

# Text Revealer: Private Text Reconstruction via Model Inversion Attacks against Transformers

Ruisi Zhang

University of California San Diego

Seira Hidano

KDDI Research, Inc.

Farinaz Koushanfar

University of California San Diego

## Abstract

Text classification has become widely used in various natural language processing applications like sentiment analysis. Current applications often use large transformer-based language models to classify input texts. However, there is a lack of systematic study on how much private information can be inverted when publishing models. In this paper, we formulate *Text Revealer* — the first model inversion attack for text reconstruction against text classification with transformers. Our attacks faithfully reconstruct private texts included in training data with access to the target model. We leverage an external dataset and GPT-2 to generate the target domain-like fluent text, and then perturb its hidden state optimally with the feedback from the target model. Our extensive experiments demonstrate that our attacks are effective for datasets with different text lengths and can reconstruct private texts with accuracy.

## 1 Introduction

Natural language processing with its application in various fields have attracted much attention in recent years. With the recent advance in transformer-based language models (LMs), BERT (Devlin et al., 2018) and its variants (Liu et al., 2019; Xie et al., 2021; Qin et al., 2022) are used to classify text datasets and achieve state-of-the-art performance. However, LMs tend to memorize data during training, which results in unintentional information leakage (Carlini et al., 2021).

Model inversion attacks (Fredrikson et al., 2015), which invert training samples from the private dataset, has long been applied in the vision domain (Yang et al., 2019; Zhang et al., 2018; Wang et al., 2021; Kahla et al., 2022). For text-based datasets, the model inversion attacks have been applied in the medical domain to infer patients' privacy information. Recent work like KART (Nakamura et al., 2020) and Lehman *et al.* (Lehman et al.,

2021) consider inverting tabular data with sensitive attributes from the medical datasets. However, they follow a fill-in-blank scheme and fail to reconstruct sentences with fluency from scratch.

In this paper, we focus on reconstructing private training data from fine-tuned LMs at inference time. It has a more general scenario where unauthorized personal data such as chats, comments, reviews, and search history may be used to train LMs. We perform a systemic study on how much private information is leaked via model inversion attacks. There are several challenges when performing model inversion attacks on NLP models. First, the candidate pixel range is 256 in images, but the candidate token range is more than 30,000. Therefore, it is harder to find the exact token for the sentence. Secondly, image inversion is more error-tolerant, i.e., error in some pixels will not affect the overall results. However, errors in some of the tokens will significantly affect the fluency of the texts. Thirdly, current text reconstruction attacks have different settings from model inversion attacks. Most of them fall into two categories: Gradient attack (Deng et al., 2021; Zhu et al., 2019), which utilize gradient during distributed training to do attack. Embedding level reconstruct attack (Xie and Hong, 2021; Pan et al., 2020), which trains a mapping function to reconstruct texts from pre-trained embeddings.

We propose *Text Revealer* to perform model inversion attack on text data. In the attack, the adversary knows the domain of the private dataset and has access to the target models. Our attack consists of two stages: in the first stage, we collect texts from the same domain as the public dataset and extract high-frequency phrases from the public dataset as templates. Then, we train a GPT-2 as the text generator on the public dataset. In the second stage, we borrow the idea from PPLM (Dathathri et al., 2019) to continuously perturb the hidden state in the GPT-2 based on the feedback from the

target model. By minimizing the cross-entropy loss, generated text distribution becomes closer to the private datasets. Experiments on the Emotion and Yelp datasets with two target models, BERT and TinyBERT, demonstrate *Text Revealer* can reconstruct private information with readable contents. In summary, our approach has the following contributions:

- We propose *Text Revealer*, the first model inversion attack for text reconstruction against text classification with transformers, to reconstruct private training data from the target models.
- Results on two transformer-based models and two datasets with different lengths have demonstrated *Text Revealer* can reconstruct private texts with accuracy.

## 2 Approach

### 2.1 Threat Model

**Adversary’s Target and Goal** The adversary’s goal is to invert memorized training data from the fine-tuned LMs. The inverted texts should meet two requirements: (1) the texts are close to the distribution of private dataset  $D_{pri}$ . (2) the texts are readable so the adversary can infer meaningful contents from them. In our paper, we use the BERT and TinyBERT as the target LMs to do text classification tasks. We choose the target LMs for the following reasons: (1) BERT and its variants achieve state-of-the-art in most text classification tasks. (2) From the ethical standpoint, evaluating how sensitive the transformer-based model is to text classification, one of the essential tasks in the NLP domain, is important for subsequent research on defense methods.

**Adversary’s Capability and Knowledge** We consider the adversary knows the domain of the dataset on which the language model (LM) is fine-tuned. The adversary also has white-box access to the LM. During the attack, given the input sentences or input embeddings, the adversary can get the prediction score over  $N$  classes with the probabilities  $P = (p_1, p_2, \dots, p_n)$ . Throughout this paper, let  $p_a(X)$  denote the prediction with a given input  $X$  for a label  $a$ .

### 2.2 Attack Construction

As shown in Figure 1, the general attack construction consists of two stages: (1) public dataset collection and analysis, and (2) word embedding perturbation.

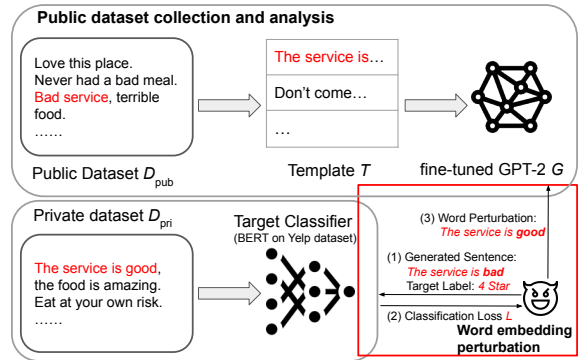


Figure 1: Attack construction pipeline.

**Public dataset collection and analysis** We first collect dataset from the same domain to form public dataset  $D_{pub}$  and perform  $n$ -gram analysis. There is no annotation in  $D_{pub}$ . As the dataset grows larger, the word frequency between private dataset  $D_{pri}$  and public dataset  $D_{pub}$  becomes closer (see Section 3.7 for the details). Therefore, we can infer some phrases may appear in the private dataset by doing the  $n$ -gram analysis and collecting high-frequency items as template  $T$ . Then, we fine-tune a GPT-2  $G$  to generate fluent texts following the distribution of  $D_{pub}$ . We train  $G$  by minimizing  $\min_G L(G, D_{pub})$ , where  $L$  is the cross entropy loss.

**Word embedding perturbation** We borrow the text perturbation idea from Plug and Play Language Model (PPLM) (Dathathri et al., 2019) but change the optimization objective to perform the model inversion attack. PPLM is a lightweight text generation algorithm that uses an attribute classifier to help the GPT-2 perturb its hidden state and guide the generation.

We perturb the hidden state of the text generator to make generated texts closer to the distribution of the private dataset. Let  $H_t$  be the current hidden state of the text generator  $G$ . Let  $L_{adv}$  be an adversarial loss to measure the distance between the generated text  $G(H_t)$  and the private dataset  $D_{pri,a}$  of the target label  $a$ . The adversary generates perturbations  $\Delta H_t$  for  $H_t$  by solving the following optimization problem:

$$\min_{\Delta H_t} L_{adv}(G(H_t + \Delta H_t), D_{pri,a}). \quad (1)$$

Since the adversary has no prior knowledge of the private dataset, we measure the loss  $L_{adv}$  with the cross-entropy calculated from the target model’s prediction  $p_a(G(H_t))$  with the generated

text  $G(H_t)$  for the target label  $a$ . This is because the cross-entropy takes a small value when the input to the model is similar to the training data. We also show that the cross-entropy is more effective for our model inversion attacks than the loss based on the whole prediction score  $P$ , such as the modified entropy (Song and Mittal, 2021) (see Section 3.6 for the details). To obtain the optimal  $\Delta H_t$ , we minimize the cross-entropy by descending the hidden state.

### 3 Experiments

#### 3.1 Datasets

**Emotion Dataset** (Saravia et al., 2018) is a sentence-level emotion classification dataset with six labels: sadness, joy, love, anger, fear, and surprise. **Yelp Dataset** (Zhang et al., 2015) is a document-level review dataset. The reviews are labeled from 1 to 5 stars indicating the user’s preference. Following the split methods in image model-inversion attacks (Zhang et al., 2020), we randomly sample 80% of the samples as public dataset and 20% of the samples as the private dataset. The average token length is 20 in the private Emotion dataset and 134 in the private Yelp dataset.

#### 3.2 Evaluation Metrics

**Recovery Rate (RR.)** is the percentage of tokens in private dataset recovered by different attack methods. We filtered punctuations, special tokens and NLTK’s stop words (Bird et al., 2009) in the private dataset. **Attack Accuracy (Acc.)** is the classification accuracy using evaluation classifier on inverted texts. According to (Zhang et al., 2020), the higher the classification accuracy is, the more private information the texts is considered to be inverted. We use BERT Large as the evaluation classifier and fine-tune it until the accuracy on the private dataset is over 95%. **Fluency** is the Pseudo Log-Likelihood (PLL) of fixed-length models on inverted texts<sup>1</sup>.

#### 3.3 Baselines

We compare our method with vanilla model inversion attack (Fredrikson et al., 2015) and vanilla text generation (Radford et al., 2019).

**Vanilla model-inversion attack (VMI)** In this setting, the adversary exploits the classification loss by adjusting text embeddings and returning the

texts minimizing the cross-entropy loss. We set the text length to the average length of the private dataset, adjust the text embeddings for 50 epochs and run the same times as the number of templates.

**Vanilla text generation (VTG)** In this setting, the adversary uses a GPT-2 model to generate the texts. The GPT-2 is trained on the public dataset and conditioned on our collected templates. We limit the maximum length to the average length of the private dataset. For the attack accuracy on VTG, we calculate the frequency of collected templates under different labels in the public dataset and use the label with the highest frequency as the target label. During the attack, only VTG requires extra annotations to benchmark its performance. For target model TinyBERT and BERT, we run VTG twice and use the results for two target models.

#### 3.4 Target Models

We use two representative transformer-based LM as our target model: (1) Tiny-BERT (Bhargava et al., 2021) and (2) BERT (Devlin et al., 2018). The TinyBERT has four layers, 312 hidden units, a feed-forward filter size of 1200, and 6 attention heads. It has 110M parameters. The BERT has 12 layers, 768 hidden units, a feed-forward filter size of 3072, and 12 attention heads. It has 4M parameters.

#### 3.5 Results

The result of our model inversion attack is summarized in Table 1. We can make the following observations: (1) VTG and *Text Revealer* tend to invert more tokens and achieve lower PLL compared with VMI. This is because VMI trains the text embeddings from scratch with random initialization, which results in not meaningful combinations of tokens. Even though many tokens can be recovered using VMI, private information still cannot be inferred from the private dataset. (2) Compared with VTG, our algorithm achieves higher recovery rates and higher attack accuracy. By perturbing the hidden state of trained GPT-2, our algorithm can infer more private information from the target model. (3) For smaller Emotion dataset, all three methods achieve high attack accuracy and can invert private information from the private dataset. However, in the larger Yelp dataset, VMI and VTG’s attack accuracy becomes near 20%. It is close to random classification because only 5 classes are in the dataset. (4) For both VMI and *Text Re-*

<sup>1</sup><https://huggingface.co/docs/transformers/perplexity>

	Emotion Dataset						Yelp Dataset					
	Tiny-BERT			BERT			Tiny-BERT			BERT		
	VMI	VTG	TR	VMI	VTG	TR	VMI	VTG	TR	VMI	VTG	TR
RR.(%)	27.54	61.23	<b>67.36</b>	34.33	59.42	<b>67.76</b>	50.89	74.47	<b>83.22</b>	75.56	73.22	<b>84.29</b>
Acc. (%)	72.04	74.50	<b>84.14</b>	73.90	73.65	<b>85.54</b>	22.56	23.34	<b>33.98</b>	21.29	22.16	<b>34.22</b>
PLL	1493	<b>43.31</b>	80.79	2628	<b>42.27</b>	94.53	2391.09	<b>16.82</b>	62.75	4010.46	<b>16.64</b>	72.78

Table 1: Results of the model-inversion attack, TR is short for our *Text Revealer*.

*vealer*, BERT’s recovery rate and attack accuracy are higher than TinyBERT. It means more private information is memorized as the transformer becomes larger.

### 3.6 Ablation Studies

**Effectiveness of fine-tuning GPT-2** In this setting, we fine-tune the GPT-2 on the public dataset and compare the results with vanilla GPT-2 on the Yelp dataset. The results are summarized in Table 2. The table shows that fine-tuned GPT-2 achieves higher recovery rate and attack accuracy. It means more sensitive information is revealed from fine-tuning.

	RR.	Acc.	PLL
fine-tuned GPT	84.29	34.22	72.78
vanilla GPT	79.17	31.29	72.21

Table 2: Model inversion with guidance from fine-tuned GPT-2 and vanilla GPT-2 on BERT.

**Effectiveness of model inversion attack** We analysis the effectiveness of word perturbation and loss function by comparing with different methods. For word perturbation, we compare *Text Revealer*’s word perturbation with Gumbel softmax (Jang et al., 2016). For Gumbel softmax, we first use fine-tuned GPT-2 to generate original sentences, and set coefficient for tokens in vocabulary. Then, we update the coefficient based on the cross-entropy loss using Gumbel softmax and update the input sentences. From the table, we can find Gumbel softmax and modified perturbation achieve similar attack accuracy. However, the recovery rate and fluency is lower than modified perturbation.

For loss function, we compare cross entropy with Modified entropy loss (Song and Mittal, 2021). The modified entropy loss makes the loss monotonically decreasing with the prediction probability of the correct label and increasing with the prediction probability of any incorrect label. In this setting, we use the same pipeline as *Text Revealer*, but change the loss to modified entropy loss to update GPT-2’s hidden state. We can find cross entropy

loss achieves best performance out of other loss functions.

Attack	RR.	Acc.	PLL
<i>Text Revealer</i>	84.29	34.22	72.78
Gumbel softmax	78.13	33.49	77.70
Modified Entropy Loss	74.23	32.74	77.25

Table 3: Model inversion with different method.

### 3.7 Analysis

**Effectiveness of template length** In this setting, we analyze how the inversion accuracy changes with the length of the template and summarize the results in Table 7. The template length is  $L$ , and we extract the first  $0.3 L$ ,  $0.5 L$ , and  $0.7 L$  tokens from the template. If the extracted token number is smaller than 1, we choose the first token from the template. For example, if the original template is "if i could give this place 0 star", then  $0.7 L$  template would be "if i could give this,"  $0.5 L$  template would be "if i could give," and  $0.3 L$  template would be "if i." As text length becomes shorter, *Text Revealer* achieves worse recovery rates and attack accuracy. It means longer template length can give GPT-2 more contexts and help it invert more private information.

Template Length	RR.	Acc.	PLL
$0.3 L$	84.17	17.60	76.45
$0.5 L$	84.03	22.49	80.85
$0.7 L$	83.62	27.32	79.71
$L$	84.29	34.22	72.78

Table 4: Model inversion with different length of the template. The percentage means the number of tokens we take since the start token.

**Inverted examples & Visualization** Following TAG (Deng et al., 2021), We display how our inverted examples is approximate the distribution of private dataset at embedding level and sentence level. We use sentences with longest matching subsequence in the private dataset as ground truth. We first use Principal Component Analysis (PCA) (Abdi and Williams, 2010) to reduce the dimension for ground truth and inverted examples

	Ground Truth	Inverted
Example 1	i don t <b>feel</b> so <b>exhausted all the time</b>	<b>feel</b> a little <b>exhausted</b> feel a little bit from <b>all the time</b>
Example 2	<b>i feel blessed amazed and</b> yes very excited	<b>i feel blessed and</b> grateful for the <b>amazing</b> new life
Example 3	<b>the Parlor has</b> one of the <b>best</b> atmospheres <b>in Phoenix!</b>	<b>The Parlor has</b> soo, the <b>best</b> Mexican restaurants <b>in Phoenix</b> Palace!
Example 4	The decor was very inviting and <b>we love Mexican food</b> (again, <b>from San Diego</b> ) <b>The chips and salsa</b> were delicious.	<b>We are from San Diego</b> I was a <b>Mexican food lover</b> , we had <b>the chips, salsa</b> was very nice.

Table 5: Iverted Examples, the first two are from the Emotion dataset, and the last two are from the Yelp dataset.

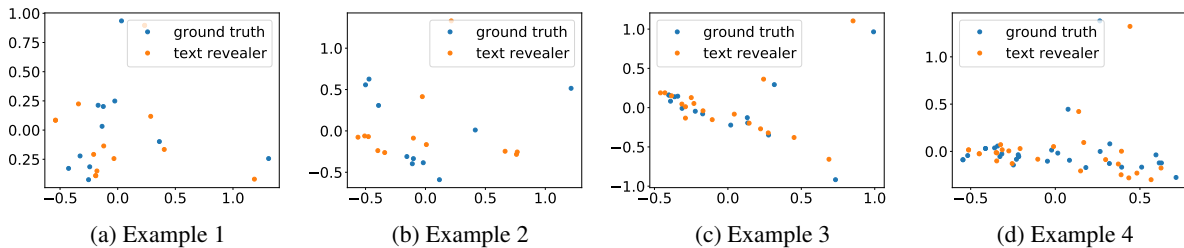


Figure 2: Visualization of inverted embeddings, the first two are from the Emotion dataset, and the last two are from the Yelp dataset.

and display the plot in Figure 2. Then, we display the inverted examples at sentence level in Table 5

**Public dataset & Private dataset** We use Kendall tau distance (Fagin et al., 2003) to measure the token distribution similarity between the private and public datasets. The Kendall tau distance is a metric to measure the top- $k$  elements correlation between two lists. The larger the distance is, the stronger correlation the two lists have. When ranking the tokens in each dataset, we filtered the special tokens, punctuations, and stop words in NLTK (Bird et al., 2009). From the results in Table 6, we can find the word ranking and frequency is close between the private and public datasets.

We also compare the accuracy of the Emotion and the Yelp Dataset in Table 6. Before fine-tuning, the classification accuracy is similar for both private and public datasets. However, after fine-tuning, the accuracy between the private dataset and the public dataset has a large gap, which means some texts are memorized (Carlini et al., 2021) in the trained transformers.

## 4 Conclusion

In this paper, we present *Text Revealer* to invert texts from the private dataset. Experiments on

	Emotion	Yelp Dataset
Distance(Top10)	0.99	0.99
Distance(Top100)	0.32	0.62
Tiny-BERT-Public(+)	84.79	60.58
Tiny-BERT-Private(+)	94.15	79.44
Tiny-BERT-Public(-)	12.81	20.38
Tiny-BERT-Private(-)	13.25	20.59
BERT-Public(+)	91.25	64.42
BERT-Private(+)	99.66	95.12
BERT-Public(-)	13.34	20.02
BERT-Private(-)	14.09	19.88

Table 6: The first cell is the ranking distance between the private and public datasets. The second and third cell is the accuracy on public and private datasets, "+" means the BERT model is finetuned on the private dataset, "-" means the plain BERT model.

different target models and datasets have demonstrated our algorithm can faithfully reconstruct private texts with accuracy. In the future, we are interested in exploring potential defense methods.

## References

Hervé Abdi and Lynne J Williams. 2010. Principal component analysis. *Wiley interdisciplinary re-*

- views: *computational statistics*, 2(4):433–459.
- Prajwal Bhargava, Aleksandr Drozd, and Anna Rogers. 2021. Generalization in nli: Ways (not) to go beyond simple heuristics.
- Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. "O'Reilly Media, Inc."
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. 2021. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2019. Plug and play language models: A simple approach to controlled text generation. *arXiv preprint arXiv:1912.02164*.
- Jieren Deng, Yijue Wang, Ji Li, Chao Shang, Hang Liu, Sanguthevar Rajasekaran, and Caiwen Ding. 2021. Tag: Gradient attack on transformer-based language models. *arXiv preprint arXiv:2103.06819*.
- 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*.
- Ronald Fagin, Ravi Kumar, and D. Sivakumar. 2003. Comparing top k lists. *SIAM Journal on Discrete Mathematics*, 17(1):134–160.
- Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. 2015. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, pages 1322–1333.
- Eric Jang, Shixiang Gu, and Ben Poole. 2016. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2019. Tinybert: Distilling bert for natural language understanding. *arXiv preprint arXiv:1909.10351*.
- Mostafa Kahla, Si Chen, Hoang Anh Just, and Ruoxi Jia. 2022. Label-only model inversion attacks via boundary repulsion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15045–15053.
- Eric Lehman, Sarthak Jain, Karl Pichotta, Yoav Goldberg, and Byron C Wallace. 2021. Does bert pre-trained on clinical notes reveal sensitive data? *arXiv preprint arXiv:2104.07762*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Yuta Nakamura, Shouhei Hanaoka, Yukihiro Nomura, Naoto Hayashi, Osamu Abe, Shuntaro Yada, Shoko Wakamiya, and Eiji Aramaki. 2020. Kart: Privacy leakage framework of language models pre-trained with clinical records. *arXiv preprint arXiv:2101.00036*.
- Xudong Pan, Mi Zhang, Shouling Ji, and Min Yang. 2020. Privacy risks of general-purpose language models. In *2020 IEEE Symposium on Security and Privacy (SP)*, pages 1314–1331. IEEE.
- Haotong Qin, Yifu Ding, Mingyuan Zhang, Qinghua Yan, Aishan Liu, Qingqing Dang, Ziwei Liu, and Xi-anlong Liu. 2022. Bibert: Accurate fully binarized bert. *arXiv preprint arXiv:2203.06390*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Elvis Saravia, Hsien-Chi Toby Liu, Yen-Hao Huang, Junlin Wu, and Yi-Shin Chen. 2018. CARER: Contextualized affect representations for emotion recognition. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
- Liwei Song and Prateek Mittal. 2021. Systematic evaluation of privacy risks of machine learning models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2615–2632.
- Kuan-Chieh Wang, Yan Fu, Ke Li, Ashish Khisti, Richard Zemel, and Alireza Makhzani. 2021. Variational model inversion attacks. *Advances in Neural Information Processing Systems*, 34.
- Keli Xie, Siyuan Lu, Meiqi Wang, and Zhongfeng Wang. 2021. Elbert: Fast albert with confidence-window based early exit. In *ICASSP 2021-IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7713–7717. IEEE.
- Shangyu Xie and Yuan Hong. 2021. Reconstruction attack on instance encoding for language understanding. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2038–2044.
- Ziqi Yang, Ee-Chien Chang, and Zhenkai Liang. 2019. Adversarial neural network inversion via auxiliary knowledge alignment. *arXiv preprint arXiv:1902.08552*.
- Weichen Zhang, Wanli Ouyang, Wen Li, and Dong Xu. 2018. Collaborative and adversarial network for unsupervised domain adaptation. In *Proceedings of*

*the IEEE conference on computer vision and pattern recognition*, pages 3801–3809.

Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28.

Yuheng Zhang, Ruoxi Jia, Hengzhi Pei, Wenxiao Wang, Bo Li, and Dawn Song. 2020. The secret revealer: Generative model-inversion attacks against deep neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Ligeng Zhu, Zhijian Liu, and Song Han. 2019. Deep leakage from gradients. *Advances in Neural Information Processing Systems*, 32.

## A Appendix

### A.1 Limitations

Our work has the following limitations. (1) From the attack side, the model inversion attack achieves lower attack accuracy as the dataset becomes larger. It would be interesting to explore better template extract methods to improve the accuracy. (2) From the ethical standpoint, evaluating how sensitive BERT based model is to classification tasks is essential to future defense methods. Due to page limits, we did not go into the defense methods. We will continue to work on it in future work.

### A.2 Ethics Statement

Our work focuses on the privacy problems in natural language processing. Even though model inversion attacks proposed in our paper can cause potential data leakage. Evaluating how sensitive BERT based model is to classification tasks is important to future defense methods.

### A.3 Hyperparameters and training details

**Template** We perform  $n$ -gram analysis on the public dataset. For Emotion dataset, we extract phrase length from 1 word to 3 words with frequency higher than 20. Then, we filter one-word adjectives and permute them with 2-3 word to form template. The total number is 220. For Yelp dataset, we extract phrase length from 3 words to 8 words with frequency higher than 20. The total number is 320. Some of the templates are summarized in Table 7

**BERT fine-tuning** The BERT and TinyBERT are trained as target models on the private dataset. For both BERT and TinyBERT, we set the batch size to 8. We use AdamW as the optimizer with initial learning rate set to  $5e-5$ . We train 5 epochs on BERT and 10 epochs on TinyBERT. Other parameters follow the original implementation in (Devlin et al., 2018; Jiao et al., 2019).

**GPT-2 fine-tuning** The GPT-2 is trained on the public dataset to help *Text Revealer* reconstruct texts. For GPT-2, we set the batch size to 8, and use AdamW as the optimizer with initial learning rate set to  $5e-5$ . We train 10 epochs for both datasets. Other parameters follow the original implementation in (Radford et al., 2019). It has 124M parameters.

**Text Revealer** When using the *Text Revealer*, we set the iteration number to the average length of the private dataset. We set the window mask to 3 and KL loss coefficient to 0. Other parameters follow the original implementation in (Dathathri et al., 2019).

**Computing Infrastructure** Our code is implemented with PyTorch. Our attack pipeline are all constructed using TITAN Xp. We fine-tune the target LMs and GPT-2 on NVIDIA RTX A6000.

### A.4 More inverted examples

We display more inverted examples in Table 8 and Table 9. We show four examples for each method in Emotion dataset and two examples for each method in Yelp dataset.



Datasets	Examples
Emotion	feeling a little, feel like, down, alone, love, good, passionate, sweet, irritable, angry, strange, nervous, surprised, amazed
Yelp	if i could give this place zero stars, i will never go, really wanted to like this place, i cant wait to go back, if you are looking for, cant say enough good things about, one of the best, i m not a fan of, this place is, if you re in the mood for, only reason i didn t give it

Table 7: Template examples.

Method	Examples
BERT-VMI	cancel the dream vascular place seeing halfway hopeful applicant Mandy. Higgins suspiciously information poollifting Cardinal faint Scots Davidfeld lanes creepy fears Call lap dick Evening Slim favor impressive wireless baptism longing Independent IN forensic Wish physically
BERT-VTG	<b>feeling a little</b> inspired by their own experience. <b>feel like productive</b> people and more than a little creepy. <b>i feel reluctant</b> to give me their phone number. <b>i feel uncomfortable</b> at all in the light of it
BERT-TR	<b>feel like glad</b> to be here i feel welcomed. <b>feel like mellow</b> at breathlessly flowing smoothly and i feel was was pretty <b>i feel afraid</b> that people dont see the movie spoiler place where <b>feel uncomfortable</b> talking to my partner or to stop feeling guilty
TinyBERT-VMI	roughly spiritual thereof 41 rebellious fragrance its caused strangled ave algeria toward puzzled kowalski bobbie sexually liz circle hiding code balls redundant persuaded carpathian yugoslav drownedpins burnley rubbish receipts borough fragmentation worthless soga
TinyBERT-VTG	<b>i feel resentful</b> and lonely, she doesn't seem to like me <b>feel like irritated</b> that my boss is complaining about something <b>feeling a little</b> rude but i wont be afraid <b>feeling a little pain</b> in my lower back and im trying to harden
TinyBERT-TR	<b>feeling a little</b> amazing artistic life positive super job productive <b>feeling mellow</b> wise just like how much she is <b>I feel thankful</b> that god gave me more blessings in return <b>feel like overwhelmed</b> or at the wrong spot saying good feeling

Table 8: Emotion dataset inverted examples, TR is short for our *Text Revealer*.

Method	Examples
BERT-VMI	<p>promising odor Daemon asked exactly platinum wonderful confession Aaron Tiger Magazine field AMC documentaries edible Name Always Benito Northeast invalid Brooks d Games exist range removal thee drag monkey revolution hostile graphicspent permitted expenses mythical newest Massey</p> <p>##wood successors rich hood Jameson Lin publishes festivities Slave ordering trans plot text Programs reluctantly promising Glad España stands combining accommodate captain elongated stranger Conservatory express receives phones bias Depending pile motivation Loving chickens owl denied bage</p>
BERT-VTG	<p><b>i wouldn t go out of my way</b> to come here, but i had a coupon for a 50 dollar pizza for one hour and a free 30 minute massage in addition to the 20 dollars i spent. when we arrived the staff greeted us and escorted us to our room. that was good service and i enjoyed that.</p> <p><b>can t wait to go back</b> and try more! the food was amazing and the service was beyond fast! I had the grilled pork chop which was one of the best dishes I have ever had in my life. the other dishes, the veal, were also very good! The dessert was just as amazing, I would have happily eaten any dessert.</p>
BERT-TR	<p><b>i wouldn t go out of my way</b> to come here, but i love the food and atmosphere well, i`m a foodie, so i know what to expect from this place. the chips are the BEST i`ve had in my lifetime. i had the beef barbacoa taco and i know what to expect from this place, it was absolutely delicious!</p> <p><b>i can t wait to go for</b> me to get my first! This is the most of my favorite place, you. I was shopping a place to take my family that and get the best. the restaurant is always clean and the staff are friendly.</p>
TinyBERT-VMI	<p>##wil apt fish pairs certification maker walker indictment thieves unity negligence 69 crook criminal fraud aligned partly negligence racist bewildered architects inspectors pageant fraud victims atheist restore harassment gasoline ny removed screenplay transformed feathers investigation regular layers random</p> <p>haiti sloping revolving simulator visitor decided marathon faerie sizable sob ty gerald decker buzz shortage suriname verity dream edpace batavia oliver meter41 winced surround fewer grumbled celtics reducing qualify secured completes brandon dimly brandonnb pointless hunts proper attemp</p>
TinyBERT-VTG	<p><b>i can t say enough good things about</b> this place i just will be back and i promise you i will never miss a thing this place is the bomb. and the staff is super nice and helpful post again.the good is amazing. they have a wide variety of beers all great.even on tap. and the service is top notch</p> <p><b>if i could give this place zero stars</b>i would. i have been going here since it opened, but in recent years this place has went downhill. the waiters here use to be very attentive and nice. now i don`t get it. everything they bring you seems to have a strange yippiefaste, which makes my blood boil.</p>
TinyBERT-TR	<p><b>you get what you pay for</b> for! i will always give my money to a place that gives me the best quality I can give. the ambiance is amazing, but the food takes forever. one friend and i went on a tuesday night, and although the place was dead, only a handful of people were eating. order 2 things: the kobe beef shortribs, the shish kabob, and the chicken shawerma.</p> <p><b>can t wait to go back!</b> my husband had the shrimp tacos and they were fantastic! he had the margaritas and they were ok, but, the atmosphere and staff was just really nice and friendly. we`ll definately go back. for the awesome margaritas!!for the delicious food! great prices!!!to our server, jason. for the amazing service!</p>

Table 9: Yelp dataset inverted examples, TR is short for our *Text Revealer*.