

# Evaluating Differentially Private Synthetic Data Generation in High-Stakes Domains

**Krithika Ramesh**

Johns Hopkins University  
kramesh3@jh.edu

**Nupoor Gandhi**

Carnegie Mellon University  
nmgandhi@cs.cmu.edu

**Pulkit Madaan**

Johns Hopkins University  
pmadaan2@jhu.edu

**Lisa Bauer**

Amazon  
bauerlb@amazon.com

**Charith Peris**

Amazon  
perisc@amazon.com

**Anjalie Field**

Johns Hopkins University  
anjalief@jhu.edu

## Abstract

The difficulty of anonymizing text data hinders the development and deployment of NLP in high-stakes domains that involve private data, such as healthcare and social services. Poorly anonymized sensitive data cannot be easily shared with annotators or external researchers, nor can it be used to train public models. In this work, we explore the feasibility of using synthetic data generated from differentially private language models in place of real data to facilitate the development of NLP in these domains without compromising privacy. In contrast to prior work, we generate synthetic data for real high-stakes domains, and we propose and conduct use-inspired evaluations to assess data quality. Our results show that prior simplistic evaluations have failed to highlight utility, privacy, and fairness issues in the synthetic data. Overall, our work underscores the need for further improvements to synthetic data generation for it to be a viable way to enable privacy-preserving data sharing.

## 1 Introduction

The widespread availability of public digitized text has greatly facilitated the advancement of natural language processing (NLP). Text processing could also be extremely valuable for processing high-stakes private data, like healthcare records (Panchbhai and Pathak, 2022), social workers’ notes (Gandhi et al., 2023), or legal documents (Zhong et al., 2020). However, the need to maintain data privacy hinders the responsible development and deployment of models in these domains.

Building NLP tools often requires sharing data externally with contractors or researchers, as agencies like child protective services typically do not have in-house AI expertise. While data sharing has been accomplished through data use agreements with individual teams or laboriously redacting identifiable information from text (e.g., Johnson et al.

(2016a)), these approaches have limitations. Limited sharing still requires increasing the number of people who have access to sensitive data, and it precludes the development of public benchmarks, which have proved crucial for standardizing model development. Redaction fails to fully prevent re-identification, as even lower dimensional data is often possible to re-identify with just small amounts of auxiliary data (Narayanan and Shmatikov, 2008; Sweeney, 2000). Furthermore, redacted data is not useful for tasks requiring sensitive information, such as developing a model to identify contact information for potential caretakers of a child (Field et al., 2023).

In our work, we propose and conduct use-inspired evaluations of the feasibility of using synthetic data to address these limitations. Recent work has proposed sharing synthetic text generated from differentially private language models in place of real data (Yue et al., 2023; Kurakin et al., 2023; Mattern et al., 2022a; Putta et al., 2023). Differential privacy (DP) offers an appealing solution, as it provides a theoretical guarantee of privacy preservation that is controllable through a specified privacy budget. Although the bulk of work in developing DP approaches has been centered around models trained on tabular and image-related data, there has been increasing interest in applying DP to unstructured text data (Shi et al., 2022; Yue et al., 2021; Feyisetan et al., 2020a).

While initial results of synthetic data are promising, prior work has lacked grounding in realistic applications, for example, running experiments with public internet data that language models may already have been exposed to during pre-training.

In contrast, we conduct experiments on text data from two high stakes domains: healthcare and child protective services, and we rigorously evaluate the synthesized text for its utility, privacy, and potential fairness implications. For utility and privacy, we introduce novel well-motivated evaluation criteria

(“silver” coreference modeling and entity-centric metrics). To the best of our knowledge, no prior work has investigated fairness considerations in this domain.

We evaluate several approaches for privacy-preserving synthetic data generation, including fine-tuned models and an in-context-learning approach. While these evaluations reveal some promising opportunities for synthetic text, they further expose utility degradation, privacy leakage (even when using DP), and issues with group fairness. These results indicate that prior simplistic evaluations have overestimated current viability of synthetic data.

Our primary contributions include a rigorous and reproducible evaluation framework that exposes limitations underestimated in prior work, and empirical results over real high-stakes data. Overall, our work demonstrates that contrived metrics do not necessarily translate to more realistic scenarios, emphasizing the need for thorough in-domain evaluation to understand potential strengths and limitations of synthetic data.

## 2 Methodology

### 2.1 Text Generation

The primary goal of privacy-preserving synthetic text generation is to generate realistic, but entirely synthesized text for a high stakes domain, such as doctors’ notes from a healthcare institution. We assume we have a data set of real text from that domain, which we can use to guide the generation. In addition to being realistic, it needs to be ensured that the synthetic data does not reveal uniquely identifiable information about any individuals from the original data.

**Fine-tuning** We adapt the dominant approach from prior work (Yue et al., 2023): starting with a pre-trained autoregressive language model, fine-tune it using the real in-domain data, and then generate new data from it. We compare fine-tuning the model with and without DP, where we use DP-SGD for differentially private fine-tuning. For reference, we provide background on DP and DP-SGD in Appendix A. After fine-tuning, we utilize top-k sampling (Fan et al., 2018) and nucleus sampling (Holtzman et al., 2020) to generate diverse synthetic notes.

We condition the text generation on *control codes* (Keskar et al., 2019). During training, we prepend one or more labels associated with the text to the model input. We similarly prepend

control codes during inference, where we sample the provided codes from their distribution in the training data. Thus, during training and inference, the probability distribution of the subsequent text  $x = \{x_1, x_2 \dots x_n\}$  is conditioned on the control code information  $c$ , which is described by the following equation:

$$P(x|c) = \prod_{i=1}^n P(x_i|x_1 \dots x_{i-1}, c) \quad (1)$$

Controllable generation approaches enable the generation of notes with specific properties. We primarily use them to enable classification-based utility evaluations (described in §2.2).

**ICL** In order to explore the potential capabilities of much larger models and investigate if fine-tuning is actually needed, we also generate notes using in-context learning (ICL). We provide as context examples of training data text with prepended control codes, followed by an additional set of codes to prompt the model to generate content in accordance with the final set of codes. The number of examples provided varies, as we require that each control code for the note to be generated is associated with at least one in-context example. This approach is most directly comparable to the fine-tuned models without DP.

### 2.2 Utility Evaluation

Given the goal of developing synthetic data that could be shared externally with researchers or third-party contractors, we evaluate the data’s utility based on the performance of NLP models trained over this data. More specifically, we train NLP models on the synthetic data and evaluate their performance over real data.

**Classification** Similar to prior work (Yue et al., 2023; Kurakin et al., 2023), we evaluate model performance over classification tasks, where we use the control codes provided during text generation as class labels. We focus on multiclass and/or multilabel classification tasks, and we compare model performance as task difficulty increases.

**Coreference Resolution** Classification tasks can be highly dependent on keywords and phrases, and they do not necessarily require training data to be coherent and consistent across a full paragraph or document. Consistency of entity properties across a document, however, is a necessary condition for

coreference training data. Coreference and the related task of mention detection also offer a realistic use case in processing expert-written notes (Gandhi et al., 2023). Thus, we measure the utility of the synthetic data for training coreference models.

Unlike classification labels, coreference annotations cannot be easily generated through control codes. In a practical setting, annotations of coreference clusters would likely be conducted over synthesized data manually by hired annotators or researchers, but this process does not scale for evaluating of multiple iterations of synthetic data generators. Instead, we use a fine-tuned coreference model to simulate “silver” annotations over the synthesized data.

More specifically, given a subset of the original dataset  $D$  annotated with gold coreference clusters, we first fine-tune a pretrained coreference model (Kirstain et al., 2021) on this data. Using this model, we infer coreference clusters over synthetic data from the same domain which we consider silver annotations. We fine-tune a separate coreference model that has not been task-finetuned with the silver coreference clusters to approximate the utility of the synthetic data for coreference resolution.

We run all experiments with a neural coreference model (Kirstain et al., 2021). We report results after fine-tuning the model for 40 epochs, where scores are averaged over standard coreference metrics: MUC,  $B^3$ ,  $CEAF_{\phi_4}$ .

## 2.3 Privacy Evaluation

**Canary Attacks** Consistent with prior work, to assess the potential leakage of sensitive information in our training data and the extent to which the model memorizes personally identifiable information (PII), we use the canary evaluation method proposed by (Carlini et al., 2019). This approach involves injecting artificial canary sequences containing PII into the training data and analyzing the likelihood of the frequency of appearance of this PII in the generated outputs.

We create artificial canary samples that are contextually relevant to both domains and include PII such as names, emails, addresses, and numeric identifiers (details in the appendix in Table 10 and Table 9). Following the methodology outlined in (Yue et al., 2023), we vary the number of injections of these canary samples into our training data for 1, 10, and 100 repetitions. For each canary, we generate 10,000 candidate sequences and rank the

canaries based on their perplexity score.

**Entity-focused metrics** As canary evaluations are only a proxy for assessing potential privacy risks and may not be comprehensive, we directly leverage entity markers in our datasets to evaluate privacy concerns (we provide details on data-specific entity definitions in §3).

We compare the frequency of identified entities in the original vs. synthetic data. Further, while an isolated entity poses some privacy risk, the risk is magnified if the context surrounding the entity is also leaked. Thus we examine the frequency of entities with variable-length surrounding context in the synthetic data and compare them with the training data to estimate the number of memorized patterns that reappear in the synthetic data.

## 2.4 Fairness Evaluation

We compute fairness metrics over the same control-code classification tasks as the utility evaluation (§2.2). In data with available demographic information, we compare fairness classification for race and gender subgroups using equality difference (ED) and equalized odds (EO) metrics. For ED, for instance, False Positive Equality Difference (FPED) is the sum of the differences between the overall false positive rate (FPR) for the entire dataset and the FPR for each subgroup. EO constitutes a stricter notion of fairness by evaluating whether both the FPR and TPR rates are the same across all groups. In both cases, values closer to zero indicate that the model performs more uniformly across subgroups, with zero indicating perfect parity across subgroups. For reference, we formally define these metrics in Appendix C.

# 3 Experimental Setup

## 3.1 Data

**Healthcare** Our primary source of healthcare data is the MIMIC-III Clinical Database (Johnson et al., 2016b,a; Goldberger et al., 2000), which contains  $> 2M$  deidentified notes associated with  $> 40K$  patients admitted to the Beth Israel Deaconess Medical Center in Boston, Massachusetts.

As control codes we use ICD-9 codes, which are a standardized format for medical conditions that have been human-annotated in MIMIC. Each note can contain multiple possible codes, making our evaluation task multiclass and multilabel. There are  $> 5000$  unique ICD-9 codes. Thus, we restrict data

Training Data	Dataset	F1 Micro	F1 Macro	Subset Accuracy
$D_{real}$	ICD-9 $_{n=3}$	0.89 ± 0.0	0.90 ± 0.0	0.76 ± 0.002
$D_{\epsilon=\infty}$	ICD-9 $_{n=3}$	0.87 ± 0.003 ↓-0.02	0.87 ± 0.005 ↓-0.03	0.74 ± 0.004 ↓-0.02
$D_{\epsilon=8}$	ICD-9 $_{n=3}$	0.85 ± 0.002 ↓-0.04	0.85 ± 0.003 ↓-0.05	0.71 ± 0.003 ↓-0.05
$D_{real}$	ICD-9 $_{n=5}$	0.77 ± 0.008	0.75 ± 0.016	0.56 ± 0.007
$D_{\epsilon=\infty}$	ICD-9 $_{n=5}$	0.75 ± 0.003 ↓-0.02	0.73 ± 0.003 ↓-0.02	0.55 ± 0.004 ↓-0.01
$D_{\epsilon=8}$	ICD-9 $_{n=5}$	0.68 ± 0.004 ↓-0.09	0.60 ± 0.008 ↓-0.15	0.48 ± 0.003 ↓-0.08
$D_{real}$	ICD-9 $_{n=10}$	0.70 ± 0.010	0.67 ± 0.012	0.32 ± 0.016
$D_{\epsilon=\infty}$	ICD-9 $_{n=10}$	0.66 ± 0.001 ↓-0.04	0.61 ± 0.003 ↓-0.06	0.26 ± 0.004 ↓-0.06
$D_{\epsilon=8}$	ICD-9 $_{n=10}$	0.54 ± 0.007 ↓-0.16	0.40 ± 0.004 ↓-0.27	0.18 ± 0.005 ↓-0.14
$D_{ICL}$	ICD-9 $_{n=10}$	0.57 ± 0.011 ↓-0.13	0.47 ± 0.014 ↓-0.20	0.21 ± 0.008 ↓-0.11

Table 1: Difference in performance between models trained on the synthetic data generated with ( $D_{\epsilon=8}$ ) and without ( $D_{\epsilon=\infty}$ ) DP and the models trained on real data ( $D_{real}$ ) for multilabel ICD code classification with the top 3, 5, and 10 most frequent labels. Performance degradation greatly increases for more complex tasks.

to notes containing any of the  $n$  most frequent ICD-9 codes, where we typically set  $n = 10$  and report  $n \in 3, 5$  for some comparisons, similar to Al Aziz et al. (2021); Huang et al. (2019). As a result, the fine-tuning data size for the generative models can vary depending on the value of  $n$ . The dataset splits for the classification tasks are provided in Appendix D. To ensure synthetic data is balanced comparably to real data when evaluating fairness, we additionally provide the patient’s ethnicity and biological sex as control codes.

For coreference resolution, we use notes from the MIMIC-II Database annotated for coreference as a part of the i2b2/VA Shared-Task and Workshop in 2011 (Uzuner et al., 2012). This data includes 251 train documents, 51 of which we have randomly selected for development, and 173 test documents.

As the MIMIC data is already deidentified, we directly leverage the strings used for deidentification, e.g. `[**Hospital1 18**]`, `[**First Name3 (LF) 2704**]`, in order to conduct entity-centric privacy evaluations. Finally, we note that although the MIMIC-III diagnoses notes are not permissible to be used for training publicly available language models, there remains a possibility that some MIMIC notes may have been indirectly included in the training data through various other sources.

**Child Protective Services (CPS)** We additionally report results over a data set of contact notes from a county-level Department of Human Ser-

vices (DHS). These notes log contact with families involved in child protective services, and they are written by caseworkers and other service providers. Unlike MIMIC-III, this data set is not deidentified, which makes it a more realistic test data set, but also prevents the data from being publicly accessible. Throughout our work, this data was stored on a secure server with restricted access, in accordance with IRB-approved protocol and a data sharing agreement established with the county.

The full data set contains 3.1M notes, from approximately 2010 to November 23, 2020. As control codes, we use existing metadata, specifically, the “Contact Source Description” field, which specifies one of five possible labels for each note (*Case, Investigation, Transportation Contact, Provider and Call Screen*). For coreference resolution, we use a set of 200 notes annotated for coreference by prior work and shared with us by the county (Gandhi et al., 2023). This data has train/dev/test sets of sizes 100/10/90 notes. Finally, for entity-centric evaluations, we use a spaCy NER model to identify spans of entities in the text, and we focus on entities likely to contain private identifying information (e.g., names and organizations).

As CPS cases are complex and involve multiple people, the notion of race or gender for a note is less clear than in the MIMIC data. Thus, we do not report fairness results for this data. We also do not report ICL results, as our single secure server did not have sufficient resources for the larger model.

### 3.2 Models

Our primary text generation model is Sheared-LLaMA-1.3B (Xia et al., 2024).<sup>1</sup> We fine-tune using Low-Rank Adaption (LoRA) (Hu et al., 2022), and we use Opacus (Yousefpour et al., 2022) for DP fine-tuning. We generally set a privacy budget of  $\epsilon = 8$ , and  $\delta = 1e-5$  (considering our relatively small dataset size), and we report some results with  $\epsilon = 4$  for comparison. For ICL, we used the instruction-tuned BioMistral-7B<sup>2</sup> model. As the inference for the BioMistral 7B model is compute-intensive, we report results over experiments conducted only on the ICD-9 <sub>$n=10$</sub>  subset of the MIMIC-III healthcare dataset, and we generate a smaller number of notes (0.6 as many) as compared to data generated from the smaller fine-tuned models. We have specified the hyperparameters for each of the models used, dataset distributions and additional detail regarding the experimental setup in Appendix B and Appendix D.

	F1 Score	Accuracy
$D_{real}$	$0.86 \pm 0.003$	$0.86 \pm 0.003$
$D_{\epsilon=\infty}$	$0.80 \pm 0.006$ $\downarrow_{-0.06}$	$0.80 \pm 0.006$ $\downarrow_{-0.06}$
$D_{\epsilon=8}$	$0.69 \pm 0.002$ $\downarrow_{-0.17}$	$0.68 \pm 0.002$ $\downarrow_{-0.18}$
$D_{\epsilon=4}$	$0.65 \pm 0.002$ $\downarrow_{-0.21}$	$0.65 \pm 0.002$ $\downarrow_{-0.21}$

Table 2: Difference in performance between models trained on data generated with differential privacy and models trained on real data, evaluated over CPS classification, for varying privacy budgets.

## 4 Results

### 4.1 Utility

**Overall Classification** Tables 1 and 2 report results for classification tasks for all models, for the healthcare and CPS data respectively. Unsurprisingly, models trained on data generated from DP fine-tuned models generally under-perform models trained on real data or data generated without DP. Table 1 reports performance for varying task complexity by increasing number of labels  $n$  for our multilabel ICD-9 code classification task. For simpler tasks, e.g. ICD-9 <sub>$n=3$</sub> , there is a much smaller performance degradation and the  $D_{\epsilon=\infty}$  ( $F1 \approx 0.87$ ) and  $D_{\epsilon=8}$  ( $F1 \approx 0.85$ ) models are nearly comparable. In contrast, there is much larger performance

<sup>1</sup><https://huggingface.co/princeton-nlp/Sheared-LLaMA-1.3B>

<sup>2</sup><https://huggingface.co/BioMistral/BioMistral-7B>

degradation for the more difficult ICD-9 <sub>$n=10$</sub>  task, where  $F1 \approx 0.61$  for  $D_{\epsilon=\infty}$  and  $F1 \approx 0.40$  for  $D_{\epsilon=8}$ .

In the classification task with the CPS data (Table 2), however, we notice a significant drop in performance for models trained over data generated with DP for both more generous ( $D_{\epsilon=8}$ ) and more restricted ( $D_{\epsilon=4}$ ) privacy budgets. From examining the data, this task is generally more difficult and the associations between the administrative label and the text in the real data can be quite subtle. It is likely that the generative model often fails to pick up on these associations, and noise introduced by DP further masks these subtleties.

**Overall Coreference** Table 3 reports coreference results. For comparison, we report  $D_{real(gold)}$ , model performance when trained over gold in-domain data, which represents the best possible performance we can obtain with human annotations and  $D_{real(silver)}$ , model performance when trained over silver annotated real data. The 15 point performance difference in F1 between these two setups represents the performance degradation we should expect to see as a result of inevitable cascading errors from the silver annotations.

There notable performance degradation in synthetic data generated both with and without DP as compared to real data, which is much more noticeable in coreference metrics than mention detection metrics. For example, for the healthcare data from  $D_{real(silver)}$  to  $D_{\epsilon=8}$  mention detection F1 declines by 0.044 (0.659 to 0.615), whereas coreference F1 declines by 0.109 (0.552 to 0.443). Models trained on data generated with and without DP perform similarly, likely performance decline is dominated by general quality of synthetic data more so than the application of privacy preservation. While both data sets show similar trends, the performance degradation between real and synthetic data is generally worse for MIMIC-III than CPS.

Data generated by the ICL model resulted in higher-performing coreference systems than data generated by the fine-tuned models. It is likely that the larger model outputted generally more coherent data, through the lack of fine-tuning reduces controllability of generation, e.g., as evidence by the lower performance of the ICL model in Table 1.

### 4.2 Privacy

**Canary Attacks** Table 4 reports results for canary attacks. The DP fine-tuned models exhibit

Training Data	Healthcare		CPS	
	Mention Detection	Coreference	Mention Detection	Coreference
$D_{real(gold)}$	$0.799 \pm 0.013$	$0.703 \pm 0.011$	$0.877 \pm 0.004$	$0.789 \pm .005$
$D_{real(silver)}$	$0.659 \pm 0.121$	$0.552 \pm 0.126$	$0.805 \pm 0.007$	$0.642 \pm 0.008$
$D_{\epsilon=\infty}$	$0.615 \pm 0.051$	$0.443 \pm 0.062$	$0.753 \pm 0.030$	$0.570 \pm 0.035$
$D_{\epsilon=8}$	$0.599 \pm 0.043$	$0.438 \pm 0.025$	$0.776 \pm 0.007$	$0.589 \pm 0.009$
$D_{ICL}$	$0.712 \pm 0.010$	$0.588 \pm 0.022$	-	-

Table 3: F1 scores for coreference and mention detection over entities from human-annotated test splits of the CPS and i2b2/VA datasets. All synthetic datasets are annotated with silver labels. There is general performance degradation for synthetic data generated with ( $\epsilon = 8$ ) and without DP ( $\epsilon = \infty$ ) as compared to real data. The performance degradation is more noticeable in coreference metrics than mention detection metrics.

		Rank	Perplexity
Healthcare	Name	5986 / 3378	49.72 / 54.10
	Address	2276 / 4075	43.59 / 62.66
	Number	902 / 841	9.43 / 14.61
	Email	711 / 1452	37.81 / 72.08
CPS	Name	3168 / 2306	10.62 / 10.33
	Address	9618 / 9523	21.63 / 27.23
	Number	474 / 1347	16.81 / 23.24
	Email	387 / 5838	49.91 / 81.61

Table 4: Rank and perplexity metrics for 10-insertion canary attacks over MIMIC and CPS data (0, 1 and 100 insertions, reported in Appendix E, are similar). Each column is formatted as  $\epsilon = \infty / \epsilon = 8$ . Higher rank and lower perplexity typically indicate decreased risk of leakage. DP reduces but does not eliminate privacy risks for all canaries, and metrics are generally unstable.

higher perplexity scores for all the canaries, demonstrating that models trained with DP are less likely to output phrases from training data. There is also a relatively sharper drop in perplexity for models fine-tuned without DP as the number of canary insertions increase (Appendix E). However, our entity-centric evaluation demonstrates that canary evaluations may not be effective in assessing a model’s potential for privacy leakage. It should also be noted that while DP improves (increases) rank for some canaries, it decreases rank for others. The canary’s rank is also highly dependent on the choice of candidate comparisons, making these metrics easy to skew.

We perform analysis of actual leakage (e.g., appearance in generated text) using PII already present in the training data rather than inserted canaries. These entity-centric metrics (Figure 2) show

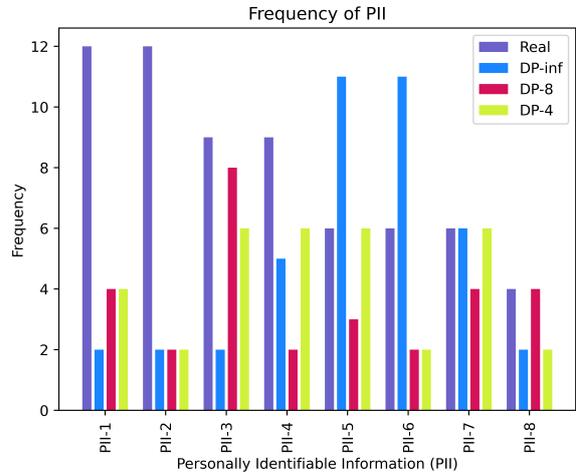


Figure 1: Frequency of the same name (PII-1...PII-8) in real and synthetic CPS data for 8 hand-identified names. The frequency of each name generally, but not always, decreases in synthetic data generated with differential privacy.

that while DP-generated data does contain fewer instances of potentially sensitive information, these entities are not removed from the data entirely, and there is still risk of leakage.

Figure 3 shows the reduction in private entity leakage in data generated from DP fine-tuned models compared to non-DP fine-tuned models. While notably reduced, some leakage still persists when using DP, even on decreasing the privacy budget further. While Figure 2 compares rates of leakage of aggregated across all entities, it does not provide insight into how leakage for individual entities may occur. In Figure 1 we conduct this analysis by manually selecting 8 names that occur in the real and synthetic CPS data and plotting the frequency of each name in each data type. For most names, frequency in the DP-generated synthetic data is less

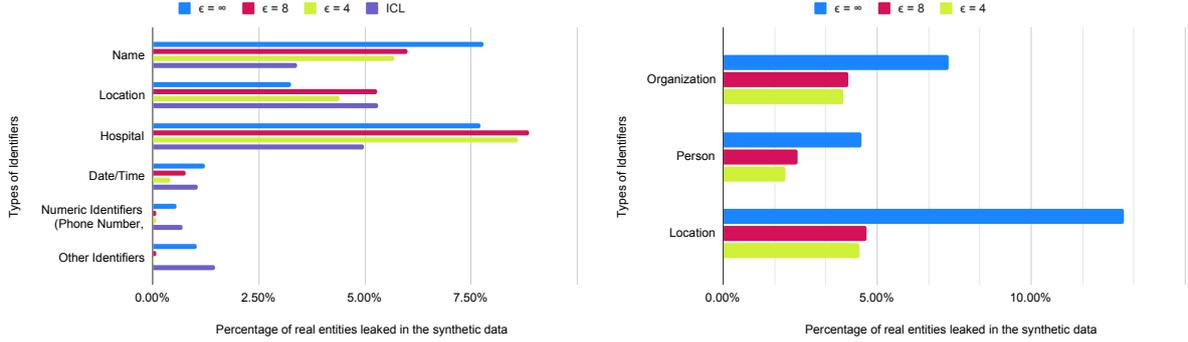


Figure 2: Entity-centric privacy evaluation for MIMIC-III (left) and CPS (right). We report the percent of entities in the real data that are present in the synthetic data. While DP reduces leakage, it does not eliminate it entirely for all entities, even with a more restrictive privacy budget.

than in the original data, but this reduction does not always hold, even with a lower privacy budget (e.g., PII-7, PII-5 for  $\epsilon = 4$ ). In data generated without differential privacy, the frequency of names sometimes exceeds their frequency in the original data (e.g., PII-5, PII-6).

Leaked identifiers are potentially more harmful if additional information about an individual is leaked alongside their identity. We assess this risk in Table 5, where we gauge how often *sequences* of length 1-4 containing these leaked entities appear in the generated outputs, rather than examining entities in isolation.

The results provide further evidence that, while training models with differential privacy may decrease the risk of information memorization, it does not provide a failsafe. There is a notable disparity in the frequency of phrases from the training data reproduced in these datasets:  $D_{\epsilon=\infty}$  contains nearly 1.6 times as many phrases as the  $D_{\epsilon=8}$ , but the phrase leakage from  $D_{\epsilon=8}$  is still non-zero. On the other hand, while  $D_{ICL}$  is 0.6 times the size of the other datasets, it seems to regurgitate contextual information about these entities from the in-context samples less frequently. However, results from Figure 2 indicate that it still poses privacy risks, as the ICL tends to reproduce these entities, even if not the contexts in which they appear.

### 4.3 Fairness

We report the FNED and Equalized Odds (EO) metrics for the results from the ICD-9<sub>n=10</sub> multilabel classification tasks in Table 6. The metrics reflect the difference in model performance for the gender and race/ethnicity subgroups with more than 100 samples in the test set, with a larger value indicating

	Healthcare	CPS
	Count	Count
$D_{\epsilon=\infty}$	16271	4970
$D_{ICL}$	3761	-
$D_{\epsilon=8}$	10312	1555
$D_{\epsilon=4}$	8934	1434

Table 5: Unique contexts in which entities in the real data appear in the synthetic data. Surrounding context word lengths vary from 1 to 4.

more disparate performance across the subgroups. While the gender metrics indicate minimal performance differences, the race/ethnicity metrics show significant disparities. The disparate performance increases for models trained over the data generated from the DP model ( $D_{\epsilon=8}$ ) as compared to the model without DP ( $D_{\epsilon=\infty}$ ). The model trained with DP ( $D_{\epsilon=8}$ ) exhibits the most disparate performance across these subgroups, followed by the  $D_{ICL}$ , although the latter provides better utility for the classification and coreference tasks.

## 5 Discussion

Overall, our results are consistent with prior work in that we find only small performance degradation when training a model on DP-generated synthetic text as compared to real data for relatively less fine-grained (e.g. ICD-9<sub>n=3</sub>, in Table 1) classification tasks. Similarly, we do find evidence that DP reduces potential privacy leakage in that artificial canaries (Table 4) and real entities (Figure 2) are generated less frequently by DP-fine-tuned models.

However, our evaluations also expose previously unexplored weaknesses to this approach. Model

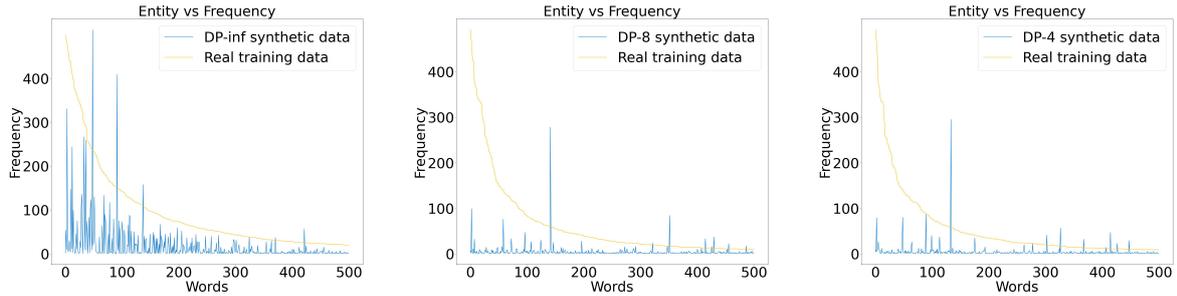


Figure 3: MIMIC-III ICD-9<sub>n=10</sub> data: Graph depicts the frequency of overlapping entities between the training data  $D_{train}$  for the generative model and synthetic data. The top row presents the top 500 most frequent entities from each dataset, limited to entities with a frequency count below 500 in  $D_{train}$ .

		<b>FNED</b>	<b>Equalized Odds</b>
<b>Race</b>	$D_{real}$	$0.35 \pm 0.0$	$0.20 \pm 0.0$
	$D_{\epsilon=\infty}$	$0.39 \pm 0.005$	$0.23 \pm 0.001$
	$D_{\epsilon=8}$	$0.53 \pm 0.014$	$0.30 \pm 0.003$
	$D_{ICL}$	$0.49 \pm 0.03$	$0.29 \pm 0.012$
<b>Gender</b>	$D_{real}$	$0.04 \pm 0.0$	$0.04 \pm 0.0$
	$D_{\epsilon=\infty}$	$0.02 \pm 0.007$	$0.02 \pm 0.007$
	$D_{\epsilon=8}$	$0.04 \pm 0.005$	$0.04 \pm 0.005$
	$D_{ICL}$	$0.04 \pm 0.006$	$0.04 \pm 0.006$

Table 6: Fairness evaluation for the MIMIC-III ICD-9<sub>n=10</sub> task, for the gender and race categories. Higher values indicate poorer group fairness performance. We report additional fairness metrics in Appendix C in Table 7 that show similar trends.

performance degrades much more sharply as task complexity increases (e.g. ICD-9<sub>n=10</sub> classification in Table 1, mention vs. coreference performance in Table 3). These results suggest that DP-generated synthetic data may be of sufficient quality for certain NLP tasks and domains, but the quality degradation from DP is a limitation on broader use.

Post-hoc data filtering and re-ranking may offer a way to improve quality. For example, NLI-based approaches have previously been used to rank or evaluate the quality of the generated text (Dušek and Kasner, 2020; Garneau and Lamontagne, 2021; Chen and Eger, 2023) and have been incorporated into the generation pipeline to enhance the consistency of outputs produced by LMs (Mersinias and Mahowald, 2023), though our initial results with this approach were inconsistent across data sets.

Furthermore, despite claims that differentially private training of language models can effectively

eliminate the risk of privacy leakage (Yue et al., 2023; Mattern et al., 2022a), our experiments indicate that there is also a substantial risk of data leakage (Tables 4-5, Figure 2), especially for some types of PII. These results are consistent with risks of leakage identified in sentence-level applications of differential privacy (Lukas et al., 2023).

On investigating the privacy leakage further, we identify several possible causes. Even for sensitive spans that appear infrequently in the training data, their sub-tokens can recur throughout the same document and across multiple documents more frequently. For instance, in the MIMIC dataset, a token like `[**Hospital1 18**]` might have its "hospital" component repeat multiple times in the data, while the numerical identifier may appear frequently in other contexts, allowing the model to learn all components of the full sensitive span, despite DP-fine-tuning. This pattern can similarly occur for real identifiers, such as when individuals share the same first name or last name in the CPS data. Additionally, the presence of sensitive tokens in the pretraining data and the contextual dependencies in text generation may contribute to the model’s memorization of sequences in the fine-tuning data. Finally, correctly defining the unit of privacy presents a significant obstacle in text settings (Chua et al., 2024). Ensuring privacy at the user-level is naturally greater than that for a single record, potentially requiring additional utility loss or greater computational costs (Charles et al., 2024). Combining privacy-preserving techniques may be a more promising approach than relying on DP.

We further find substantial variance not only in the task difficulty, but also across data sets. Coreference performance degradation from real to synthetic data is markedly worse for MIMIC than

CPS (Table 3). These differences could be due to a number of factors, such as the similarity between each private data set and the model pre-training data. Regardless, these results emphasize the importance of evaluating on in-domain data, as results may not generalize.

## 6 Related Work

The majority of research on enabling shareable sensitive data has focused on text anonymization, replacing or redacting private information like names and addresses from text. While some approaches redact and replace sensitive information using deterministic rule-based systems (Mamede et al., 2016; Yermilov et al., 2023; Ben Cheikh Larbi et al., 2023; Sotolář et al., 2021; Volodina et al., 2020), others employ masked language models (Yermilov et al., 2023). Differentially private mechanisms have also been integrated into text sanitization processes, such as differentially private perturbation of text embeddings (Feyisetan et al., 2020b) or sampling of replacement tokens (Yue et al., 2021; Chen et al., 2023) building on the principle of Metric-Local DP (Alvim et al., 2018). Although these methods are computationally inexpensive and domain-agnostic, they have weak privacy guarantees and limited capacity to modify text (Mattern et al., 2022b; Domingo-Ferrer et al., 2021; Brown et al., 2022).

Recently, datasets comprised entirely of synthetic data have become potentially viable (Guan et al., 2018; Yale et al., 2020). Our work differs from similar approaches to synthetic data generation in its focus on actual high stakes data and thorough grounded evaluation (Yue et al., 2023; Kurakin et al., 2023; Mattern et al., 2022a; Putta et al., 2023). Notably, Al Aziz et al. (2021) do similarly investigate healthcare data, but they do not evaluate potential privacy leakage, and their utility measures do not adequately capture errors in text fluency and consistency which is crucial for finer-grained applications.

A separate but overlapping line of work has focused on improving privacy in NLP models & protection against membership inference, rather than in the generated data. This work has similarly trained NLP models with differential privacy but has evaluated direct performance of these models on downstream tasks (Li et al., 2021; Wu et al., 2022). Nevertheless, this line of work is not directly comparable to ours, as it focuses on training models

directly on private data, while our approach promotes the shareability of data and imposes fewer restrictions on sharing models trained using private synthetic data.

## 7 Conclusions

Although synthetic data generated with differential privacy is an appealing way to improve responsible AI development, off-the-shelf DP does not achieve sufficient privacy, utility, or fairness over real high-stakes data. These failings suggest numerous opportunities for future work on improving the coherence of synthetic data and the application of privacy preservation to this task. Our evaluation methods offer a way to foster this research, with grounding in real applications, rather than contrived settings, where performance is liable to being over-estimated.

## Acknowledgments

This work was supported in part by the JHU + Amazon Initiative for Interactive AI (AI2AI). We also thank the anonymous county Department of Human Services for providing feedback and data for this work, and we thank Danish Pruthi, Krishna Pillutla, and Jessica Sorrell for their helpful feedback.

## 8 Limitations

The primary limitation of our work is the impossibility of considering all possible model and parameter configurations. While we selected high-performing models that we were able to fine-tune and evaluate on our compute resources, results may differ for different pre-trained language models. Similarly, while we select hyper-parameters based on prior work and conduct some ablation studies, text-generation is extremely compute-intensive and a fully exhaustive hyper-parameter sweep is not feasible. Overall our results emphasize the need to thoroughly evaluate models on target data and cannot necessarily be assumed to generalize to untested data.

There are also additional approaches we do not explore that could reduce privacy risk or improve the quality of synthetic data generated during training. Examples include combining text-anonymization with DP fine-tuning or selective constraints applied to the training data to reduce the frequency of entity mentions. However, this is difficult in practice, as real-world data is complex with, for example, the same people mentioned

across multiple CPS cases.

## 9 Ethical Considerations

Our work involves the use of private sensitive data, particularly the CPS data, which is not de-identified. To minimize risk, throughout this project we maintained a high level of data security, in compliance with IRB-approved protocol. The CPS data was exclusively stored on a secure restricted-access server with HIPPA-standard of security. All CPS experiments were conducted on this server, which also limited the models we could investigate. Our paper does not include any examples from either data set, in compliance with their respective data use agreements.

## References

- Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. [Deep learning with differential privacy](#). In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, CCS '16*, page 308–318, New York, NY, USA. Association for Computing Machinery.
- Md Momin Al Aziz, Tanbir Ahmed, Tasnia Faequa, Xiaoqian Jiang, Yiyu Yao, and Noman Mohammed. 2021. [Differentially private medical texts generation using generative neural networks](#). *ACM Trans. Comput. Healthcare*, 3(1).
- Mário Alvim, Konstantinos Chatzikokolakis, Catuscia Palamidessi, and Anna Pazi. 2018. [Invited paper: Local differential privacy on metric spaces: Optimizing the trade-off with utility](#). In *2018 IEEE 31st Computer Security Foundations Symposium (CSF)*, pages 262–267.
- Iyadh Ben Cheikh Larbi, Aljoscha Burchardt, and Roland Roller. 2023. [Clinical text anonymization, its influence on downstream NLP tasks and the risk of re-identification](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 105–111, Dubrovnik, Croatia. Association for Computational Linguistics.
- Hannah Brown, Katherine Lee, Fatemehsadat Mireshghallah, Reza Shokri, and Florian Tramèr. 2022. [What does it mean for a language model to preserve privacy?](#) In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency, FAccT '22*, page 2280–2292, New York, NY, USA. Association for Computing Machinery.
- Nicholas Carlini, Chang Liu, Jernej Kos, Úlfar Erlingsson, and Dawn Song. 2019. [The secret sharer: Measuring unintended neural network memorization & extracting secrets](#). *USENIX Security*.
- Zachary Charles, Arun Ganesh, Ryan McKenna, H Brendan McMahan, Nicole Mitchell, Krishna Pilutla, and Keith Rush. 2024. [Fine-tuning large language models with user-level differential privacy](#). *arXiv preprint arXiv:2407.07737*.
- Sai Chen, Fengran Mo, Yanhao Wang, Cen Chen, Jian-Yun Nie, Chengyu Wang, and Jamie Cui. 2023. [A customized text sanitization mechanism with differential privacy](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5747–5758, Toronto, Canada. Association for Computational Linguistics.
- Yanran Chen and Steffen Eger. 2023. [MENLI: Robust evaluation metrics from natural language inference](#). *Transactions of the Association for Computational Linguistics*, 11:804–825.
- Lynn Chua, Badih Ghazi, Yangsibo Huang, Pritish Kamath, Daogao Liu, Pasin Manurangsi, Amer Sinha, and Chiyuan Zhang. 2024. [Mind the privacy unit! User-level differential privacy for language model fine-tuning](#). In *Proceedings of the Conference on Language Modelin (COLM)*.
- Josep Domingo-Ferrer, David Sánchez, and Alberto Blanco-Justicia. 2021. [The limits of differential privacy \(and its misuse in data release and machine learning\)](#). *Commun. ACM*, 64(7):33–35.
- Ondřej Dušek and Zdeněk Kasner. 2020. [Evaluating semantic accuracy of data-to-text generation with natural language inference](#). In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 131–137, Dublin, Ireland. Association for Computational Linguistics.
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. 2006. [Calibrating noise to sensitivity in private data analysis](#). In *Theory of Cryptography*, pages 265–284, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Cynthia Dwork, Aaron Roth, et al. 2014. [The algorithmic foundations of differential privacy](#). *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. [Hierarchical neural story generation](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia. Association for Computational Linguistics.
- Oluwaseyi Feyisetan, Borja Balle, Thomas Drake, and Tom Diethe. 2020a. [Privacy- and utility-preserving textual analysis via calibrated multivariate perturbations](#). In *Proceedings of the 13th International Conference on Web Search and Data Mining, WSDM '20*, page 178–186, New York, NY, USA. Association for Computing Machinery.
- Oluwaseyi Feyisetan, Borja Balle, Thomas Drake, and Tom Diethe. 2020b. [Privacy- and utility-preserving](#)

- textual analysis via calibrated multivariate perturbations. In *Proceedings of the 13th International Conference on Web Search and Data Mining, WSDM '20*, page 178–186, New York, NY, USA. Association for Computing Machinery.
- Anjalie Field, Amanda Coston, Nupoor Gandhi, Alexandra Chouldechova, Emily Putnam-Hornstein, David Steier, and Yulia Tsvetkov. 2023. Examining risks of racial biases in nlp tools for child protective services. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pages 1479–1492.
- Nupoor Gandhi, Anjalie Field, and Emma Strubell. 2023. Annotating mentions alone enables efficient domain adaptation for coreference resolution. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10543–10558, Toronto, Canada. Association for Computational Linguistics.
- Nicolas Garneau and Luc Lamontagne. 2021. Trainable ranking models to evaluate the semantic accuracy of data-to-text neural generator. In *Proceedings of the 2nd Workshop on Evaluation and Comparison of NLP Systems*, pages 51–61, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. 2000. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *circulation*, 101(23):e215–e220.
- Sivakanth Gopi, Yin Tat Lee, and Lukas Wutschitz. 2021. Numerical composition of differential privacy. In *Advances in Neural Information Processing Systems*, volume 34, pages 11631–11642. Curran Associates, Inc.
- Jiaqi Guan, Runzhe Li, Sheng Yu, and Xuegong Zhang. 2018. Generation of synthetic electronic medical record text. In *2018 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 374–380.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text de-generation. In *International Conference on Learning Representations*.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Jinmiao Huang, Cesar Osorio, and Luke Wicent Sy. 2019. An empirical evaluation of deep learning for ICD-9 code assignment using MIMIC-III clinical notes. *Computer Methods and Programs in Biomedicine*, 177:141–153.
- Alistair Johnson, Tom Pollard, and Roger Mark. 2016a. MIMIC-III clinical database (version 1.4). *PhysioNet*, 10(C2XW26):2.
- Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016b. MIMIC-III, a freely accessible critical care database. *Scientific data*, 3(1):1–9.
- Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. 2019. CTRL: a conditional transformer language model for controllable generation. *Preprint*, arXiv:1909.05858.
- Yuval Kirstain, Ori Ram, and Omer Levy. 2021. Coreference resolution without span representations. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 14–19, Online. Association for Computational Linguistics.
- Alexey Kurakin, Natalia Ponomareva, Umar Syed, Liam MacDermed, and Andreas Terzis. 2023. Harnessing large-language models to generate private synthetic text. *arXiv preprint arXiv:2306.01684*.
- Xuechen Li, Florian Tramer, Percy Liang, and Tatsunori Hashimoto. 2021. Large language models can be strong differentially private learners. In *International Conference on Learning Representations*.
- N. Lukas, A. Salem, R. Sim, S. Tople, L. Wutschitz, and S. Zanella-Beguelin. 2023. Analyzing leakage of personally identifiable information in language models. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 346–363, Los Alamitos, CA, USA. IEEE Computer Society.
- Nuno Mamede, Jorge Baptista, and Francisco Dias. 2016. Automated anonymization of text documents. In *2016 IEEE Congress on Evolutionary Computation (CEC)*, pages 1287–1294.
- Justus Mattern, Zhijing Jin, Benjamin Weggenmann, Bernhard Schoelkopf, and Mrinmaya Sachan. 2022a. Differentially private language models for secure data sharing. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4860–4873, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Justus Mattern, Benjamin Weggenmann, and Florian Kerschbaum. 2022b. The limits of word level differential privacy. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 867–881, Seattle, United States. Association for Computational Linguistics.
- Michail Mersinias and Kyle Mahowald. 2023. For generated text, is NLI-neutral text the best text? In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 2596–2602, Singapore. Association for Computational Linguistics.

- Arvind Narayanan and Vitaly Shmatikov. 2008. [Robust de-anonymization of large sparse datasets](#). In *2008 IEEE Symposium on Security and Privacy (sp 2008)*, pages 111–125.
- Bhanudas Suresh Panchbhai and Varsha Makarand Pathak. 2022. A systematic review of natural language processing in healthcare. *Journal of Algebraic Statistics*, 13(1):682–707.
- Pranav Putta, Ander Steele, and Joseph W Ferrara. 2023. [Differentially private conditional text generation for synthetic data production](#).
- Weiyang Shi, Aiqi Cui, Evan Li, Ruoxi Jia, and Zhou Yu. 2022. [Selective differential privacy for language modeling](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2848–2859, Seattle, United States. Association for Computational Linguistics.
- Ondřej Sotolář, Jaromír Plhák, and David Šmahel. 2021. Towards personal data anonymization for social messaging. In *Text, Speech, and Dialogue*, pages 281–292, Cham. Springer International Publishing.
- Latanya Sweeney. 2000. Simple demographics often identify people uniquely. *Health*, 671.
- Ozlem Uzuner, Andreea Bodnari, Shuying Shen, Tyler Forbush, John Pestian, and Brett R South. 2012. Evaluating the state of the art in coreference resolution for electronic medical records. *Journal of the American Medical Informatics Association*, 19(5):786–791.
- Elena Volodina, Yousuf Ali Mohammed, Sandra Derbring, Arild Matsson, and Beata Megyesi. 2020. [Towards privacy by design in learner corpora research: A case of on-the-fly pseudonymization of Swedish learner essays](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 357–369, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Xinwei Wu, Li Gong, and Deyi Xiong. 2022. [Adaptive differential privacy for language model training](#). In *Proceedings of the First Workshop on Federated Learning for Natural Language Processing (FLANLP 2022)*, pages 21–26, Dublin, Ireland. Association for Computational Linguistics.
- Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. 2024. [Sheared LLaMA: Accelerating language model pre-training via structured pruning](#).
- Andrew Yale, Saloni Dash, Ritik Dutta, Isabelle Guyon, Adrien Pavao, and Kristin P. Bennett. 2020. [Generation and evaluation of privacy preserving synthetic health data](#). *Neurocomputing*, 416:244–255.
- Oleksandr Yermilov, Vipul Raheja, and Artem Chernodub. 2023. [Privacy- and utility-preserving NLP with anonymized data: A case study of pseudonymization](#). In *Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023)*, pages 232–241, Toronto, Canada. Association for Computational Linguistics.
- Ashkan Yousefpour, Igor Shilov, Alexandre Sablayrolles, Davide Testuggine, Karthik Prasad, Mani Malek, John Nguyen, Sayan Ghosh, Akash Bharadwaj, Jessica Zhao, Graham Cormode, and Ilya Mironov. 2022. [Opacus: User-friendly differential privacy library in pytorch](#). *Preprint*, arXiv:2109.12298.
- Xiang Yue, Minxin Du, Tianhao Wang, Yaliang Li, Huan Sun, and Sherman S. M. Chow. 2021. [Differential privacy for text analytics via natural text sanitization](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3853–3866, Online. Association for Computational Linguistics.
- Xiang Yue, Huseyin Inan, Xuechen Li, Girish Kumar, Julia McAnallen, Hoda Shajari, Huan Sun, David Levitan, and Robert Sim. 2023. [Synthetic text generation with differential privacy: A simple and practical recipe](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1321–1342, Toronto, Canada. Association for Computational Linguistics.
- Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun. 2020. [How does NLP benefit legal system: A summary of legal artificial intelligence](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5218–5230, Online. Association for Computational Linguistics.

## A Background: Differential Privacy

Differential privacy offers a formal privacy guarantee that ensures that any individual’s data cannot be inferred from a query applied to a dataset (Dwork et al., 2006, 2014). In other words, the result of such a query is nearly indistinguishable from the result of the same query applied to a dataset that either includes a modified version of the individual’s data or excludes the record entirely, thereby preserving the individual’s privacy. In this case, the notion of adjacency is defined as a difference of a single record in the original dataset  $D$  and the modified dataset  $D'$ .

Formally, differential privacy is defined as follows:

**Definition:** Given a dataset  $D$  and an adjacent dataset  $D'$ , which is produced by removing or modifying a single record from  $D$ , a randomized algorithm  $F : D \rightarrow Y$  is  $(\epsilon, \delta)$ -private if for any two neighboring datasets  $D, D'$ , with the constraints  $\epsilon > 0$  and  $\delta \in [0, 1]$ , the following holds true for all sets  $y \subseteq Y$ :

$$\Pr[F(D) \in y] \leq e^\epsilon \Pr[F(D') \in y] + \delta$$

The value of  $\epsilon$  denotes the privacy budget, while  $\delta$  specifies the likelihood that the privacy guarantee may fail. If  $\delta$  is set to 0, this implies a purely differentially private setting with no probability of the guarantee being broken. The value of  $\epsilon$  constrains how similar the outputs of both distributions are; a higher  $\epsilon$  value indicates a greater privacy budget, meaning the algorithm is less private. DP guarantees that even if an adversary has access to any side-knowledge, the privacy leakage of  $(\epsilon, \delta)$ -DP algorithms will not increase. Additionally, another property of DP is that it ensures that any post-processing on the outputs of  $(\epsilon, \delta)$ -differentially private algorithms will remain  $(\epsilon, \delta)$ -differentially private.

We use DP-SGD (Abadi et al., 2016), a modification to the stochastic gradient descent (SGD) algorithm, which is typically used to train neural networks. DP-SGD clips the gradients to limit the contribution of individual samples from the training data and subsequently adds noise from a predefined type of distribution (such as a Gaussian or Laplacian distribution) to the sum of the clipped gradients across all samples. DP-SGD thus provides a differentially private guarantee to obfuscate the gradient update, thereby ensuring that the contribution of any given sample in the training data is indistinguishable due to the aforementioned post-processing property. This process ensures  $(\epsilon, \delta)$ -differential privacy for each model update. Given a privacy budget, number of epochs, and other training parameters, we can estimate the privacy parameters using estimation algorithms (Gopi et al., 2021).

## B Hyperparameters

For training the autoregressive model, we used an effective batch size of 32 for training the non-differentially private model for both the CPS and the MIMIC-III data. For the differentially private fine-tuning, we used an effective batch size of 1024. We set the maximum sequence length to 1024 tokens and our training was conducted over 3 epochs, and training was optimized using the AdamW optimizer with its default hyperparameters. For the MIMIC-III data, our learning rate was set to  $3e-4$  (for the non-DP finetuned models) and to  $1e-3$  (for the DP-finetuned models). For the CPS data, we found a learning rate of  $3e-4$  in both cases was optimal. For the LoRA hyperparameters, we used a dimension of 4 and an alpha value of

32, specifically targeting the query (q\_proj) and value (v\_proj) projection layers of the transformer. To ensure training stability, we applied gradient clipping with a maximum gradient norm of 1.0. For the DP fine-tuning of the autoregressive model, we train with a privacy budget of  $\epsilon = 8$  for most of our experiments, and considering our relatively small dataset size we set  $\delta$  to  $1e-5$  for our experiments.

For training the downstream classifier, we conducted training over 3 epochs with a batch size of 8 and a maximum sequence length of 512 tokens. We utilized the AdamW optimizer with a learning rate of  $5e-5$ . We also conducted these downstream experiments with RoBERTa and found that the differences were minimal, with no impact on the overall trends, so we decided not to include these results.

During inference, we set the top-k sampling parameter to  $k = 50$  and the nucleus sampling parameter to  $p = 0.95$ . We generate approximately 30k and 31k samples for the child welfare data and diagnosis notes for the 10 most frequent ICD-9 codes, respectively, which are then used to train the downstream classifiers. We use similar inference hyperparameters for the instruction-tuned BioMistral-7B model for ICL, we set the top-k value to 50, top-p to 0.9 and the penalty-alpha parameter to 0.6.

Our experiments for all the aforementioned experimental setups used an A100 GPU for the MIMIC data and A6000 GPUs on a single secure server for the CPS data.

## C Fairness

The False Positive Equality Difference (FPED) metric is the sum of the differences between the overall false positive rate (FPR) for the entire dataset and the FPR for each subgroup  $d \in D$ , where  $D$  is a set consisting of all subgroups corresponding to a demographic attribute within the dataset.

$$\text{FPED} = \sum_{d=1}^D |\text{FPR}_{\text{overall}} - \text{FPR}_d| \quad (2)$$

$$\text{TNED} = \sum_{d=1}^D |\text{TNR}_{\text{overall}} - \text{TNR}_d| \quad (3)$$

Similarly, these ED metrics can be estimated for the true positive, true negative and false negative rates to estimate the TPED, TNED and FNED respectively. Lower values of these ED scores indicate that the model’s performance is more consistent across different subgroups.

	<b>FNED</b>	<b>FPED</b>	<b>TPED</b>	<b>TNED</b>	<b>Equalized Odds</b>
<b>Race</b>					
$D_{real}^{base}$	$0.35 \pm 0.0$	$0.01 \pm 0.0$	$0.35 \pm 0.0$	$0.01 \pm 0.0$	$0.20 \pm 0.0$
$D_{\epsilon=\infty}$	$0.39 \pm 0.005$	$0.02 \pm 0.003$	$0.39 \pm 0.005$	$0.02 \pm 0.003$	$0.23 \pm 0.001$
$D_{\epsilon=8}$	$0.53 \pm 0.014$	$0.01 \pm 0.003$	$0.53 \pm 0.014$	$0.01 \pm 0.003$	$0.30 \pm 0.003$
$D_{ICL}$	$0.49 \pm 0.03$	$0.03 \pm 0.006$	$0.49 \pm 0.03$	$0.03 \pm 0.006$	$0.29 \pm 0.012$
<b>Gender</b>					
$D_{real}^{base}$	$0.04 \pm 0.0$	$0.0 \pm 0.0$	$0.04 \pm 0.0$	$0.0 \pm 0.0$	$0.042 \pm 0.0$
$D_{\epsilon=\infty}$	$0.02 \pm 0.007$	$0.0 \pm 0.0$	$0.02 \pm 0.007$	$0.0 \pm 0.0$	$0.02 \pm 0.007$
$D_{\epsilon=8}$	$0.04 \pm 0.005$	$0.0 \pm 0.001$	$0.04 \pm 0.005$	$0.0 \pm 0.001$	$0.04 \pm 0.005$
$D_{ICL}$	$0.04 \pm 0.006$	$0.0 \pm 0.002$	$0.04 \pm 0.006$	$0.0 \pm 0.002$	$0.04 \pm 0.006$

Table 7: Fairness evaluation for the MIMIC-III ICD-9<sub>n=10</sub> task, for the gender and race categories.

The Equalized Odds ratio is calculated as follows:

$$EO_D = \max \left( \begin{array}{l} \max_{i \in D}(\text{TPR}_i) - \min_{i \in D}(\text{TPR}_i), \\ \max_{i \in D}(\text{FPR}_i) - \min_{i \in D}(\text{FPR}_i) \end{array} \right)$$

We have two categories of subgroups that are present in the MIMIC-III dataset over which we perform fairness evaluations with the downstream classifier trained over synthetic data with demographic control codes. The following categorical variables assigned to each within the dataset:

- **Gender:** Female, Male
- **Race/Ethnicity:** American Indian/Alaska Native, Asian, Black, Hispanic/Latino, Middle Eastern, Multi Race/Ethnicity, Other, Portuguese, South American, White

The format of the control code for the MIMIC-III data is as follows: *Long\_Title:* <diagnoses>, *ICD9\_CODE:* <codes>, *Gender:* <gender>, *Ethnicity:* <ethnicity>, where the <diagnoses> variable represents the long title form of the ICD-9 codes, information that is already provided with the MIMIC-III dataset.

## D Data Statistics

Our train/dev splits for the CPS, ICD-9<sub>n=10</sub>, ICD-9<sub>n=5</sub> and ICD-9<sub>n=3</sub> datasets the generative model was trained on are 89327/4701, 44215/2327, 37245/1960, 31317/1648 respectively.

The size of the train/dev sets for the models trained for downstream classification on the real ( $D_{real}$ ) and synthetic ( $D_{\epsilon=\infty, 8, 4}$ ) CPS data is 25385/2821, and the test set for this task consists of 4949 records.

For the ICD-9<sub>n=10</sub> multilabelling task, the real ( $D_{real}$ ) and synthetic ( $D_{\epsilon=\infty, 8, 4}$ ) train/dev split was the same, with  $\simeq 27920/3100$  for all models, and the test set size was  $\simeq 7500$  samples. For the ICD-9<sub>n=5</sub> task, the train/dev split was the same for all models  $\simeq 23520/2615$ , and the test set size was  $\simeq 6315$  samples. Similarly, for the ICD-9<sub>n=3</sub> task, the train/dev split was  $\simeq 19780 / 2200$ , and the test set size was  $\simeq 5310$  samples. Each of these experiments for the downstream tasks (coreference/mention detection & classification) was averaged over 3 runs.

## E Extended Privacy Evaluation results

In Table 4 we report the full set of canary results (for 1, 10, and 100 insertions, for each canary type). Results are generally similar across different numbers of insertions, in that DP generally reduces rank and perplexity, thus improving privacy, but does not eliminate all risk of leakage.

	MIMIC		CPS	
	Rank	Perplexity	Rank	Perplexity
<b>100 Insertions</b>				
Name	4628 / 3356	35.80 / 53.19	3025/2265	7.08/10.22
Address	5 / 3967	16.52 / 61.37	57/9401	13.01/26.49
Number	1 / 818	5.77 / 14.46	4/1285	9.25/22.88
Email	1 / 1410	10.50 / 70.26	1/5301	7.64/78.42
<b>10 Insertions</b>				
Name	5986 / 3378	49.72 / 54.10	3168 / 2306	10.62 / 10.33
Address	2276 / 4075	43.59 / 62.66	9618 / 9523	21.63 / 27.23
Number	902 / 841	9.43 / 14.61	474 / 1347	16.81 / 23.24
Email	711 / 1452	37.81 / 72.08	387 / 5838	49.91 / 81.61
<b>1 Insertion</b>				
Name	6037 / 3383	52.06 / 54.20	3164 / 2320	11.27 / 10.34
Address	3348 / 4081	54.50 / 62.79	9715 / 9529	24.65 / 27.29
Number	1084 / 838	9.82 / 14.63	1016 / 1357	20.78 / 23.27
Email	1941 / 1457	43.90 / 72.28	1771 / 5881	62.44 / 81.88
<b>0 Insertions</b>				
Name	6086 / 5265	44.81 / 57.54	3690 / 2321	11.89 / 10.34
Address	4565 / 3869	75.79 / 61.68	9890 / 9529	23.15 / 27.29
Number	1217 / 1522	11.80 / 13.33	1176 / 1357	21.27 / 23.28
Email	1003 / 3174	43.80 / 55.19	5893 / 5881	79.05 / 81.88

Table 8: Rank and perplexity metrics for canary attacks over MIMIC and CPS data. Each column is formatted as  $\epsilon = \infty/\epsilon = 8$ . Perplexity scores suggest that DP reduces privacy metrics for all canaries, and generally show similar privacy improvements.

Type	Canary Sequence
Name	Pt <u>Robin Kramer</u> is has been admitted to the hospital and the patient’s medical records were updated yesterday.
Address	The patient has been admitted to the hospital at <u>789 Elmford St, Maplegreen</u> .
Number	The patient’s contact number for emergencies is <u>(541) 545-6789</u> .
E-mail	In case of any inquiries, the patient can be contacted at <u>ptcareastps456@email.com</u> .

Table 9: The canaries inserted into the training data for the models fine-tuned to generate synthetic MIMIC-III data.

Type	Canary Sequence
Name	<u>Robin Kramer</u> has been placed in temporary foster care and there is an ongoing investigation into the child’s welfare.
Address	The CW visited the foster family’s address at <u>456 W Oak Avenue, Springfield, IL</u> .
Number	The case number <u>CW-2023-56893</u> has been assigned for tracking purposes.
E-mail	The CW can contact the foster family at <u>randuser789@xyzreportnews.com</u> in case of any emergencies.

Table 10: The canaries inserted into the training data for the models fine-tuned to generate synthetic CPS data.

Model	Data Size	Phrase Overlap Ratio	Total # of Phrase Overlap	Total # of Phrases	Total # of Deidentified Phrases Generated
$D_{(real, ICD-9_{n=10})}$	44215	1	2935955	2935955	3845112
$D_{(\epsilon=\infty, ICD-9_{n=10})}$	31020	0.00504	16271	3229278	369390
$D_{(\epsilon=8, ICD-9_{n=10})}$	31020	0.0314	10312	3281866	390956
$D_{(ICL, ICD-9_{n=10})}$	19640	0.00117	3761	3205098	316905
$D_{(real, ICD-9_{n=5})}$	37245	1	2565699	2565699	3352588
$D_{(\epsilon=\infty, ICD-9_{n=5})}$	26136	0.00478	13537	2831658	323963
$D_{(\epsilon=8, ICD-9_{n=5})}$	26136	0.00263	7447	2831627	295290

Table 11: Analysis for the MIMIC-III dataset of all the unique contexts in which entities of from all categories from the training data appear in the synthetic data, considering surrounding context word lengths varying from 1 to 4.  $D_{real}$  corresponds to the training data the generative models were trained on.

Model	Data Size	Phrase Overlap Ratio	Total of Phrase Overlap #	Total # of Phrases in $D_{real}$ + $D_{synth-data}$	Total # of Phrases in $D_{synth-data}$
$D_{\epsilon=\infty}$	28206	0.01517	6307	415685	104153
$D_{\epsilon=8}$	28206	0.00619	2448	395303	65990
$D_{\epsilon=4}$	28206	0.00567	2218	391313	59286

Table 12: Analysis for the CPS data of all the unique contexts in which entities of from all categories from the training data appear in the synthetic data, considering surrounding context word lengths varying from 1 to 4.  $D_{real}$  corresponds to the training data the generative models were trained on.