

# Towards Differentially Private Text Representations

Lingjuan Lyu

National University of Singapore  
lyulj@comp.nus.edu.sg

Xuanli He

Monash University  
xuanli.he1@monash.edu

Yitong Li

The University of Melbourne  
yitongl4@student.unimelb.edu.au

Tong Xiao

Northeastern University  
xiaotong@mail.neu.edu.cn

## ABSTRACT

Most deep learning frameworks require users to pool their local data or model updates to a trusted server to train or maintain a global model. The assumption of a trusted server who has access to user information is ill-suited in many applications. To tackle this problem, we develop a new deep learning framework under an untrusted server setting, which includes three modules: (1) embedding module, (2) randomization module, and (3) classifier module. For the randomization module, we propose a novel local differentially private (LDP) protocol to reduce the impact of privacy parameter  $\epsilon$  on accuracy, and provide enhanced flexibility in choosing randomization probabilities for LDP. Analysis and experiments show that our framework delivers comparable or even better performance than the non-private framework and existing LDP protocols, demonstrating the advantages of our LDP protocol.

## CCS CONCEPTS

• **Computing methodologies** → **Artificial intelligence; Natural language processing; Neural networks**; • **Security and privacy** → **Domain-specific security and privacy architectures; Privacy protections.**

## KEYWORDS

Privacy-preserving; neural representations; natural language processing.

## ACM Reference Format:

Lingjuan Lyu, Yitong Li, Xuanli He, and Tong Xiao. 2020. Towards Differentially Private Text Representations. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '20)*, July 25–30, 2020, Virtual Event, China. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/3397271.3401260>

## 1 INTRODUCTION

The proliferation of deep learning (DL) has led to notable success in natural language processing (NLP), meanwhile, a series of privacy and efficiency challenges arise [3, 9, 11]. In NLP tasks, the input text often provides sufficient clues to portray the authors, such as

their genders, ages, and other important attributes. Concretely, sentiment analysis tasks often impose privacy-related implications on the authors whose text is used to train models, and user attributes can be easily detectable from online review data, as evidenced by [9]. Private information can take the form of key phrases explicitly contained in the text. However, it can also be implicit [16]. For example, the input representation after the embedding layer, or the intermediate hidden representation may still carry sensitive information which can be exploited for adversarial usages. It has been justified that an attacker can recover private variables with higher than chance accuracy, using only the hidden representation [3, 11]. Such attack would occur in scenarios where end users send their learned representations to the cloud for grammar correction, translation, or text analysis tasks [11].

To protect privacy, previous efforts resorted to a trusted aggregator to ensure centralized DP (CDP) [1]. On the other hand, when participants are reluctant to directly share their crowd-sourced data with the server, federated learning becomes a promising learning paradigm that pushes model training to the edge [14]. However, running complex deep neural networks (DNNs) with millions of parameters comes with resource limitations and user experience penalties. Moreover, most federated learning frameworks still assume a trusted aggregator who can have access to local model parameters or gradients [14]. The recent work pointed out the limitation of the trusted server and the associated privacy issues [12, 13]. Without an untrusted server, Shokri and Shmatikov [17] proposed to blur local model gradients by adding noise using differential privacy. However, their privacy bounds are given per-parameter, the gigantic amount of model parameters prevents their technique from providing a meaningful privacy guarantee. The other cryptograph-based methods can be resource-hungry or overly complex for users [2]. More recently, Li et al. [11] and Coavoux et al. [3] proposed to train deep models with adversarial learning. However, both works provide only empirical privacy, without any formal privacy guarantees.

To address the aforementioned problems, we are inspired to take a different approach by utilizing LDP. **Our contributions** include:

- We are the first to train on differentially private crowd-sourced representations for NLP tasks. We propose a novel LDP protocol to preserve the privacy of the extracted representation from user inputs. It offers enhanced flexibility in choosing the randomization probabilities in LDP.
- Experimental results on various NLP tasks show that our framework delivers comparable or even better performance than the non-private framework, and our LDP protocol demonstrates advantages over the existing LDP protocols.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*SIGIR '20, July 25–30, 2020, Virtual Event, China*

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-8016-4/20/07...\$15.00

<https://doi.org/10.1145/3397271.3401260>

## 2 PRELIMINARIES AND RELATED WORK

### 2.1 Local Differential Privacy (LDP)

For the scenario where data are sourced from multiple individuals, while the server is untrusted, LDP [7] is needed to enable data owners to perturb their private data before publication. LDP roots in randomized response [20], and it has been deployed in many real-world applications such as Google’s Chrome browser, Apple’s iOS, and US Census Bureau. A formal definition of LDP is provided in Definition 2.1.

**DEFINITION 2.1.** *A randomized algorithm  $\mathcal{A}$  satisfies  $\epsilon$ -LDP if and only if for any two input tuples  $v$  and  $v'$ , we have*

$$\Pr\{\mathcal{A}(v) = o\} \leq \exp(\epsilon) \cdot \Pr\{\mathcal{A}(v') = o\}$$

for  $\forall o \in \text{Range}(\mathcal{A})$ , where  $\text{Range}(\mathcal{A})$  denotes the set of all possible outputs of the algorithm  $\mathcal{A}$ .

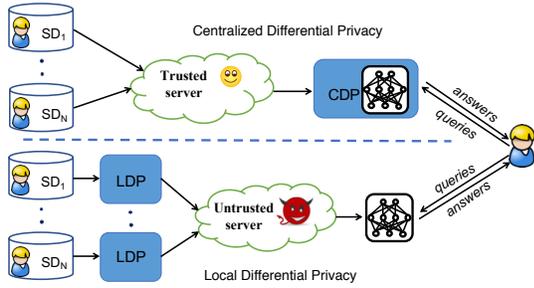


Figure 1: Deep Learning with CDP and LDP.

Compared to CDP [1, 8], LDP offers a stronger level of protection. As illustrated in Figure 1, in DL with CDP, the trusted server owns the data of all users [1], and the server implements CDP algorithm before answering queries from end users. This approach can pose a privacy threat to data owners when the server is untrusted. Moreover, DL algorithms with CDP are inherently computationally complex [1]. By contrast, in DL with LDP, data owners are willing to contribute their data for social good, but do not fully trust the server, so it necessitates data perturbation before releasing it to the server for further learning.

### 2.2 LDP Protocols

The most relevant LDP protocol is called Unary Encoding (UE), which consists of two steps [19]:

**Encoding.** Any single input  $v$  is encoded into a  $d$ -bit vector ( $d$  is domain size), where only the  $v$ -th bit equals to 1, i.e.,  $\vec{B} = \text{Encode}(v)$ , such that  $B[v]=1$  and  $B[i]=0$  for  $i \neq v$ . Hence each  $d$ -bit vector contains  $d - 1$  zeros and only 1 one, and the maximum difference between two adjacent binary vectors is 2, i.e., sensitivity  $\Delta f = 2$ .

**Perturbing.** Each bit with value 1 is preserved with probability  $p$ , thus,  $\text{Perturb}(\vec{B})$  outputs  $\vec{B}'$  as

$$\Pr[B'[i] = 1] = \begin{cases} p, & \text{if } B[i] = 1 \\ q, & \text{if } B[i] = 0 \end{cases} \quad (1)$$

Here two key parameters in perturbation are  $p = \Pr\{1 \rightarrow 1\}$ , the probability that 1 remains 1 after perturbation, and  $q = \Pr\{0 \rightarrow 1\}$ , the probability that 0 is flipped to 1.

Depending on the choice of  $p$  and  $q$ , UE based LDP protocols can be classified into [19]:

**Symmetric Unary Encoding (SUE):** SUE assumes the probability that a bit of 1 is preserved ( $p$ ) equals the probability that a bit of 0 is preserved ( $1 - q$ ), i.e.,  $p + q = 1$ ,  $p = \frac{e^{\epsilon/\Delta f}}{1 + e^{\epsilon/\Delta f}}$ ,  $q = \frac{1}{1 + e^{\epsilon/\Delta f}}$ .

**Optimized Unary Encoding (OUE):** OUE optimizes SUE by using the optimized choices of  $p, q$  for  $\epsilon$ -LDP. Setting  $p$  and  $q$  can be viewed as splitting  $\epsilon$  into  $\epsilon_1 + \epsilon_2$  such that  $\frac{p}{1-p} = e^{\epsilon_1}$  and  $\frac{1-q}{q} = e^{\epsilon_2}$ . That is,  $\epsilon_1$  and  $\epsilon_2$  are the privacy budgets spent on transmitting 1’s and 0’s respectively. If there are more 0’s than 1’s in the encoded representation, it is reasonable to allocate as much privacy budget for transmitting the 0 bits as possible to maintain utility. In the extreme, setting  $\epsilon_1 = 0$  and  $\epsilon_2 = \epsilon$  gives  $p = \frac{1}{2}$  and  $q = \frac{1}{1 + e^{\epsilon/\Delta f}}$ .

## 3 DEEP LEARNING WITH LDP

### 3.1 Optimized Multiple Encoding (OME)

As both SUE and OUE are dependent on the domain size  $d$ , which may not scale well when  $d$  is large. To remove the dependence on  $d$ , we propose a new LDP protocol called Optimized Multiple Encoding (OME). The key idea is to map each real value  $v_i$  of the embedding vector into a binary vector with a fixed size  $l$ . Therefore, for the extracted embedding vector  $\vec{v} = \{v_1, v_2, \dots, v_r\}$  with  $r$  elements, changing all elements of  $\vec{v}$  results in  $\Delta f = 2r$  in both SUE and OUE, and  $\Delta f = rl$  in OME. To enhance flexibility and utility in OME, we follow the intuition behind OUE [19] to perturb 0 and 1 differently.

In particular, we introduce a randomization factor  $\lambda$  to adjust the randomization probabilities in OME. As implied in Theorem 1, by increasing  $\lambda$ , we can decrease  $q$ , thus increasing the probability of keeping the original 0’s. For the value of  $p$ , we increase the probability of preserving the original 1’s for half of the bit vector while decreasing the corresponding probability for the other half. In this way, OME maintains both privacy and utility.

**THEOREM 1.** *For any inputs  $v, v'$  and any encoded bit vector  $B$  with sensitivity  $rl$ , OME provides  $\epsilon$ -LDP given*

$$p = \Pr\{1 \rightarrow 1\} = \begin{cases} \frac{\lambda}{1+\lambda}, & \text{for } i \in 2n \\ \frac{1}{1+\lambda^3}, & \text{for } i \in 2n + 1 \end{cases} \quad (2)$$

$$q = \Pr\{0 \rightarrow 1\} = \frac{1}{1 + \lambda e^{\frac{\epsilon}{rl}}} \quad (3)$$

**PROOF.** Let  $v$  and  $\vec{B}$  represent an input and its encoded bit representation. Given that  $\vec{B}$  has a sensitivity of  $rl$ , the privacy budget  $\epsilon$  needs to be divided by the sensitivity for each bit. By setting

$$p = \Pr\{1 \rightarrow 1\} = \begin{cases} \frac{\lambda}{1+\lambda}, & \text{for } i \in 2n \\ \frac{1}{1+\lambda^3}, & \text{for } i \in 2n + 1 \end{cases} \quad q = \Pr\{0 \rightarrow 1\} = \frac{1}{1 + \lambda e^{\frac{\epsilon}{rl}}}$$

$$1-p = \Pr\{1 \rightarrow 0\} = \begin{cases} \frac{1}{1+\lambda}, & \text{for } i \in 2n \\ \frac{\lambda^3}{1+\lambda^3}, & \text{for } i \in 2n + 1 \end{cases} \quad 1-q = \Pr\{0 \rightarrow 0\} = \frac{\lambda e^{\frac{\epsilon}{rl}}}{1 + \lambda e^{\frac{\epsilon}{rl}}}$$

Then for any inputs  $v, v'$ , we have

$$\begin{aligned} \frac{\Pr\{\vec{B}|v\}}{\Pr\{\vec{B}|v'\}} &= \frac{\prod_{i=1}^{r_l} \Pr\{B[i]|v\}}{\prod_{i=1}^{r_l} \Pr\{B[i]|v'\}} = \frac{\prod_{i \in 2n} \Pr\{B[i]|v\}}{\prod_{i \in 2n} \Pr\{B[i]|v'\}} \times \frac{\prod_{i \in 2n+1} \Pr\{B[i]|v\}}{\prod_{i \in 2n+1} \Pr\{B[i]|v'\}} \\ &\leq \left( \frac{\Pr\{1 \rightarrow 1\}}{\Pr\{1 \rightarrow 0\}} \times \frac{\Pr\{0 \rightarrow 0\}}{\Pr\{0 \rightarrow 1\}} \right)_{i \in 2n}^{r_l/2} \times \left( \frac{\Pr\{1 \rightarrow 1\}}{\Pr\{1 \rightarrow 0\}} \times \frac{\Pr\{0 \rightarrow 0\}}{\Pr\{0 \rightarrow 1\}} \right)_{i \in 2n+1}^{r_l/2} \\ &= \left( \frac{\lambda}{1+\lambda} \times \frac{\frac{\lambda e^{\frac{\epsilon}{r_l}}}{1+\lambda e^{\frac{\epsilon}{r_l}}}}{\frac{1}{1+\lambda}} \right)^{r_l/2} \times \left( \frac{\frac{1}{1+\lambda^3}}{\frac{\lambda^3}{1+\lambda^3}} \times \frac{\frac{\lambda e^{\frac{\epsilon}{r_l}}}{1+\lambda e^{\frac{\epsilon}{r_l}}}}{\frac{1}{1+\lambda e^{\frac{\epsilon}{r_l}}}} \right)^{r_l/2} = e^\epsilon \end{aligned}$$

□

### 3.2 Framework Realization

As shown in Figure 2, the general setting for our proposed deep learning with LDP consists of three main modules: (1) embedding module outputs a 1-D real representation with length  $r$ ; (2) randomization module produces local differentially private representation; and (3) classifier module trains on the randomized binary representations to generate a differentially private classifier as per the post-processing invariance of DP [8]. The detailed training process is summarised in Algorithm 1.

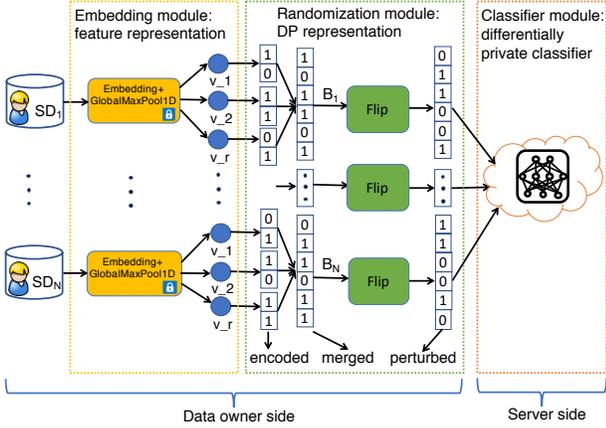


Figure 2: General setting for deep learning with LDP.

#### Algorithm 1 Deep Learning with LDP

- 1: Embedding: Each user maps its input into a 1-D embedding vector (with length  $r$ ) using a pretrained embedding module (GloVe word embeddings, BERT, etc), z-score normalization is applied to avoid large values;
- 2: Encoding: Each element  $v_i \in \mathbb{R}^r$  of the normalized representation is mapped into a binary vector of size  $l = 1 + m + n$ . The first 1 bit represents the sign of the input (1 for negative and 0 for positive), the integer and fraction part are represented by the remaining  $m$  and  $n$  bits respectively;
- 3: Merging and perturbation: Each user merges all the  $r$  binary vectors each with length  $l$  into a long binary vector with total length  $rl$ , which is then randomized by using our proposed OME in Theorem 1;
- 4: Train target model: The server trains a differentially private classifier on all the received noisy representations.

## 4 PERFORMANCE EVALUATION

For performance evaluation, we focus on a range of NLP tasks: 1) sentiment analysis; 2) intent detection; and 3) paraphrase identification. In these tasks, the original sentences might carry some sensitive information such as name entities or monetary descriptions. These private information should be protected, meanwhile the performance for these tasks should not be heavily penalised.

**Datasets.** For sentiment analysis, we use three datasets: IMDB, Amazon, and Yelp, derived from [10], where each review is labelled with a binary sentiment (positive vs. negative). For all sentiment datasets, we perform a train:test split into 8:2.

Intent detection aims to classify each query into seven intents. We derive Intent dataset from [4], which consists of 13,784 training examples and 700 test examples in total.

For paraphrase identification, we use Microsoft Research Paraphrase Corpus (MRPC) from [6]. This task decides whether given two sentences are semantically equivalent. Following Wang et al. [18], we partition this data into train/test (3.7k/1.7k).

**Model and Training.** To extract the intermediate features, we use GloVe word embeddings with a dimension size of 50 [15] for sentiment analysis, and use the pretrained BERT-base [5] for both Intent and MRPC.

For binary encoding of the extracted features, we use 10 bits (1 bit for the sign, 4 bits for the integer part, and 5 bits for the fraction part) to represent each element of the embedding vector.

The classifier module is a multi-layer perceptron with 128 hidden units for sentiment analysis and 768 units for Intent and MRPC. For all datasets, we train the models for 50 epochs using SGD optimizer with learning rate 0.01, decay  $10^{-6}$ , momentum 0.9, and a batch size of 32, and apply a dropout rate of 0.5 on the representation. For each dataset, we average the results over 20 runs.

**Experimental Results.** We first compare our local differentially private NN (LDPNN) with the non-private NN (NPNN), where the randomization module is removed. Table 1 shows that our LDPNN delivers comparable or even better results than the NPNN across various privacy budgets  $\epsilon$  when the randomization factor  $\lambda \geq 50$ . We hypothesise that LDP acts as a regularization technique to avoid overfitting. We conjecture another important reason is that the *enlarged feature space* through encoding produces more powerful representation than the conventional 1-D output of the embedding layer. As LDPNN performance is directly related to the randomization probabilities  $p$  and  $q$ , the higher  $p$  and the lower  $q$ , the lower the randomization of the binary vector, and the better performance will be expected. When the embedding size  $r$  and encoding size  $l$  are fixed,  $p$  is determined by the randomization factor  $\lambda$ , and  $q$  is determined by both the privacy budget  $\epsilon$  and randomization factor  $\lambda$ , as indicated in Equation 2 and Equation 3. Hence we next investigate how  $\epsilon$  and  $\lambda$  impact the model accuracy.

**Impact of  $\epsilon$ .** Contrary to the heuristic study in deep learning with CDP [1], from Table 1, we observe that accuracies are relatively stable when the privacy budget  $\epsilon$  is changed within a wide range of values. The reason lies in the large sensitivity of the encoded binary representation. When  $\lambda$  is a constant, the large sensitivity  $rl$  in OME weakens the effect of  $\epsilon$  on the randomization probabilities, as evidenced by Figure 3 (left),  $p = \{p_1, p_2\}$  and  $q$  of OME keep

**Table 1: Accuracy of NPNN and LDPNN using OME.**

Parameter	Value	Accuracy [%]				
		IMDb	Amazon	Yelp	Intent	MRPC
NPNN		70.67	67.50	66.00	94.17	68.38
$\epsilon$ ( $\lambda = 100$ )	0.5	65.33	68.00	64.54	91.28	66.91
	1	67.33	69.50	66.73	91.35	67.15
	5	67.80	71.00	67.50	91.57	70.10
	10	69.33	72.50	68.10	91.87	70.15
$\lambda$ ( $\epsilon = 1$ )	1	48.00	45.50	42.50	13.00	64.95
	10	64.00	65.46	57.81	85.43	66.66
	50	66.67	69.21	66.50	90.57	66.91
	100	67.33	69.50	66.73	91.35	67.15

nearly consistent when  $\epsilon$  changes. This also explains high accuracy even under a very tight privacy budget (e.g.  $\epsilon=0.5$ ).

We also compare with the other two LDP protocols (SUE and OUE in Section 2.2) on sentiment analysis task. Table 2 shows that our OME significantly outperforms both SUE and OUE. The reason lies in the optimized randomization probabilities of OME, as shown in Figure 3 (left),  $p$  and  $q$  in both SUE and OUE are fluctuating around 0.5, causing low accuracies.

**Table 2: Comparison with other LDP protocols.**

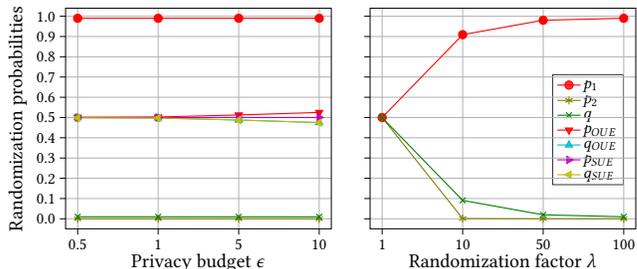
LDP protocols ( $\epsilon = 1$ )	Accuracy [%]		
	IMDb	Amazon	Yelp
SUE	55.33	45.50	48.00
OUE	50.67	54.00	51.00
OME( $\lambda = 100$ , ours)	<b>67.33</b>	<b>69.50</b>	<b>66.73</b>

**Impact of  $\lambda$ .** For randomization factor  $\lambda$ , according to Equation 3, without the randomization factor  $\lambda$ , i.e.,  $\lambda = 1$ , lower  $\epsilon$  values and higher  $rl$  values will result in higher  $q$  values, which may compromise utility. To alleviate this problem, OME calibrates the value of  $\lambda$  to adjust randomization probabilities. As observed in Figure 3 (right), with the increasing  $\lambda$ , OME can largely decrease  $q$  – the probability of perturbing 0 to 1. Although the probability  $p_2$  of preserving the original 1’s decreases for half of the bit vector, the corresponding probability  $p_1$  increases for the other half. This partially explains why OME can maintain both privacy and utility.

Overall, all these results show that our OME offers the reduced impact of the privacy budget  $\epsilon$  on model accuracy, and significantly outperforms the most state-of-the-art LDP protocols.

## 5 CONCLUSION

We formulated a new deep learning framework, which allows data owners to send differentially private representations for further learning on the untrusted servers. A novel LDP protocol was proposed to adjust the randomization probabilities of the binary representation while maintaining both high privacy and accuracy under various privacy budgets. Experimental results on a range of NLP tasks confirm the effectiveness and superiority of our framework.



**Figure 3: Randomization probabilities change with: (left)  $\epsilon$ , (right)  $\lambda$ , where  $p_1 = \Pr\{1 \rightarrow 1\}$  for  $i \in 2n$ ,  $p_2 = \Pr\{1 \rightarrow 1\}$  for  $i \in 2n + 1$ ,  $q = \Pr\{0 \rightarrow 1\}$  in OME.**

## REFERENCES

- [1] Martín Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep learning with differential privacy. In *Proceedings of CCS*. ACM, 308–318.
- [2] Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. 2017. Practical Secure Aggregation for Privacy-Preserving Machine Learning. In *Proceedings of CCS*. ACM, 1175–1191.
- [3] Maximin Coavoux, Shashi Narayan, and Shay B Cohen. 2018. Privacy-preserving neural representations of text. In *Proceedings of EMNLP*. 1–10.
- [4] Alice Coucke, Alaa Saade, Adrien Ball, Théodore Bluche, Alexandre Caulier, David Leroy, Clément Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, et al. 2018. Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces. *arXiv preprint arXiv:1805.10190* (2018).
- [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [6] William B Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*.
- [7] John C Duchi, Michael I Jordan, and Martin J Wainwright. 2013. Local privacy and statistical minimax rates. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*. IEEE, 429–438.
- [8] Cynthia Dwork and Aaron Roth. 2014. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science* 9, 3–4 (2014), 211–407.
- [9] Dirk Hovy, Anders Johannsen, and Anders Søgaard. 2015. User review sites as a resource for large-scale sociolinguistic studies. In *Proceedings of WWW*. 452–461.
- [10] Dimitrios Kotzias, Misha Denil, Nando De Freitas, and Padhraic Smyth. 2015. From group to individual labels using deep features. In *Proceedings of the SIGKDD*. ACM, 597–606.
- [11] Yitong Li, Timothy Baldwin, and Trevor Cohn. 2018. Towards robust and privacy-preserving text representations. In *Proceedings of ACL*. 25–30.
- [12] Lingjuan Lyu, Han Yu, and Qiang Yang. 2020. Threats to Federated Learning: A Survey. *arXiv preprint arXiv:2003.02133* (2020).
- [13] Lingjuan Lyu, Jiangshan Yu, Karthik Nandakumar, Yitong Li, Xingjun Ma, Jiong Jin, Han Yu, and Kee Siong Ng. 2020. Towards Fair and Privacy-Preserving Federated Deep Models. *IEEE TPD* 31, 11 (2020), 2524–2541.
- [14] H Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, et al. 2017. Communication-efficient learning of deep networks from decentralized data. *AISTATS* (2017).
- [15] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of EMNLP*. 1532–1543.
- [16] Daniel Preotiuc-Pietro, Vasileios Lampsos, and Nikolaos Aletras. 2015. An analysis of the user occupational class through Twitter content. In *Proceedings of ACL*, Vol. 1. 1754–1764.
- [17] Reza Shokri and Vitaly Shmatikov. 2015. Privacy-preserving deep learning. In *Proceedings of CCS*. ACM, 1310–1321.
- [18] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018b. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461* (2018b).
- [19] Tianhao Wang, Jeremiah Blocki, Ninghui Li, and Somesh Jha. 2017. Locally differentially private protocols for frequency estimation. In *USENIX Security*. 729–745.
- [20] Stanley L Warner. 1965. Randomized response: A survey technique for eliminating evasive answer bias. *J. Amer. Statist. Assoc.* 60, 309 (1965), 63–69.