

ER-AE: Differentially Private Text Generation for Authorship Anonymization

Haohan Bo

McGill University, Canada
haohan.bo@mail.mcgill.ca

Benjamin C. M. Fung

McGill University, Canada
ben.fung@mcgill.ca

Steven H. H. Ding

Queen's University, Canada
ding@cs.queensu.ca

Farkhund Iqbal

Zayed University, UAE
farkhund.iqbal@zu.ac.ae

Abstract

Most of privacy protection studies for textual data focus on removing explicit sensitive identifiers. However, personal writing style, as a strong indicator of the authorship, is often neglected. Recent studies, such as SynTF, have shown promising results on privacy-preserving text mining. However, their anonymization algorithm can only output numeric term vectors which are difficult for the recipients to interpret. We propose a novel text generation model with a two-set exponential mechanism for authorship anonymization. By augmenting the semantic information through a REINFORCE training reward function, the model can generate differentially private text that has a close semantic and similar grammatical structure to the original text while removing personal traits of the writing style. It does not assume any conditioned labels or paralleled text data for training. We evaluate the performance of the proposed model on the real-life peer reviews dataset and the Yelp review dataset. The result suggests that our model outperforms the state-of-the-art on semantic preservation, authorship obfuscation, and stylometric transformation.

1 Introduction

Privacy has become a vital issue in online data gathering and public data release. Various machine learning models and privacy preservation algorithms have been studied for relational data (Johnson et al., 2018), network graph data (Chen et al., 2014), and transactional data (Li et al., 2012). Some of them have been successfully adopted in real-life applications such as telemetry collection (Cortés et al., 2016). However, the studies on privacy protection for textual data are still preliminary. Most related works only focus on replacing the sensitive key phrases in the text (Vasudevan and John, 2014) without considering the author's writing style, which is indeed a strong indicator of a

person's identity. Even though some textual data, such as double-blind academic reviews, is released anonymously, the adversaries may recover the author's identity using the personal traits in writing. Stylometric techniques (Koppel et al., 2011) can identify an author of the text from 10,000 candidates. They are effective across online posts, articles, emails, and reviews (Ding et al., 2015, 2017). Nevertheless, traditional text sanitization methods (Narayanan and Shmatikov, 2008) focus on anonymizing the contents, such as patient information, instead of the writing style, so they are ineffective against writing style analysis. The original author can be easily re-identified even if protected by these traditional approaches (Iqbal et al., 2008, 2010, 2013; Schmid et al., 2015).

Only a few recent studies focus on authorship anonymization, aiming to hide the personal traits of writing style in the given textual data. *Anonymouth* (McDonald et al., 2012) is a semi-automatic framework that offers suggestions to users to change their writing style. Yet, this framework is not practical since it requires two datasets as a reference to compare the change in writing style. Also, the user has to make all the final modification decisions. *SynTF* (Weggenmann and Kerschbaum, 2018) represents a line of research that protects the privacy of the numeric vector representation of textual data. It adopts the exponential mechanism for a privacy guarantee, but the output is only an opaque term frequency vector, not an interpretable text in natural language. Furthermore, its token substitution approach does not consider the grammatical correctness and semantic.

Style transfer is another line of research that tries to generate text with controllable attributes (Shen et al., 2017; Hu et al., 2017; Sennrich et al., 2016). Representative models (Hu et al., 2017) can control the sentiment and tense of the generated text. However, they do not modify the personal traits in writing. Their applications on sentiment and word-

reordering correspond to the content of the text more than the writing style. We argue that their definition of styles, such as sentiment or tense, is different from the personal linguistic writing characteristics that raise privacy concern. *A4NT* (Shetty et al., 2018) is a generative neural network that sanitizes the writing style of the input text. However, it requires text samples to be labeled with known author identities. It is not applicable to many textual data publishing scenarios. Additionally, according to the samples provided in the paper, it has difficulties keeping the same semantic meaning between the original and the generated text. Without using any privacy model, *A4NT* does not provide any privacy guarantee.

To address the aforementioned issues, we propose an *Embedding Reward Auto-Encoder (ER-AE)* to generate differentially private text. Relying on differential privacy, it protects the author’s identity through text indistinguishability without assuming any specific labels, any parallel data or any assumption on the attacker. It guards the privacy of the data against the worst information disclosure scenario. *ER-AE* receives the original text as input and generates a new text using the two-set exponential mechanism. We propose a *REINFORCE* (Sutton et al., 2000) embedding reward function to augment the semantic information during the text generation process. The model can keep the generated text a close semantic and sentiment similarity to the original while providing a guarantee that one can hardly recover the original author’s identity. Unlike the aforementioned authorship anonymization works, *ER-AE* produces human-friendly text in natural language. Our key contributions are summarized as follows:

- The first differentially private authorship anonymization model that can generate human-friendly text in natural language, instead of a numeric vector.
- A novel two-set exponential mechanism to overcome the large output space issue while producing meaningful results.
- A novel combination of a differential privacy mechanism with a sequential text generator, providing a privacy guarantee through a sampling process.
- A new *REINFORCE* reward function that can augment the semantic information through external knowledge, enabling better preservation of the semantic similarity in the data synthesis

process.

- Comprehensive evaluations on two real-life datasets, namely *NeurIPS & ICLR peer reviews* and *Yelp product reviews*, show that *ER-AE* is effective in obfuscating the writing style, anonymizing the authorship, and preserving the semantics of the original text.

All the source code and data are publicly accessible for reproducibility and transferability.¹

2 Related Work

Differential Privacy. Recently, differential privacy has received a lot of attention in the machine learning community. The deep private auto-encoder (Phan et al., 2016) is designed to preserve the training data privacy. Their purpose is to guarantee that publishing the trained model does not reveal the privacy of individual records. Our purpose is different. We publish the differentially private data generated by the model, rather than the model itself. Most existing models for differentially private data release, such as Chen et al. (2014), focus on different types of data rather than text. One recent work (Weggenmann and Kerschbaum, 2018) aims to protect privacy in text data using the exponential mechanism. However, it releases the term frequency vectors instead of a readable text. This approach limits the utility of published data to only the applications that assume term frequency as features. In contrast, our goal is to generate differentially private text in a natural language without compromising individual privacy.

Writing Style Transfer. Studies on writing style transfer try to change the writing style revealed from the text according to a given author. Shetty et al. (2018) design a GAN to transfer Obama’s text to Trump’s style. A sequence to sequence (seq2seq) model is proposed by Jhamtani et al. (2017) to transfer modern English into Shakespearean English. Shetty et al. (2017) design a model with a cross-alignment method to control the text sentiment while preserving semantic. These models can also be applied to writing style anonymization. However, these studies require the data to be labeled with authorship identity. They assume a number of known authors. In contrast, ours does not assume any label information.

Writing Style Obfuscation. Writing style obfuscation studies try to hide the identity of an au-

¹<https://github.com/McGill-DMaS/Authorship-Anonymization>

thor. Anonymouth (McDonald et al., 2012) is a tool that utilizes *JStylo* to generate writing attributes. It gives users suggestions on which way they can anonymize their text according to two reference datasets. (Kacmarcik and Gamon, 2006) also propose a similar architecture to anonymize text. However, instead of directly changing the text, they all work on the term frequency vector, whose real-life utility is limited. Compared with semi-automatic methods that require users to make a decision, our approach provides an end-to-end solution that directly learns from data.

3 Preliminaries and Problem Definition

Adjacency is a key notion in differential privacy. One of the commonly used adjacency definitions is that two datasets D_1 and D_2 are adjacent if D_2 can be obtained by modifying one record in D_1 (Dwork et al., 2010). Differential privacy (Dwork et al., 2006) is a framework that provides a rigorous privacy guarantee on a dataset. It demands *inherent randomness* of a sanitization algorithm or generation function:

Definition 1. Differential Privacy. Two datasets are considered as adjacent if there is only one single element is different. Let privacy budget $\epsilon > 0$, a randomized algorithm $\mathcal{A} : D^n \rightarrow Z$, and the image of \mathcal{A} : $im(\mathcal{A})$. The algorithm \mathcal{A} is said to preserve ϵ -differential privacy if for any two adjacent datasets $D_1, D_2 \in D^n$, and for any possible set of output $Z \in im(\mathcal{A})$:

$$Pr[\mathcal{A}(D_1) \in Z] \leq e^\epsilon \cdot Pr[\mathcal{A}(D_2) \in Z] \quad \square$$

It guarantees that the result from a given algorithm \mathcal{A} is not sensitive to a change of any individual record in D . ϵ denotes the privacy budget, the allowed degree of sensitivity. A large ϵ implies a higher risk to privacy. However, ϵ is a relative value that implies different degrees of risk given different problems (Weggenmann and Kerschbaum, 2018). Some studies (Sala et al., 2011) use a large ϵ , while the others (Chen et al., 2014) use a smaller value.

Adversary Scenario. Generally in an authorship identification problem, one assumes that the attacker holds an anonymous text authored by one of the suspects from the dataset. The attacker aims to infer the true author of the anonymous text based on a set of reference texts from each suspect. However, this scenario assumes certain information on the applicable dataset, such as author labels and the number of reference text samples. Therefore,

following (Weggenmann and Kerschbaum, 2018), we define that any two pieces of text as adjacent datasets.

Adjacency. Any two pieces of text can be considered *adjacent* in the strictest scenario that datasets D_1 and D_2 both have only one record, and D_2 can be obtained by editing one record in D_1 following Definition 1. With differential privacy, we can have text indistinguishability: one cannot distinguish the identity of any text to another. In our case, the identity of a text corresponds to the author who wrote the text. Along with this, the attacker would fail in the original authorship identification scenario since the anonymous text is indistinguishable from the rest of the dataset.

Our definition follows Weggenmann and Kerschbaum (2018)’s idea, leading to the strictest and most conservative definition of adjacency.

Definition 2. Differentially Private Text Generation. Let \mathbb{D} denote a dataset that contains a set of texts where $x \in \mathbb{D}$ is one of them. $|x|$, the length of the text, is bound by l . Given \mathbb{D} with a privacy budget ϵ , for each x the model generates another text \tilde{x}_{dp} that satisfies ϵl -differential privacy. \square

Following the above definitions, any two datasets that contain only one record are probabilistically indistinguishable w.r.t. a privacy budget ϵ . It directly protects the identity of an individual record, disregarding whether some of the records belong to the same author or not. It assumes that every record is authored by a different author, which is the strictest situation. Technically, the proposed text generation approach protects the writing style by reorganizing the text, replacing tokens with different spelling, removing the lexical, syntactical and idiosyncratic features of the given text. The above definition is based on SynTF (2018), but our target is readable text rather than numeric vectors, which is more challenging.

4 ER-AE for Differentially Private Text Generation

Figure 1 depicts the overall architecture of our proposed ER-AE model, which consists of an encoder and a generator. The encoder receives a sequence of tokens as input and generates a latent vector to represent the semantic features. The generator, which is incorporated with the two-set exponential mechanism, can produce differentially private text according to the latent vector. ER-AE is trained

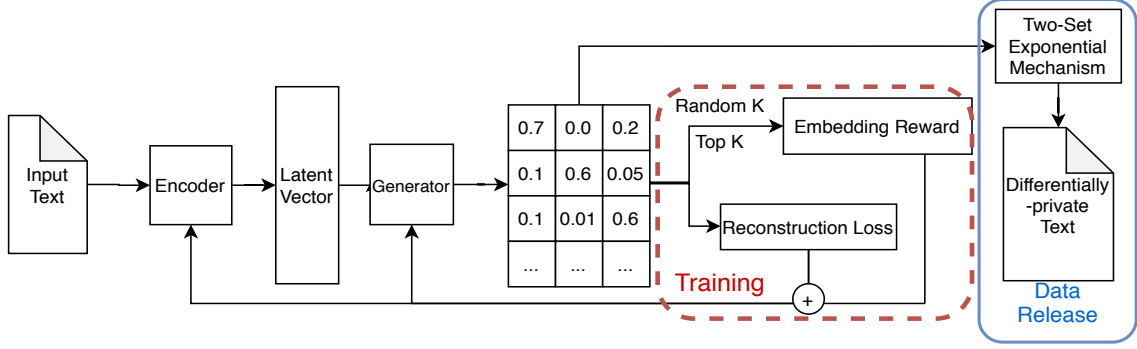


Figure 1: Overall architecture of ER-AE.

by combining a reconstruction loss function and a novel embedding loss function.

Algorithm 1 Generation Procedure of ER-AE

Input: Text: x , Parameters: θ , Encoder: $E_\theta(\cdot)$, Generator: $G_\theta(\cdot)$, Privacy budget: ϵ .
 Produce the latent vector: $E_\theta(x)$.
 Get probabilities of new tokens: $Pr[\tilde{x}] \leftarrow G_\theta(E_\theta(x))$.
for $i \leftarrow 1$ to length of x **do**
 Build two candidate token sets based on $Pr[\tilde{x}]$: S, O .
 Apply exponential mechanism to choose token set: T .
 Randomly sample new i -th token from T : $\tilde{x}_{dp}[i]$.
end for
Output: Differentially Private Text: \tilde{x}_{dp} .

Our ER-AE model starts with a basic sequence-to-sequence (seq2seq) auto-encoder structure. Given a text x , its tokens $\langle x_1 \dots x_l \rangle$ are firstly converted into a sequence of embedding vectors $\langle Em(x_1) \dots Em(x_l) \rangle$ by $Em: \mathcal{V} \rightarrow \mathbb{R}^{m_1}$, where \mathcal{V} is the vocabulary across the dataset and m_1 is the embedding dimension. On its top, we apply a bi-directional recurrent neural network with *Gated Recurrent Unit (GRU)* (Cho et al., 2014) that leverages both the forward and backward information. GRU achieves a comparable performance to LSTM with less computational overhead (Cho et al., 2014). Then, the produced final state vectors from both directions, \mathbf{s}_f and \mathbf{s}_b , are concatenated and linearly transformed to be a latent vector $E(x)$. m is the hidden state dimension for the GRU function.

$$E(x) = \mathbf{W}_h \times \text{concat}(\mathbf{s}_f, \mathbf{s}_b), \quad (1)$$

where $\mathbf{s}_f, \mathbf{s}_b \in \mathbb{R}^m$, $\mathbf{W}_h \in \mathbb{R}^{h \times 2m}$.

The generator is another recurrent neural network with GRU. It generates a text token-by-token. For each timestamp i , it calculates a logit weight z_{iv} for every candidate token $v \in \mathcal{V}$, conditioned on the latent vector, last original token x_{i-1} , and

the last hidden state \mathbf{s}_{i-1} of the GRU function.

$$z_{iv} = \mathbf{w}_v^\top \text{GRU}(E(x), Em(x_{i-1}), \mathbf{s}_{i-1}) + b_v$$

Let \tilde{x}_i denote the random variable for the generated token at timestamp i . Its probability mass function is proportional to each candidate token's weight z_{ti} . This is modeled through a typical softmax function:

$$Pr[\tilde{x}_i = v] = \exp(z_{iv}) / \sum_{v' \in \mathcal{V}} \exp(z_{iv'}) \quad (2)$$

For each timestamp i , a typical seq2seq model generates text by applying $\text{argmax}_{v \in \mathcal{V}} Pr[\tilde{x}_i = v]$. However, this process does not protect the privacy of the original data.

4.1 Differentially Privacy Text Sampling with Two-Set Exponential Mechanism

To protect an individual's privacy and hide the authorship of the original input text, we couple differential privacy mechanism with the above sampling process in the generator. The *exponential mechanism* (McSherry and Talwar, 2007) can be applied to both numeric and categorical data (Fernandes et al., 2018). It is effective in various sampling process for discrete data. It guarantees privacy protection by injecting noise into the sampling process:

Definition 3. Exponential Mechanism. Let \mathcal{M} and \mathcal{N} be two enumerable sets. Given a privacy budget $\epsilon > 0$, a rating function $\rho: \mathcal{M} \times \mathcal{N} \rightarrow \mathbb{R}$. The probability density function of the random variable $\varepsilon_{\epsilon, \rho}(m)$, $Pr[\varepsilon_{\epsilon, \rho}(m) = n]$ is:

$$\frac{\exp\left(\frac{\epsilon}{2\Delta} \rho(m, n)\right)}{\sum_{n'} \exp\left(\frac{\epsilon}{2\Delta} \rho(m, n')\right)} \quad (3)$$

where Δ , the sensitivity, means the maximum difference of rating function values between two adjacent datasets, and $m \in \mathcal{M}, n \in \mathcal{N}$. \square

However, according to Weggenmann and Kerschbaum (2018), the exponential mechanism requires a large privacy budget to produce meaningful results while the output space is large, the vocabulary size in our case. It’s nontrivial to randomly sample a good result directly among 20,000 candidates.

To tackle the large output space issue, inspired by *subsampling exponential mechanism* (Lantz et al., 2015), we propose a *two-set exponential mechanism* to produce meaningful results with a better privacy protection. Instead of using a database independent distribution, we use a model-based distribution to generate subsets of tokens.

Definition 4. Two-Set Exponential Mechanism.

Let \mathcal{V} be an enumerable set with size s . Given the model-based probabilities of each item in \mathcal{V} , $Pr[v]$ for $v \in \mathcal{V}$, an item set S of size k is built by repeatedly sampling proportional to $Pr[v]$ with replacement. Other items are denoted as set O , where $\mathcal{V} = O \cup S$, $\emptyset = O \cap S$. Let $\mathcal{N} = \{S, O\}$. An item set C_{dp} is chosen from \mathcal{N} through the exponential mechanism with a rating function ρ : $\sum_{v \in C} Pr[v] / \sum_{C' \in \mathcal{N}, v' \in C'} Pr[v']$. Given $\epsilon > 0$, $N \in \mathcal{N}$, the probability density function (PDF) of the random variable $\epsilon_{\epsilon, \rho}(C)$, $Pr[\epsilon_{\epsilon, \rho}(C) = N]$, is:

$$\frac{\exp\left(\frac{\epsilon}{2\Delta}\rho(C, N)\right)}{\sum_{N' \in \mathcal{N}} \exp\left(\frac{\epsilon}{2\Delta}\rho(C, N')\right)} \quad (4)$$

After choosing the set, C_{dp} , an item is randomly picked from the chosen set: $v \sim Random(C_{dp})$. Thus, given $v, w \in \mathcal{V}$, $Pr[\epsilon_{\epsilon, \rho}(v) = w]$ is:

$$Pr[wS] * Pr[\epsilon_{\epsilon, \rho}(C) = S|wS] * Pr[w|wS, S] + Pr[wO] * Pr[\epsilon_{\epsilon, \rho}(C) = O|wO] * Pr[w|wO, O],$$

where $Pr[wS]$ and $Pr[wO]$ are respectively the probability of w in set S and O , $Pr[wO] = (1 - Pr[w])^k$. \square

Theorem 1. Two-Set Exponential Mechanism.²

Given a privacy budget $\epsilon > 0$ and the size of output space s , two-set exponential mechanism is $(\epsilon + \ln(s))$ -differentially private. \square

By plugging our model with this mechanism, we have the probability mass function for $\epsilon_{\epsilon, \rho_i}(\tilde{x}_i)$: $Pr[\epsilon_{\epsilon, \rho_i}(\tilde{x}_i) = tk]$. This function models the disturbed probability distribution for all the alternative token tk to replace the original variable. According to Theorem 4, sampling from $\epsilon_{\epsilon, \rho_i}(\tilde{x}_i)$ for each

²The proof is provided in Appendix B

timestamp i is $(\epsilon + \ln(s))$ -differentially private. Recall that in Definition 1, the timestamp is bound by l . To generate text \tilde{x}_{dp} , the generator samples a token for timestamp i through the chosen set T_i :

$$\tilde{x}_{dp}[i] \sim Random(T_i) \quad \text{for } i \in [1, l] \quad (5)$$

The *composition theorem* (Dwork et al., 2014) is an extension to differential privacy. By repeating n ϵ -differentially private algorithms, the complete process achieves an ϵn -differential privacy. Algorithm 1 shows the differentially private text generation of ER-AE. As proved in Appendix A:

Theorem 2. Differentially Private Text Sampling.

Given a privacy budget $\epsilon > 0$, a sequence length $l > 0$, the generator’s sampling function in Eq. 5 is $(\epsilon + \ln(s)) * l$ -differentially private. \square

4.2 Initial Grammar and Semantic Preservation

To generate a human-friendly text that has a close semantic to the original one, we need to have a high-quality rating function ρ_i for Eq. 4. This is achieved by training the ER-AE model’s encoder to extract semantic information, and its generator to learn the relationships among the tokens for prediction. We follow an unsupervised learning approach since we do not assume any label information. First, we adopt the reconstruction loss function:

$$\mathcal{L}_{recon} = \sum_{x_i \in x, x \in \mathbb{D}} -\log Pr[\tilde{x}_i = x_i] \quad (6)$$

It maximizes the probability of observing the original token x_i itself for the random variable \tilde{x}_i . In the recent controllable text generation models, the reconstruction loss plays an important role to preserve grammar structure and semantics of input data (Shetty et al., 2018) when combined with the other loss.

4.3 REINFORCE Training for Semantic Augmentation

Diving into the optimization aspect of the softmax function, the reconstruction loss function above encourages the model to produce a higher probability on the original token while ignoring the rest candidates. It does not consider the other tokens that may have a similar meaning under a given context. This issue significantly limits the variety of usable alternative tokens. Additionally, this loss function relies on a single softmax function for multi-object

learning, it cannot provide the expressiveness required by the language model (Yang et al., 2018). We inspect the candidates and in most of the cases, only the top-ranked token fits the context in the text. This is problematic because the mechanism for our sampling process also relies on the other candidates to generate text. To address the above issue, we propose a novel embedding reward function using the pre-trained word embeddings. Word representation learning models show that discrete text tokens’ semantic can be embedded into a continuous latent vector space. The distance between word embedding vectors can be a reference to measure the similarity between different words. To encourage our rating function ρ_i to learn richer and better substitute tokens, we propose a reward function that leverages the semantics learned from the other corpus. The text dataset to be anonymized and released can be small, and the extra semantic knowledge learned from the other corpus can provide additional reference for our rating function. This reward function is inspired by the Policy Gradient loss function (Sutton et al., 2000), \mathcal{L}_{embed} is:

$$- \sum_{x_i \in x, x \in \mathbb{D}} \left(\sum_{v \in \mathbb{E}_k(\tilde{x}_i)} \log(\text{Pr}[\tilde{x}_i = v])\gamma(x_i, v) + \sum_{w \sim \mathbb{V}_k} \log(\text{Pr}[\tilde{x}_i = w])\gamma(x_i, w) \right)$$

Generally, this reward function assigns credits to the under-rated tokens in the reconstruction loss function. Recall that \mathbb{D} is the original dataset and x is one of its texts. At time step i , this reward function first assigns rewards to the top- k selected tokens, denoted as $\mathbb{E}_k(\tilde{x}_i)$, according to probability estimates for random variable \tilde{x}_i in Eq. 2. The rewards are proportional to their semantic relationship to the original token x_i . It is defined as a function $\gamma : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$, $\gamma(w, v)$ is:

$$\min(\text{cosine}(Em(w), Em(v)), 0.85) \quad (7)$$

The *min* function avoids the generator focusing only on the original token. By assigning rewards to $\mathbb{E}_k(\tilde{x}_i)$, it encourages the other candidates having a close semantic to the targeted one, but it may fail to reach infrequent tokens. Therefore, in the second part of the reward function, we encourage the model to explore less frequent tokens by random sampling candidates as \mathbb{V}_k . This design balances the exploitation (top- k) and the exploration (\mathbb{V}_k) in reinforcement learning.

During training, the model is firstly pre-trained by minimizing the reconstruction loss in Eq. 6 through the *Adam* optimizer, and adopts the embedding reward loss later. The total loss is

$$\mathcal{L} = \lambda_{recon} \times \mathcal{L}_{recon} + \lambda_{embed} \times \mathcal{L}_{embed} \quad (8)$$

Specifically, the reconstruction loss can lead the model to generate grammatically correct text, and the embedding reward loss encourages the model to focus more on semantically similar tokens. The balance of the two loss functions are controlled by λ_{recon} and λ_{embed} .

5 Experiment

All the experiments are carried out on a Windows Server equipped with two Xeon E5-2697 CPUs (36 cores), 384 GB of RAM, and four NVIDIA TITAN XP GPU cards. We evaluate ER-AE on two different datasets with respect to its effectiveness for privacy protection and utility preservation.

- **Yelp Review Dataset**³: All the reviews and tips from the top 100 reviewers ranked by the number of published reviews and tips. It contains 76,241 reviews and 200,940 sentences from 100 authors.
- **Academic Review Dataset**: All the public reviews from NeurIPS (2013-2018) and ICLR (2017) based on the original data and the web crawler provided by (Kang et al., 2018). It has 17,719 reviews, 268,253 sentences, and the authorship of reviews is unknown.

Each dataset is divided into 70/10/20 for train/dev/evaluation respectively. As mentioned in the related work discussion, most of the controllable text generation and style transfer studies rely on known authorship or other labels. Other generation models such as paraphrasing, however, hold an essentially different goal and cannot provide a privacy guarantee on the generated data. They are not applicable to our problem. Therefore, we pick SynTF (Weggenmann and Kerschbaum, 2018) and different generation and sampling models for evaluation:

- **Random Replacement (Random-R)**: This method generates a new text by replacing each token in the text by randomly picking substitution from the vocabulary.
- **AE with Differential Privacy (AE-DP)**: Extended version of AE with the added two-set exponential mechanism for text generation. It does not include the embedding reward.

³http://www.yelp.com/dataset_challenge

Table 1: Results for each evaluation metric on both datasets. \uparrow indicates the higher the better. \downarrow indicates the lower the better.

Model	Yelp (100-author)			Conferences' Dataset	
	USE \uparrow	Authorship \downarrow	Stylometric \uparrow	USE \uparrow	Stylometric \uparrow
Original text	1	0.5513	0	1	0
Random-R	0.1183	0.0188	62.99	0.1356	65.624
AE-DP	0.6163	0.097	11.443	0.614	9.859
SynTF (2018)	0.1955	0.0518	26.3031	0.2161	25.95
ER-AE (ours)	0.7548	0.0979	13.01	0.7424	9.838

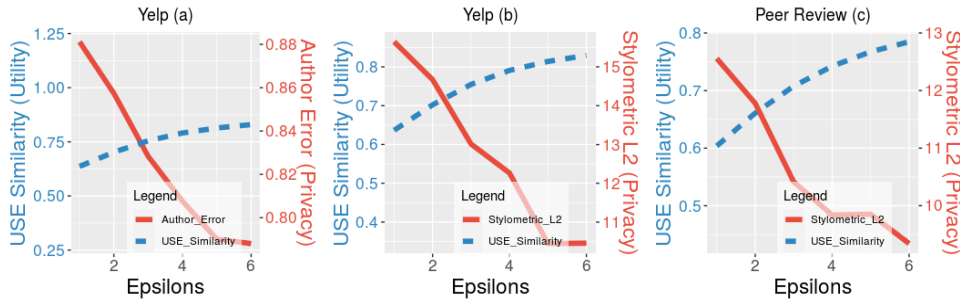


Figure 2: Privacy v.s. Utility. Comparing USE similarity (utility), authorship identification error rate (privacy) and Stylometrics L2 distance (privacy) for different ϵ s on applicable datasets.

Table 2: The intermediate result of top five words and their probabilities at that the third and the fourth generation steps.

Input: there are several unique hot dog entrees to choose.	
	several
AE-DP	several 0.98 , those 0.007, some 0.003 various 0.002, another 0.001
ER-AE	many 0.55, some 0.20, several 0.14 different 0.04, numerous 0.03
	unique
AE-DP	unique 0.99 , different 0.0001, new 3.1e-05, nice 2.5e-05, other 2.1e-05
ER-AE	unique 0.37 , great 0.21, amazing 0.15, wonderful 0.1, delicious 0.05

Table 3: The estimated probability of a good candidate sampled with different mechanisms.

Input: there are several unique hot dog entrees to choose.		
	several	unique
Exponential Mechanism	0.0017	0.00091
Two-Set Exponential Mechanism	0.7411	0.6794

- **SynTF (Weggenmann and Kerschbaum, 2018):** We directly generate the tokens through SynTF’s differentially private sampling function, without further extraction of the frequency vector.

SynTF is a state-of-the-art generation model that satisfies differential privacy property on textual data. The other two simple baselines are for ablation test purposes.

For ER-AE, we adopted a two-layers stacked GRU network for both the encoder and the generator. There are 512 cells in each GRU layer. The vocabulary size is 20,000, separately built for each dataset. All the word embeddings in our model come from the pre-trained *BERT* embeddings provided by (Devlin et al., 2019), which has a dimension of 768 for each embedding. The maximum input length of our model is 50, the learning rate is 0.001, the k for embedding reward loss function is 5, the λ_{recon} is 1, the λ_{embed} is 0.5, and the batch size is 128. The k in two-set exponential mechanism is 5. ER-AE is implemented in TensorFlow (Abadi et al., 2016), and it uses the tokenizer in the NLTK library. Some traditional tricks for text generation, such as beam search, are not mentioned because they are incompatible with differential privacy. All the models are evaluated from three aspects: semantic preservation, privacy protection, and stylometric changes:

- **Semantic Preservation (USE):** A pre-trained *Universal Sentence Embedding (USE)* model⁴ from Google. It can embed a sentence into a latent vector that represents its semantics. It is widely used for supervised NLP tasks such as sentiment analysis. We measure the degree of semantic preservation using the cosine similarity between the latent vector of the original text and

⁴<https://tfhub.dev/google/universal-sentence-encoder/1>

Table 4: Sample sentences generated by models.

Input	the play place is pretty fun for the little ones .
Random-R	routing longtime 1887 somalia pretty anatomical shallow the dedicated drawer rosalie
AE-DP	employer play lancaster mute fish fun for wallace little chandler .
SynTF	conditioned unique catherine marquis governing skinny garment hu vivid . insists
ER-AE	the play place is pretty nice with the little ones !
Input	i also ordered a tamarind margarita and it was great .
Random-R	substantial char recommended excavation tamarind coil longitudinal recover verify great housed
AE-DP	intersection also ordered service tamarind drooling scratched denis monkfish motions .
SynTF	carnage spence unsigned also clinging said originated beacon liking strike accomplishments
ER-AE	i also requested a tamarind margarita and it were great .
Input	i 'm not complaining because you do get exactly what you pay for .
Random-R	substantial char recommended excavation tamarind coil longitudinal recover verify great housed
AE-DP	comic-book 'm not mins because you donnelly get exactly tenderloin nerves bottomless for aldo box
SynTF	leaf penetrated amounted jolted courageous socket fades unwilling tu judges regional numbering
ER-AE	i 'm not disappointing because you do make occult what you pay for .
Input	the manuscript is well written is provides good insight into the problem .
AE-DP	the fig2c is well 1102-103 wish provides horseshoe insight into the problem compositionality
SynTF	ness voice incoming depending entrances somehow priscilla rows romantic oblivious mall
ER-AE	the manuscript is well edited has provides excellent insight into the problem .
Input	in particular , the generality of the approach is very well presented .
SynTF	wife pierced rotate specialist probe elects prussian beatty eccentric sweating .
ER-AE	in particular , this generality of an approach is very well written well

one of the generated text.

- **Privacy Protection (Authorship):** A state-of-the-art authorship identification neural network model (Sari et al., 2017) to identify the authorship of generated text. The model is firstly trained on the training dataset, and the performance is evaluated on the testing set. The author’s privacy is protected if s/he cannot be identified using authorship identification techniques.
- **Stylometric Changes:** Well-established stylistic context-free features such as text length and a number of function words. We adopt *Stylo-Matrix* (Ding et al., 2017) for an aggregation of features in (Iqbal et al., 2013; Zheng et al., 2006). The feature vector change is measured by the difference in L2 norm.

Quantitative Evaluation (Table 1). With a low utility (USE) score around 0.2 for both datasets, SynTF, and Random-R generate grammatically incorrect text and completely change the meaning of the original one. In contrast, ER-AE without semantic augmentation through REINFORCE training, denoted as AE-DP, achieves a much higher utility score of around 0.61. The full model ER-AE, with an ϵ of 3, achieves the highest utility score of 0.75 for Yelp reviews and 0.74 for peer reviews. AE-DP, SynTF, and ER-AE all significantly reduce the chance of a successful authorship identification attack from 55% to lower than 10% in the Yelp data and introduce a variation in stylometric features

of more than 10 in magnitude in the peer review dataset. They are all effective and competitive on removing the personal writing trait from the text data, but as mentioned above, AE-DP achieves the best and a much higher utility score. Although Random-R performs better on privacy protection, its generated texts are irrelevant to the original. Overall, with a competitive performance on anonymization, ER-AE performs significantly better than all of the other models on utility.

Impact of Embedding Reward. Table 4 shows that the embedding reward plays an important role in selecting semantically similar candidates for substitution. AE-DP assigns a large probability to the original token and a tiny probability to the others. If applied with the mechanism, it is more likely to pick a semantically irrelevant token. ER-AE shows a smoother distribution and assigns higher probabilities to top-ranked semantically relevant tokens. Its generated candidates are better.

Case Study. Table 4 shows that both SynTF and Random-R cannot generate human-friendly text. Due to the issue of reconstruction loss function [6], AE-DP cannot substitute token with similarly semantic tokens and destroys the semantic meaning. ER-AE, powered by embedding reward, can substitute some tokens with semantically similar ones: “written” is replaced by “editted”, and the whole sentence still makes sense. Besides, it can preserve the grammatical structure of the input. However,

due to some missing information from word embeddings, the model would fail to generate good candidates for sampling. The third sample replaces “exactly” with “occult”. ER-AE still performs way better than other models.

Utility vs. Privacy. The privacy budget ϵ controls the trade-off between privacy and utility. A larger ϵ implies better utility but less protection on privacy. However, this is a relative value that implies different degrees of risk given different problems (Weggenmann and Kerschbaum, 2018; Fernandes et al., 2018). As proved by Weggenmann and Kerschbaum (2018) a higher ϵ is intrinsically necessary for a large output space, in our case the vocabulary, to generate relevant text. In fact, we have already significantly reduced the optimal ϵ value of 42.5 used by Weggenmann and Kerschbaum (2018) to around 13, given the same dataset. One possible way to lower the bound of ϵ is to directly factor in authorship and utility, such as topics, into the privacy model. However, it limits applicability to datasets.

Exponential Mechanism vs. Two-Set Exponential Mechanism. In Table 3, we estimated the probability of a meaningful token (among top 5 semantically similar tokens) is sampled based on the intermediate probabilities in Table 2. Given a large output space of 20,000, the exponential mechanism is not likely to sample a meaningful token with a probability of 0.01 %. However, the two-set exponential mechanism dramatically improves it from 0.01 % to around 70%. Our generator has a much higher chance to generate meaningful results with a similar privacy budget.

6 Conclusion

In this paper, we propose a novel model, ER-AE, to protect an individual’s privacy for text data release. We are among the first to fuse the differential privacy mechanisms into the sequence generation process. We demonstrate the effectiveness of our model on the Yelp review dataset and two peer reviews datasets. However, we also find that ER-AE performs not very well on long texts due to the privacy budget accounting issue. Our future research will focus on improving long texts generation with better budget allocation scheme.

7 Ethical Considerations

Our model outperforms others on authorship obscuration and semantic preservation. Similar to other

text generation tools, this model may be abused to generate fake reviews, but this can be assuaged by using fake review detection methods. This research work directly contributes to the area of privacy protection and indirectly promotes freedom of speech and freedom of expression in cyberspace.

Acknowledgments

This research is supported in part by Natural Sciences and Engineering Research Council of Canada (NSERC) Discovery Grants (RGPIN-2018-03872), Canada Research Chairs Program (950-232791), and Provost Research Fellowship Award (R20093) and Research Incentive Funds (R18055) from Zayed University, United Arab Emirates. The Titan Xp used for this research was donated by the NVIDIA Corporation.

References

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. Tensorflow: A system for large-scale machine learning. In *Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.
- Rui Chen, Benjamin C. M. Fung, Philip S. Yu, and Bipin C. Desai. 2014. Correlated network data publication via differential privacy. *The VLDB Journal—The International Journal on Very Large Data Bases*, 23(4).
- Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder–decoder approaches. *Syntax, Semantics and Structure in Statistical Translation*.
- Jorge Cortés, Geir E Dullerud, Shuo Han, Jerome Le Ny, Sayan Mitra, and George J Pappas. 2016. Differential privacy in control and network systems. In *IEEE 55th Conference on Decision and Control (CDC)*. IEEE.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*.
- Steven H. H. Ding, Benjamin C. M. Fung, and Mourad Debbabi. 2015. A visualizable evidence-driven approach for authorship attribution. *ACM Transactions on Information and System Security (TISSEC)*, 17(3).

- Steven H. H. Ding, Benjamin C. M. Fung, Farkhund Iqbal, and William K Cheung. 2017. Learning stylistic representations for authorship analysis. *IEEE transactions on cybernetics*, (99).
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. 2006. Calibrating noise to sensitivity in private data analysis. Springer.
- Cynthia Dwork, Aaron Roth, et al. 2014. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3–4).
- Cynthia Dwork, Guy N Rothblum, and Salil Vadhan. 2010. Boosting and differential privacy. In *2010 IEEE 51st Annual Symposium on Foundations of Computer Science*. IEEE.
- Natasha Fernandes, Mark Dras, and Annabelle McIver. 2018. Author obfuscation using generalised differential privacy. *arXiv preprint arXiv:1805.08866*.
- Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P Xing. 2017. Toward controlled generation of text. In *Proceedings of the 34th International Conference on Machine Learning*. JMLR.
- Farkhund Iqbal, Hamad Binsalleeh, Benjamin C. M. Fung, and Mourad Debbabi. 2010. Mining writeprints from anonymous e-mails for forensic investigation. *Digital Investigation*, 7(1-2).
- Farkhund Iqbal, Hamad Binsalleeh, Benjamin C. M. Fung, and Mourad Debbabi. 2013. A unified data mining solution for authorship analysis in anonymous textual communications. *Information Sciences*, 231.
- Farkhund Iqbal, Rachid Hadjidj, Benjamin C. M. Fung, and Mourad Debbabi. 2008. A novel approach of mining write-prints for authorship attribution in e-mail forensics. *Digital Investigation*, 5(1):S42–S51.
- Harsh Jhamtani, Varun Gangal, Eduard Hovy, and Eric Nyberg. 2017. Shakespearizing modern language using copy-enriched sequence to sequence models. In *Proceedings of the Workshop on Stylistic Variation*.
- Noah Johnson, Joseph P Near, and Dawn Song. 2018. Towards practical differential privacy for sql queries. *Proceedings of the VLDB Endowment*, 11(5).
- Gary Kacmarcik and Michael Gamon. 2006. Obfuscating document stylometry to preserve author anonymity. In *Proceedings of the COLING/ACL on Main conference poster sessions*. Association for Computational Linguistics.
- Dongyeop Kang, Waleed Ammar, Bhavana Dalvi, Madeleine van Zuylen, Sebastian Kohlmeier, Eduard Hovy, and Roy Schwartz. 2018. [A dataset of peer reviews \(peerread\): Collection, insights and nlp applications](#). In *Meeting of the North American Chapter of the Association for Computational Linguistics (NAACL)*, New Orleans, USA.
- Moshe Koppel, Jonathan Schler, and Shlomo Argamon. 2011. Authorship attribution in the wild. *Language Resources and Evaluation*, 45(1).
- Eric Lantz, Kendrick Boyd, and David Page. 2015. Subsampled exponential mechanism: Differential privacy in large output spaces. In *Proceedings of the 8th ACM Workshop on Artificial Intelligence and Security*.
- Ninghui Li, Wahbeh Qardaji, Dong Su, and Jianneng Cao. 2012. Privbasis: Frequent itemset mining with differential privacy. *Proceedings of the VLDB Endowment*, 5(11).
- Andrew WE McDonald, Sadia Afroz, Aylin Caliskan, Ariel Stolerman, and Rachel Greenstadt. 2012. Use fewer instances of the letter “i”: Toward writing style anonymization. In *International Symposium on Privacy Enhancing Technologies Symposium*. Springer.
- Frank McSherry and Kunal Talwar. 2007. Mechanism design via differential privacy. In *Proceedings of the 48th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*. IEEE.
- Arvind Narayanan and Vitaly Shmatikov. 2008. Robust de-anonymization of large datasets (how to break anonymity of the netflix prize dataset). *University of Texas at Austin*.
- NhatHai Phan, Yue Wang, Xintao Wu, and Dejing Dou. 2016. Differential privacy preservation for deep auto-encoders: an application of human behavior prediction. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence*.
- Alessandra Sala, Xiaohan Zhao, Christo Wilson, Haitao Zheng, and Ben Y Zhao. 2011. Sharing graphs using differentially private graph models. In *Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference*. ACM.
- Yunita Sari, Andreas Vlachos, and Mark Stevenson. 2017. Continuous n-gram representations for authorship attribution. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*.
- Michael Schmid, Farkhund Iqbal, and Benjamin C. M. Fung. 2015. E-mail authorship attribution using customized associative classification. *Digital Investigation (DIIN)*, 14(1).
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Controlling politeness in neural machine translation via side constraints. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.

Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2017. Style transfer from non-parallel text by cross-alignment. In *Advances in neural information processing systems*.

Rakshith Shetty, Bernt Schiele, and Mario Fritz. 2018. A4NT: author attribute anonymity by adversarial training of neural machine translation. In *Proceedings of the 27th USENIX Conference on Security Symposium*. USENIX Association.

Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. 2000. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*.

Veena Vasudevan and Ansamma John. 2014. A review on text sanitization. *International Journal of Computer Applications*, 95(25).

Benjamin Weggenmann and Florian Kerschbaum. 2018. Syntf: Synthetic and differentially private term frequency vectors for privacy-preserving text mining. In *Proceedings of the 41st ACM International Conference on Research & Development in Information Retrieval (SIGIR)*. ACM.

Zhilin Yang, Zihang Dai, Ruslan Salakhutdinov, and William W Cohen. 2018. Breaking the softmax bottleneck: A high-rank rnn language model. In *International Conference on Learning Representations*.

Rong Zheng, Jiexun Li, Hsinchun Chen, and Zan Huang. 2006. A framework for authorship identification of online messages: Writing-style features and classification techniques. *Journal of the American society for information science and technology*, 57(3).

A Proof of Differentially Private Text Sampling.

Theorem 3. Differentially Private Text Sampling. *Given a privacy budget $\epsilon > 0$, a sequence length $l > 0$, the generator’s sampling function in Eq.7 is $(\epsilon + \ln(s)) * l$ -differentially private.* \square

Proof. At the generation stage, for each timestamp i , our model generates a token by sampling from Eq. 7, which follows the form of exponential mechanism. This process achieves $(\epsilon + \ln(s))$ -differential privacy as in Definition 4. Every input of the generator is the original input data x_{i-1} (see Eq.2). Eq.7 satisfies the sequential composition theorem. By repeating this process l times, the complete sampling function provides $(\epsilon + \ln(s)) * l$ -differential privacy. \tilde{x}_{dp} is $(\epsilon + \ln(s)) * l$ -differentially private. \square

B Proof of Two-Set Exponential Mechanism.

Theorem 4. Two-Set Exponential Mechanism. *Given a privacy budget $\epsilon > 0$ and the size of output space s , two-set exponential mechanism is $(\epsilon + \ln(s))$ -differentially private.* \square

Proof. Given tokens sets S and O , $\mathcal{V} = O \cup S$, $\emptyset = O \cap S$, and $\mathcal{N} = \{S, O\}$. Let the choice, C , on S and O be ϵ -differentially private, the sampling of an item in the chosen set be totally random. With $Pr[tk, tkN, N|x] = Pr[tkN|x] * Pr[\epsilon_{\epsilon, \rho}(C) = N|tkN, x] * Pr[tk|N, tkN, x]$, where $N \in \mathcal{N}$, $i \in [1, l]$, $tk \in \mathcal{V}$, for any $x' \sim x$:

$$\frac{Pr[\epsilon_{\epsilon, \rho}(\tilde{x}_i) = tk|x]}{Pr[\epsilon_{\epsilon, \rho}(\tilde{x}_i) = tk|x']} = \frac{Pr[tk, tkS, S|x]}{Pr[tk, tkS, S|x'] + Pr[tk, tkO, O|x']} + \frac{Pr[tk, tkO, O|x]}{Pr[tk, tkS, S|x'] + Pr[tk, tkO, O|x']}.$$

For the first part, denoted as P_S , with \mathcal{V} of size s , by dividing the numerator and denominator with $Pr[\epsilon_{\epsilon, \rho}(C) = S|tkS, x] * Pr[tk|S, tkS, x]$, we can get:

$$P_S = \frac{Pr[tk, tkS, S|x]}{Pr[tk, tkS, S|x'] + Pr[tk, tkO, O|x']} \geq \frac{Pr[tkS|x]}{\exp(\epsilon) * s}$$

since

$$\frac{Pr[\epsilon_{\epsilon, \rho}(C) = S|tkS, x']}{Pr[\epsilon_{\epsilon, \rho}(C) = S|tkS, x]} \leq \exp(\epsilon)$$

$$\frac{Pr[\epsilon_{\epsilon, \rho}(C) = O|tkO, x']}{Pr[\epsilon_{\epsilon, \rho}(C) = S|tkS, x]} \leq \exp(\epsilon)$$

$$Pr[tk|S, tkS, x'] / Pr[tk|S, tkS, x] \leq s$$

$$Pr[tk|O, tkO, x'] / Pr[tk|S, tkS, x] \leq s.$$

For the second part, denoted as P_O , similarly, we have $P_O \geq P[tkO|x] / (\exp(\epsilon) * s)$. Then,

$$\frac{Pr[\epsilon_{\epsilon, \rho}(\tilde{x}_i) = tk|x]}{Pr[\epsilon_{\epsilon, \rho}(\tilde{x}_i) = tk|x']} = P_S + P_O \geq \frac{Pr[tkS|x] + Pr[tkO|x]}{\exp(\epsilon) * s}$$

Since $Pr[tkS|x] + Pr[tkO|x] = 1$, the equation can be written as: $Pr[\epsilon_{\epsilon, \rho}(\tilde{x}_i) = tk|x'] / Pr[\epsilon_{\epsilon, \rho}(\tilde{x}_i) = tk|x] \leq \exp(\epsilon) * s = \exp(\epsilon + \ln(s))$. Therefore, the two-set exponential mechanism satisfies $(\epsilon + \ln(s))$ -differential privacy. \square