{8C^2 \ln(1.25/\delta)}{\epsilon^2} I_d)$ is drawn from a Gaussian distribution.

In the following we propose an improved mechanism namely high dimensional truncated Laplacian mechanism. Before that we first recall the probability density function of the one-dimensional truncated Laplacian distribution, which could be written as the following with some appropriate constants α , A and B:

$$f_{TLap}(x) = \begin{cases} \frac{1}{B}e^{-\alpha|x|}, & \text{for } x \in [-A, A]\\ 0, & \text{otherwise.} \end{cases}$$
(3)

In our mechanism, we add high dimensional truncated Laplacian noise to the clipped embedding vector. Here each coordinate of the noise is i.i.d. sampled from a truncated Laplacian distribution with some specific α , A and B.

It is notable that although using truncated Laplacian noise to ensure DP has been studied in (Geng et al., 2020; Sommer et al., 2021), all of them only consider the case where d = 1 and their methods cannot extend to the case where d > 1. For example, (Geng et al., 2020) only shows adding noise with density function (3) for $A = \frac{\Delta_1}{\epsilon} \log(1 + \frac{e^{\epsilon}}{2\delta})$ and $\alpha = \frac{\epsilon}{\Delta_1}$ can ensure (ϵ, δ) -DP. Compared with Theorem 3 we can see the constant A is more complicated and the proof is also different. Thus, our mechanism cannot be considered as a trivial extension of the one-dimensional case. Secondly, while the Laplacian mechanism can guarantee ϵ -DP, the truncated one can only ensure (ϵ, δ) -DP. However, as we will see below, our mechanism is superior to Laplacian mechanism for utility. It is also notable that we need to assume $\epsilon < 2\delta^{\frac{1}{d}}\sqrt{d}$, this is reasonable since we always wish ϵ to be as small as possible, as large ϵ indicates the algorithm is no longer private. If we want large $\epsilon > 2\delta^{\frac{1}{d}}\sqrt{d}$, we can use the trick of adding dummy dimension to the vector to increase its dimensionality manually and then projecting back to the original space after adding noise. In the following we will show our mechanism has lower variance than the Laplacian and Gaussian mechanisms.



Figure 2: Privacy Test. Curves of the value N_w with privacy budget ϵ for Yelp dataset.

4 Theoretical Sensitivity Analysis 2

In the last section, we introduce our truncated laplacian mechanism, we will analyze its sensitivity and proof our claim in this section.

Theorem 3 Suppose $\text{CLIPEmb}(w) \in \mathbb{R}^d$ is the clipped embedding vector with threshold C. Define 227

222

194

196

197

198

199

201

202

203

204

206

207

210

211

212

213

214

215

216

217

218

219

221

223

230

231

241

242

243

244

245

246

247

248

249

252

254

258

261

263

 $\Delta_{\infty} = 2C$ and $\Delta_1 = 2\sqrt{d}C$. For $\epsilon \leq 2\delta^{\frac{1}{d}}\sqrt{d}$, if

$$\alpha = \frac{\epsilon}{\Delta_1}, A = -\frac{\Delta_1}{\epsilon} \log(1 - \frac{\epsilon}{2\delta^{\frac{1}{d}}\sqrt{d}})$$
$$B = \frac{2(1 - e^{-\alpha A})}{\alpha} = \frac{\Delta_{\infty}}{\delta^{\frac{1}{d}}},$$

then the mechanism $\mathcal{A}(w) = \text{CLIPEmb}(w) + \eta$ is (ϵ, δ) -DP, where $\eta = (\eta_1, \dots, \eta_1)$ and each η_i has the density function as in (3) with the above parameters.

Proof 1 (Proof of Theorem 3) Consider a pair of tokens w, w'. Let perturbed encoderl $r_1 = \text{CLIPEmb}(w) + \eta_1$, also let $r_2 =$ $\text{CLIPEmb}(w') + \eta_2 = \text{CLIPEmb}(w) + \Delta_s + \eta_2$, where $\|\Delta_s\|_1 \leq \Delta_1$ and $\|\Delta_s\|_{\infty} \leq \Delta_{\infty}$ which are due to the clipping operation.

Let us denote the set of possible values of r_k by S_k for k = 1, 2.

Define $\mathcal{U} = [-C - A, C + A]^d$. Note that for any subset $\mathcal{V} \subseteq \mathcal{U} - (\mathcal{S}_1 \cup \mathcal{S}_2), \mathbb{P}(r_1 \in \mathcal{V})) =$ $\mathbb{P}(r_2 \in \mathcal{V}) = 0$, hence (ϵ, δ) -DP is satisfied for this part. We need to ensure (ϵ, δ) -DP is satisfied for all elements in $\mathcal{S}_1 \cup \mathcal{S}_2$ too.

First, consider an element $s \in S_1 \cap S_2$ *. Then:*

$$f(r_1 = s) = f(\eta_1 = s - \text{CLIPEmb}(\mathbf{s}))$$

Similarly:

$$f(r_2 = s) = f(\eta_2 = s - \text{CLIPEmb}(\mathbf{s}) - \boldsymbol{\Delta}_{\boldsymbol{s}})$$

Using the above equations:

$$\exp(-\alpha\Delta_1) \le \exp(-\alpha \|\Delta_s\|_1)$$
$$\le \frac{\mathbb{P}(r_1 = s)}{\mathbb{P}(r_2 = s)} \le \exp(\alpha \|\Delta_s\|_1) \le \exp(\alpha\Delta_1)$$

From the above equation, setting setting $\alpha = \epsilon/\Delta_1$ ensures pure ϵ -DP for all $s \in S_1 \cap S_2$. With this, it follows that for any $\mathcal{V} \subseteq S_1 \cap S_2$:

$$e^{-\epsilon}\mathbb{P}\left(r_{2}\in\mathcal{V}\right)\leq\mathbb{P}\left(r_{1}\in\mathcal{V}\right)\leq e^{\epsilon}\mathbb{P}\left(r_{2}\in\mathcal{V}\right).$$

by setting $\alpha = \epsilon / \Delta_1$.

Now consider an element $s \in S_2 - S_1$. Clearly, $f(r_1 = s) = 0$. Also:

$$\max_{s \in \mathcal{S}_2 - \mathcal{S}_1} \mathbb{P}\left(r_2 = s\right) \le \frac{1}{B}$$

But notice that volume $(S_2 - S_1) \leq \Delta_{\infty}^d$. This follows from the fact that for every coordinate, there

are at most Δ_{∞} levels that can be attained by r_2 264 but not by r_1 . Thus, for any $\mathcal{T} \subseteq S_2 - S_1$, we have 265

$$\mathbb{P}(r_1 \in \mathcal{T}) = 0 \text{ and } \mathbb{P}(r_2 \in \mathcal{T}) \le \left(\frac{\Delta_{\infty}}{B}\right)^d$$
 266

267

277

282

283

284

Similarly, for any $\mathcal{T} \subseteq S_1 - S_2$, we have

$$\mathbb{P}(r_2 \in \mathcal{T}) = 0 \text{ and } \mathbb{P}(r_1 \in \mathcal{T}) \le \left(\frac{\Delta_{\infty}}{B}\right)^d.$$
 263

Now, let us now consider some general $\mathcal{T} \subseteq S_1 \cup$ S_2 . Let $\mathcal{T}_0 = \mathcal{T} \cap (S_1 \cup S_2), \mathcal{T}_1 = \mathcal{T} \cap (S_1 - S_2)$ and $\mathcal{T}_2 = \mathcal{T} \cap (S_2 - S_1)$. It is easy to see that $\mathcal{T} = \mathcal{T}_0 \cup \mathcal{T}_1 \cup \mathcal{T}_2$ and that $\mathcal{T}_0, \mathcal{T}_1$ and \mathcal{T}_2 are pairwise-disjoint. Then: 269 270 270 270 270 271 272 273

$$\mathbb{P}(r_{1} \in \mathcal{T}) = \mathbb{P}(r_{1} \in \mathcal{T}_{0}) + \mathbb{P}(r_{1} \in \mathcal{T}_{1}) + \mathbb{P}(r_{1} \in \mathcal{T}_{2})$$

$$\leq e^{\epsilon} \mathbb{P}(r_{2} \in \mathcal{T}_{0}) + \left(\frac{\Delta_{\infty}}{B}\right)^{d} + 0$$

$$\leq e^{\epsilon} \mathbb{P}(r_{2} \in \mathcal{T}) + \left(\frac{\Delta_{\infty}}{B}\right)^{d}.$$
(4)

Thus, we can set $\delta = (\frac{\Delta_{\infty}}{B})^d$. Obviously, this result 275 is only useful if $B > \Delta_{\infty}$.

For each coordinate

$$\int_{x \in \mathbb{R}} f_{\text{TLap}}(x) dx = \int_0^A 2\frac{1}{B} e^{-\alpha |x|} dx$$
$$= \frac{2}{B\alpha} \left(1 - e^{-\alpha A}\right) = 1$$

We can solve $B = \frac{2(1-e^{-\alpha A})}{\alpha}$. Thus, take 279 $B = \frac{\Delta_{\infty}}{\delta^{\frac{1}{d}}}$, we can see $A = -\frac{1}{\alpha}\log(1-\frac{\alpha\Delta_{\infty}}{2\delta^{\frac{1}{\delta}}}) =$ 280

$$\frac{\Delta_1}{\epsilon} \log(1 - \frac{\epsilon}{2\sqrt{d\delta^{\frac{1}{\delta}}}}).$$
28

Theorem 4 The variance of mechanism \mathcal{A} in Theorem 3 is lower than the variance of Laplacian mechanism and Gaussian mechanism when $\delta \leq \frac{1}{c^d}$.

Proof 2 (Proof of Theorem 4) We first show the variance of our mechanism \mathcal{A} is bounded by $2\frac{d\Delta_1^2}{\epsilon^2}$. We can easily see that the variance is $\mathbb{E} || \mathcal{A}(w) - 287$ $w ||_2^2 = dV$ with $V = \int_{x \in \mathbb{R}} f_{\text{TLap}}(x) |x|^2 dx$, so 288

$$\int x^{2} f(x) dx$$

$$= 2\frac{1}{B} \int_{0}^{A} e^{-\alpha x} x^{2} dx$$

$$= 2\frac{1}{B} \int_{0}^{A} -\frac{1}{\alpha} x^{2} d\left(e^{-\alpha x}\right)$$

$$= 2\frac{1}{B} (-\frac{1}{\alpha}) A^{2} e^{-\alpha A} + 2\frac{1}{B} \int_{0}^{A} \frac{1}{\alpha} e^{-\alpha x} 2x dx$$
(5)

$$\int_{0}^{A} \frac{1}{\alpha} e^{-\alpha x} 2x dx$$

$$= -\int_{0}^{A} \frac{1}{\alpha^{2}} \cdot 2x d \left(e^{-\alpha x}\right)$$

$$= -\frac{1}{\alpha^{2}} 2A e^{-\alpha A} + \int_{0}^{A} \frac{2}{\alpha^{2}} e^{-\alpha x} dx$$

$$= -\frac{1}{\alpha^{2}} 2A \cdot e^{-\alpha A} + \frac{2}{\alpha^{3}} \left(1 - e^{-2\alpha A}\right)$$
(6)

Thus, we have

١

and

$$\begin{aligned} W &= -2\frac{1}{\alpha}\frac{1}{B}A^2e^{-\alpha A} - 4\frac{1}{\alpha^2}\frac{1}{B}Ae^{-\alpha A} \\ &+ 4\frac{1}{\alpha^3}\frac{1}{B}\left(1 - e^{-\alpha A}\right) \\ &= -2\frac{1}{\alpha}\frac{1}{B}Ae^{-\alpha A}\left(A + 2\frac{1}{\alpha}\right) + 2\frac{\Delta_1^2}{\varepsilon^2} \\ &< 2\frac{\Delta_1^2}{\varepsilon^2}. \end{aligned}$$
(7)

Thus, in total we have $\mathbb{E} \| \mathcal{A}(w) - w \|_2^2 \leq \frac{2d\Delta_1^2}{\epsilon^2} = \frac{8d^2C^2}{\epsilon^2}$. Next for Laplacian mechanism in Theorem 1 we

Next for Laplacian mechanism in Theorem 1 we have $\mathbb{E}[||\mathcal{A}_{lap}(w) - w||_2^2] = \frac{2d\Delta_1^2}{\epsilon^2}$. Thus the variance of high dimensional truncated Laplacian is always lower than Laplacian.

Similarly, the variance of Gaussian mechanism in Theorem 4 is $\frac{8C^2d(\ln 1.25+\ln 1/\delta)}{\epsilon^2}$, we can easily see that our mechanism has lower variance when $\delta \leq \frac{1}{\epsilon^d}$.

5 Experiments

In this section, we conduct experiments for our method based on two parts: DP text re-write for fine-tuning (private embedding) and downstream tasks (sentiment analysis). In all experiments, we compare our Truncated Laplace (TrLaplacian) mechanism with Gaussian and Laplacian mechanisms.



Figure 3: Privacy Test. Curves of N_w value w.r.t. privacy budget ϵ for Yahoo dataset.

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

329

331

332

333

334

335

337

338

339

340

341

342

343

344

5.1 Experimental Setup

Datasets: For the DP text re-write task, we use the Yelp¹ and Yahoo (Yang et al., 2019) datasets. The Yelp Open dataset is a subset of Yelp business, review, and user data with a training size of 8,539 and a testing size of 2,174. The Yahoo dataset contains 14,180 news articles and 34,022 click events. All data are collated to obtain a training, validation, and testing set segmented by sentences. For downstream tasks, we use the SST-2 dataset (Socher et al., 2013), from which we use 68221 heavily polarized reviews from the Internet Movie Database. We divide the SST-2 dataset into an 80:20 ratio for training and testing. The training set consists of 54,576 reviews, with 30,362 positive reviews and 24,214 negative reviews. The testing set consists of 13,645 reviews, with 7,651 positive reviews and 5,994 negative reviews. The statistics of the datasets are presented in Table 7 in Appendix.

Baseline: For DP text re-write, although Krishna et al. (2021) uses the Laplacian mechanism to the sentence level DP instead of word level as in Definition 1. However, as Habernal (2021) mentioned, the approach in (Krishna et al., 2021) is not DP. Thus, here we will not compare with their method, and we will use the Laplacian and Gaussian mechanisms for the clipped embedding as baseline methods. For private fine-tuning, as we mentioned previously, all the previous methods only focus on metric DP instead of the original DP in Definition 1. Thus, our method is incomparable with theirs, and we will still use Laplacian and Gaussian mechanisms as baselines.

293

292

- 29
- 296

297 298

29

301

302

¹https://www.yelp.com/dataset/

Table 1: **Privacy Test.** Performance under GloVe Embedding initialization for the non-private case ($\epsilon = \infty$) and the three mechanisms, where the privacy budget ranges from 0.05 to 0.5. \uparrow means a higher value under this metric indicates better results, and \downarrow means the opposite. The best performance is **bolded**.

		Original		Gau	ssian			Lapl	acian		TrLaplacian			
Privac	y budget ϵ	∞	0.05	0.1	0.2	0.5	0.05	0.1	0.2	0.5	0.05	0.1	0.2	0.5
	Loss↓	2.95	51.25	26.66	9.92	5.97	51.43	37.86	15.35	7.31	2.89	2.86	2.84	3.04
	Rouge1↑	92.37	14.01	59.52	83.61	89.06	13.02	43.30	75.77	86.98	92.44	92.43	92.41	92.25
Yahoo	BLEU↑	8.501	9.286	8.418	8.489	8.499	9.132	8.287	8.474	8.493	8.499	8.500	8.497	8.504
	$N_w\uparrow$	0.703	0.072	0.511	0.595	0.628	0.066	0.334	0.566	0.642	0.706	0.682	0.666	0.662
	BERT-S↑	0.975	0.849	0.908	0.955	0.963	0.839	0.889	0.942	0.959	0.976	0.971	0.971	0.971
	Loss↓	3.07	34.67	21.62	10.61	5.98	36.00	34.64	14.86	7.38	2.98	2.99	3.02	2.94
	Rouge1↑	89.40	15.97	48.89	76.48	84.97	12.60	14.68	66.62	82.08	89.45	89.47	89.34	89.54
Yelp	BLEU↑	8.934	8.976	8.850	8.926	8.930	8.607	8.916	8.913	8.928	8.931	8.935	8.936	8.936
	$N_w\uparrow$	0.706	0.144	0.381	0.608	0.694	0.052	0.138	0.525	0.646	0.705	0.721	0.722	0.725
	BERT-S↑	0.973	0.874	0.895	0.943	0.964	0.855	0.874	0.927	0.952	0.971	0.973	0.971	0.972

Table 2: Utility Test. Comparison of classification accuracy with three embedding methods (Random, GloVe and fastText) for different mechanisms under various privacy budget via sentiment analysis task over the SST-2 dataset.

	Random(seed = 42)				GloVe			fastText			
Privacy budget ϵ	TrLaplace	Laplace	Guassian	TrLaplace	Laplace	Guassian	TrLaplace	Laplace	Guassian		
0.05	86.04	85.97	84.93	88.68	88.57	88.57	89.37	89.37	89.40		
0.1	85.44	84.89	84.06	88.95	88.25	88.24	89.51	89.50	89.30		
0.2	86.18	85.45	85.90	88.93	88.51	88.76	89.45	89.35	89.19		
0.5	86.33	85.55	85.34	88.88	88.48	88.60	89.51	89.40	89.18		

Evaluation Metrics: We use the loss of crossentropy to measure the performance of language models. Specifically, cross-entropy is mainly used to determine how similar the actual output is to the expected output. Smaller model loss indicates less noise added to perturb the text. Additionally, we will use Rouge1 and BLEU scores. Rouge1 (Lin, 2004) calculates recall using standard results and the number of 1-grams co-occurring in the auto-generated text. Similarly, BLEU (Papineni et al., 2002) measures the similarity between standard results and automatically generated text. Rouge1 measures word-level accuracy, while BLEU measures sentence fluency. Moreover, we use BERTScore (Zhang* et al., 2020) to measure the semantic similarity of the perturbed sentence with the original one. To measure the privacypreserving ability, we use the percentage of N_w (Feyisetan et al., 2020), which is the number of words that are not replaced. Thus, under the same privacy budget, larger N_w will be better (we want to change fewer words for accuracy).

Setup: As an embedding can be considered as
an initialization of the model, here we will consider three different initialization: Random embedding (Wieting and Kiela, 2019), GloVe (Pennington et al., 2014) and fastText (Bojanowski et al., 2017).
We conduct experiments on these embeddings and

the subsequent fine-tuning in the DP model via our mechanism. Each pre-trained word embedding is a 300-dimensional vector, and the size of considered vocabulary is 10^4 . For privacy budget, we set $\delta = \frac{1}{4^d}$, and we consider both the high privacy regime where $\epsilon \in \{0.05, 0.1, 0.2, 0.5\}$ and the low privacy regime $\epsilon \in \{1, 5, 10, 20\}$. For large ϵ we will use our previous dummy dimension trick $(d = 500 \text{ for } \epsilon = 10 \text{ and } d = 1700 \text{ for } \epsilon = 20).$ 373

374

375

376

378

379

381

382

384

385

386

390

391

392

394

395

396

397

398

399

5.2 Privacy Experiment on Embedding

We first show the results on private embedding. Specifically, we use GloVe or fastText for initialization, and then use three different private embedding mechanisms with different privacy budgets. Noted that large $\epsilon > 10$ is meaningless for privacy, we concentrated more on a small privacy budget in the main context. Fig. 1 and 5 show the text after projecting the clipped and perturbed embedding back to the word domain in step 4 of Algorithm 1 for different mechanisms when $\epsilon = 0.1$. We can see our method (TrLaplace) outperforms the other two methods from both privacy and semantic perspectives, while the Gaussian mechanism fails to obfuscate the time, and the Laplacian mechanism totally replaces the time by another word, which destroys the structure of the sentence.

Tab. 1 and Tab. 3 are the results on different



Figure 4: Privacy-Utility Test. Curves of Loss, Rouge1 and BERTScore with different privacy budget ϵ for Yelp (Upper) and Yahoo (Lower) datasets.

Table 3: **Privacy Test.** Performance under GloVe Embedding initialization for the non-private case ($\epsilon = \infty$) and the three mechanisms, where the privacy budget ranges from 1 to 20. \uparrow means a higher value under this metric indicates better results, and \downarrow means the opposite. The best performance is **bolded**.

		Original		Gau	ssian			Laplacian				TrLaplacian			
Privac	y budget ϵ	∞	1	5	10	20	1	5	10	20	1	5	10	20	
	Loss↓	2.95	4.28	3.01	3.03	2.98	4.93	3.24	3.05	3.13	2.85	2.97	2.92	2.81	
	Rouge1↑	92.37	90.97	92.27	92.16	92.19	90.02	92.09	92.28	92.26	92.41	92.35	92.24	92.45	
Yahoo	BLEU↑	8.501	8.501	8.501	8.499	8.500	8.503	8.501	8.502	8.500	8.498	8.501	8.499	8.499	
	N_w \uparrow	0.703	0.637	0.680	0.664	0.672	0.660	0.658	0.675	0.655	0.674	0.670	0.702	0.680	
	BERT-S↑	0.975	0.968	0.973	0.971	0.972	0.966	0.970	0.971	0.971	0.974	0.972	0.975	0.974	
	Loss↓	3.07	4.74	3.14	3.13	2.97	5.02	3.30	3.66	3.17	2.93	3.03	3.00	2.98	
	Rouge1↑	89.40	86.63	89.13	89.27	89.80	86.43	89.04	88.15	89.23	89.68	89.40	89.37	89.60	
Yelp	BLEU↑	8.934	8.933	8.936	8.933	8.944	8.931	8.932	8.933	8.934	8.934	8.931	8.934	8.938	
	$N_w\uparrow$	0.706	0.708	0.725	0.708	0.739	0.691	0.721	0.704	0.699	0.724	0.700	0.712	0.740	
	BERT-S↑	0.973	0.969	0.975	0.975	0.975	0.964	0.969	0.969	0.968	0.975	0.971	0.976	0.976	

Table 4: Utility Test. Comparison of classification accuracy with three embedding methods (Random, GloVe and fastText) and different mechanisms under various privacy budget via sentiment analysis task in SST-2 dataset.

	Rand	lom(seed =	= 42)		GloVe			fastText			
Privacy budget ϵ	TrLaplace	Laplace	Guassian	TrLaplace	Laplace	Guassian	TrLaplace	Laplace	Guassian		
1	85.99	84.05	85.36	89.01	88.61	88.62	89.19	89.18	89.08		
5	85.90	85.27	85.31	88.76	88.76	88.47	89.46	89.43	89.20		
10	85.27	84.98	84.57	89.15	88.52	88.48	89.68	89.45	89.53		
20	85.75	85.44	84.12	88.75	88.40	88.57	89.45	89.40	89.24		

Table 5: **Privacy Test.** Performance under fastText Embedding initialization for the non-private case ($\epsilon = \infty$) and three mechanisms (Gaussian, Laplacian and TrLaplacian) on Yelp dataset. The privacy budget ranges from 0.05 to 20. \uparrow means a higher value under this metric indicates better results, and \downarrow means the opposite. The best performance is **bolded**.

	Original		Gau	ssian				Lapla	acian			TrLaplacian			
Privacy budget ϵ	∞	0.05	0.1	0.2	0.5	0.0	5	0.1	0.2	0.5	0.05	0.1	0.2	0.5	
Loss↓	3.35	35.01	29.33	9.31	4.50	36.2	3	29.69	17.15	5.58	1.20	1.20	1.26	1.23	
Rouge1↑	87.8	12.72	28.68	77.95	86.90	10.9	9	27.96	58.97	85.16	92.43	92.67	92.29	92.43	
BLEU↑	8.929	8.226	8.745	8.918	8.931	8.99	8	8.681	8.898	8.931	8.937	8.938	8.937	8.938	
N_w \uparrow	0.713	0.138	0.232	0.661	0.765	0.05	8	0.225	0.484	0.753	0.813	0.807	0.804	0.813	
BERT-S↑	0.967	0.864	0.873	0.945	0.966	0.85	7	0.867	0.908	0.962	0.981	0.978	0.979	0.978	
	Original		Gau	ssian				Lapla	acian			TrLap	olacian		
Privacy budget ϵ	∞	1	5	10	20	1		5	10	20	1	5	10	20	
Loss↓	3.35	3.10	1.68	1.48	1.29	3.6)	1.55	1.53	1.51	1.22	1.25	1.28	1.27	
Rouge1↑	87.8	89.47	92.06	92.40	92.49	88.1	7	91.87	91.90	91.91	92.42	92.35	92.34	92.31	
BLEU↑	8.929	8.936	8.937	8.936	8.936	8.93	5	8.937	8.936	8.934	8.938	8.939	8.937	8.938	
$N_w\uparrow$	0.713	0.794	0.809	0.804	0.813	0.75	8	0.801	0.795	0.792	0.807	0.802	0.800	0.808	
BERT-S↑	0.967	0.976	0.977	0.978	0.980	0.96	7	0.978	0.976	0.977	0.979	0.978	0.978	0.980	

400 metrics regarding private embedding with Glove initialization and Tab. 5 is with fastText initializa-401 tion. We also present the detailed trends w.r.t ϵ for 402 three mechanisms in Fig. 4. When $\epsilon < 1$, from 403 Tab. 1 we can see that for both Yahoo and Yelp, 404 the loss of Gaussian and Laplacian mechanisms 405 will be catastrophically large while our mechanism 406 has a much smaller loss. From Tab. 3 we can 407 see we have almost the same phenomenon when 408 in the low privacy regime. Moreover, for Rouge1, 409 Trlaplacian also leads the other two mechanisms 410 for both datasets, which means our mechanism led 411 to steady superiority from lexical/syntactic aspects. 412 For BLEU, the gap between all three mechanisms 413 to the non-private case becomes small for both two 414 datasets. But our method still has a slight advantage 415 compared with the other two. For N_w value, we 416 can see in Fig. 2 and Fig. 3, our mechanism outper-417 forms the other two mechanisms by changing less 418 percentage of words to achieve the same privacy 419 level, which indicates our method can exactly find 420 sensitive words without hurting other words, thus 421 keeps semantic properties. For BERTScore, our 422 mechanism is almost the same as the non-private 423 case, while there is a larger gap for others. It is no-424 table that, in almost all experiments our mechanism 425 is the best, and the Gaussian mechanism is better 426 than the Laplacian mechanism, which matches our 427 theorem. However, it becomes less obvious when ϵ 428 is large. The main reason is that when ϵ is enough 429 430 large the noise will be sufficiently small and becomes nearly negligible, which can also be sup-431 ported by the proof of Theorem 4. For evaluation 432 metrics, our mechanism may even be better than 433

the non-private case, this may be due to small noise that could improve generalization, which is similar to adversarial training. 434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

6 Utility of Private Fine-tuning

Due to space limitations, the discussion on the Utility of Private Fine-tuning has been moved to the appendix.

Table 6: Results on SST-2 data for classification task with GloVe initialization under $\epsilon = 10$, where 0/1 represents the label and support is the size for each class.

Mechanism	Label	Support	Precision	Recall	F1-score
Guassian			0.87	0.87	0.87
TrLaplace	0	2976	0.88	0.87	0.87
Laplace			0.89	0.84	0.86
Guassian			0.90	0.90	0.90
TrLaplace	1	3847	0.90	0.91	0.90
Laplace			0.88	0.92	0.90

7 Conclusions

We introduce a novel method called the high dimensional truncated Laplacian mechanism for private embedding, which extends the one-dimensional case to the high-dimensional case. Theoretical analysis demonstrates that our method exhibits lower variance compared to existing private word embedding techniques. Experiments show that even in the high privacy regime, our approach incurs only a minimal loss in utility compared to the non-private case, which maintains privacy while preserving the quality of embeddings for promising performance.

453 Limitations

First, the word level DP has the disadvantages of 454 length constraints and linear growth of privacy bud-455 get (Mattern et al., 2022). However, such limita-456 tions are rooted from the definition of DP, instead 457 of our mechanism. Secondly, to ensure DP guaran-458 459 tees, in this paper our mechanism involves clipping embedding vectors and adding calibrated noises, 460 which inevitably introduce errors to the outputs of 461 the task at hand. And these errors may affect dif-462 ferent groups of individuals differently and may 463 464 cause unfairness issues. However, we still need to mention that, such unfairness issues are mainly due 465 to the definition of DP, rather than our method, as 466 DP machine learning algorithms will always have a 467 disparate impact on model accuracy (Bagdasaryan 468 469 et al., 2019).

References

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

504

- Balamurugan Anandan, Chris Clifton, Wei Jiang, Mummoorthy Murugesan, Pedro Pastrana-Camacho, and Luo Si. 2012. t-plausibility: Generalizing words to desensitize text. *Trans. Data Priv.*, 5(3):505–534.
- Eugene Bagdasaryan, Omid Poursaeed, and Vitaly Shmatikov. 2019. Differential privacy has disparate impact on model accuracy. In Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pages 15453–15462.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomás Mikolov. 2017. Enriching word vectors with subword information. *Trans. Assoc. Comput. Linguistics*, 5:135–146.
- Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. 2019. The secret sharer: Evaluating and testing unintended memorization in neural networks. In 28th USENIX Security Symposium, USENIX Security 2019, Santa Clara, CA, USA, August 14-16, 2019, pages 267–284. USENIX Association.
- 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 30th USENIX Security Symposium, USENIX Security 2021, August 11-13, 2021, pages 2633–2650. USENIX Association.
- Ricardo Silva Carvalho, Theodore Vasiloudis, and Oluwaseyi Feyisetan. 2021a. BRR: preserving privacy of text data efficiently on device. *CoRR*, abs/2107.07923.

Ricardo Silva Carvalho, Theodore Vasiloudis, and Oluwaseyi Feyisetan. 2021b. TEM: high utility metric differential privacy on text. *CoRR*, abs/2107.07928. 505

506

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam D. Smith. 2006a. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography, Third Theory of Cryptography Conference, TCC* 2006, New York, NY, USA, March 4-7, 2006, Proceedings, volume 3876 of Lecture Notes in Computer Science, pages 265–284. Springer.
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam D. Smith. 2006b. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography, Third Theory of Cryptography Conference, TCC* 2006, New York, NY, USA, March 4-7, 2006, Proceedings, volume 3876 of Lecture Notes in Computer Science, pages 265–284. Springer.
- Cynthia Dwork and Aaron Roth. 2014. The algorithmic foundations of differential privacy. *Found. Trends Theor. Comput. Sci.*, 9(3-4):211–407.
- Oluwaseyi Feyisetan, Borja Balle, Thomas Drake, and Tom Diethe. 2020. Privacy- and utility-preserving textual analysis via calibrated multivariate perturbations. In WSDM '20: The Thirteenth ACM International Conference on Web Search and Data Mining, Houston, TX, USA, February 3-7, 2020, pages 178– 186. ACM.
- Oluwaseyi Feyisetan, Tom Diethe, and Thomas Drake. 2019. Leveraging hierarchical representations for preserving privacy and utility in text. In 2019 IEEE International Conference on Data Mining, ICDM 2019, Beijing, China, November 8-11, 2019, pages 210–219. IEEE.
- Quan Geng, Wei Ding, Ruiqi Guo, and Sanjiv Kumar. 2020. Tight analysis of privacy and utility tradeoff in approximate differential privacy. In *The 23rd International Conference on Artificial Intelligence and Statistics, AISTATS 2020, 26-28 August 2020, Online* [*Palermo, Sicily, Italy*], volume 108 of *Proceedings* of Machine Learning Research, pages 89–99. PMLR.
- Ivan Habernal. 2021. When differential privacy meets NLP: The devil is in the detail. In *Proceedings of the* 2021 Conference on Empirical Methods in Natural Language Processing, pages 1522–1528, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ivan Habernal. 2022. How reparametrization trick broke differentially-private text representation learning. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022, pages 771–777. Association for Computational Linguistics.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations,*

617

618

656

657

658

659

660

661

662

562 563 ICLR 2015, San Diego, CA, USA, May 7-9, 2015,

Satyapriya Krishna, Rahul Gupta, and Christophe

Dupuy. 2021. ADePT: Auto-encoder based differ-

entially private text transformation. In Proceedings

of the 16th Conference of the European Chapter of

the Association for Computational Linguistics: Main Volume, pages 2435–2439, Online. Association for

Chin-Yew Lin. 2004. ROUGE: A package for auto-

Justus Mattern, Benjamin Weggenmann, and Florian Kerschbaum. 2022. The limits of word level differ-

Radford M. Neal. 2003. Slice sampling. The Annals of

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-

Jing Zhu. 2002. Bleu: a method for automatic evalu-

ation of machine translation. In Proceedings of the

40th Annual Meeting of the Association for Compu-

tational Linguistics, pages 311–318, Philadelphia,

Pennsylvania, USA. Association for Computational

Jeffrey Pennington, Richard Socher, and Christopher D.

Manning. 2014. Glove: Global vectors for word

representation. In Proceedings of the 2014 Confer-

ence on Empirical Methods in Natural Language Pro-

cessing, EMNLP 2014, October 25-29, 2014, Doha,

Qatar, A meeting of SIGDAT, a Special Interest Group

Ildikó Pilán, Pierre Lison, Lilja Øvrelid, Anthi Pa-

padopoulou, David Sánchez, and Montserrat Batet.

2022. The text anonymization benchmark (TAB): A

dedicated corpus and evaluation framework for text

Swaroop Ramaswamy, Rajiv Mathews, Kanishka Rao, and Françoise Beaufays. 2019. Federated learning

for emoji prediction in a mobile keyboard. CoRR,

Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership inference attacks

against machine learning models. In 2017 IEEE Sym-

posium on Security and Privacy, SP 2017, San Jose,

CA, USA, May 22-26, 2017, pages 3-18. IEEE Com-

Richard Socher, Alex Perelygin, Jean Wu, Jason

Chuang, Christopher D. Manning, Andrew Ng, and

Christopher Potts. 2013. Recursive deep models for

semantic compositionality over a sentiment treebank.

In Proceedings of the 2013 Conference on Empiri-

of the ACL, pages 1532-1543. ACL.

anonymization. CoRR, abs/2202.00443.

Association for Computational Linguistics.

ential privacy. CoRR, abs/2205.02130.

matic evaluation of summaries. In Text Summariza-

tion Branches Out, pages 74-81, Barcelona, Spain.

Conference Track Proceedings.

Computational Linguistics.

Statistics, 31(3):705 – 767.

Linguistics.

abs/1906.04329.

puter Society.

- 565
- 566
- 568

- 573 574
- 576
- 578

- 580 581

594

584

590

604

606

610 611

> 612 613 614

cal Methods in Natural Language Processing, pages 1631-1642, Seattle, Washington, USA. Association 615 for Computational Linguistics. 616

- David M. Sommer, Lukas Abfalterer, Sheila Zingg, and Esfandiar Mohammadi. 2021. Learning numeric optimal differentially private truncated additive mechanisms. CoRR, abs/2107.12957.
- John Wieting and Douwe Kiela. 2019. No training required: Exploring random encoders for sentence classification. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net.
- Nan Xu, Oluwaseyi Feyisetan, Abhinav Aggarwal, Zekun Xu, and Nathanael Teissier. 2021a. Densityaware differentially private textual perturbations using truncated gumbel noise. In Proceedings of the Thirty-Fourth International Florida Artificial Intelligence Research Society Conference, North Miami Beach, Florida, USA, May 17-19, 2021.
- Zekun Xu, Abhinav Aggarwal, Oluwaseyi Feyisetan, and Nathanael Teissier. 2021b. On a utilitarian approach to privacy preserving text generation. CoRR, abs/2104.11838.
- Ze Yang, Can Xu, Wei Wu, and Zhoujun Li. 2019. Read, attend and comment: A deep architecture for automatic news comment generation. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 5077-5089, Hong Kong, China. Association for Computational Linguistics.
- Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. Clinical text classification with rule-based features and knowledge-guided convolutional neural networks. BMC Medical Informatics Decis. Mak., 19-S(3):31-39.
- 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 Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021, pages 3853–3866, Online. Association for Computational Linguistics.
- Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In International Conference on Learning Representations.

Α More Details and Experiments

Dataset The statistics of dataset are shown in tab 7.

Table 7: Dataset statistics.

Dataset	Avg. Length (tokens)	Train Size (neg/pos)	Test Size (neg/pos)
Yahoo	181	8539/8673	2174/2189
Yelp	19	3610/3310	909/912
SST-2	10	24214/30362	5994/7651

10

Table 8: **Time Cost.** Comparison of the time cost of each epoch (seconds) under GloVe Embedding initialization for the non-private case and three mechanisms (Gaussian, Laplacian and TrLaplacian), the privacy budget ranges from 0.05 to 20.

			ε <	< 1			$\epsilon \ge$	<u>≥</u> 1	
Priva	cy budget ϵ	0.05	0.1	0.2	0.5	1	5	10	20
	Non-private				11				
Vahaa	Gaussian	111	113	111	111	111	111	111	111
1 21100	Laplacian	111	113	111	111	111	111	111	111
	TrLaplacian	123	123	123	123	123	123	123	123
	Non-private				1	11			
Voln	Gaussian	38	37	38	38	37	37	37	37
reip	Laplacian	38	37	37	37	37	37	37	37
	TrLaplacian	46	41	46	42	42	42	42	42

Implementation Details Models in this paper are implemented based on the PyTorch ² and TensorFlow ³ with their libraries. Experiments are conducted on NVIDIA GeForce RTX 3090 GPUs. To implement our mechanism, we use the acceptancerejection sampling method (Neal, 2003) to sample a point from the high dimensional truncated Laplace distribution from the Laplace distribution, only by rejecting the samples outside the interval.

670

674

675

678

679

681

693

697

698

For text re-write, we use the auto-encoder model. The embedding is initialized with the 300-dimensional pre-trained Random, GloVe, and fast-Text word embedding. We use one-layer BiLSTM with dropout for encoder, and using setup: dropout rate 0.5, Adam (Kingma and Ba, 2015) with an initial learning rate of 0.0005 and betas (0.5, 0.999), batch size 1024, and number of training epochs 50. For the downstream classification task over the IMDB data, we use Adam with an initial learning rate of 10^{-3} , dropout rate 0.2. We set the maximum number of epochs to be 20.

B Utility of Private Fine-tuning

we present the classification accuracy results for private fine-tuning across various embeddings and privacy levels in Tab. 2, Tab. 4 and Tab. 6. It is evident that our mechanism consistently outperforms the other two methods for all embeddings. Furthermore, our approach achieves results that are comparable to the non-private case, where the accuracy scores are 90.14 for Random, 90.19 for GloVe, and 90.19 for fastText in non-private cases. Importantly, the efficacy of our approach will become even more pronounced when dealing with larger datasets. This can be attributed to the minimal amount of noise that the TrLaplacian mechanism

²https://pytorch.org/

requires, thereby preserving the utility of the em-699 bedding. Tab. 6 shows that for class 0, our method 700 achieves significant improvement in accuracy com-701 pared with the other two methods. And for class 702 1, the precision of our method is higher than the 703 others. In Tab. 8 in the Appendix, we show the 704 time cost of each epoch for each experiment, and 705 we can see that compared with Laplacian and Gaus-706 sian mechanisms, our method does not need too 707 much additional time, which means our mechanism 708 is also efficient. 709

Comparison Semantic Problem of Private Embedding										
Original :	do	not	co	me here!	foo	od p	oisoning al	ert!	(→Neg.)	
Trlaplace:	do	not	co	me here!	foo	od p	oisoning al	ert!	(→Neg.)	
Laplace:	this	place	is	awesom	e!	love	this place	! (→Pos.)	
Gaussian:	do	not	go	here! fo	ood	glori	ous <unk></unk>	! (-	→Pos.)	

Figure 5: Another example of text re-write with different mechanisms with $\epsilon = 0.1$. The Gaussian and Laplacian mechanism destroyed semantic properties of original sentence.

³https://www.tensorflow.org/