
GS-WGAN: A Gradient-Sanitized Approach for Learning Differentially Private Generators

Dingfan Chen¹

Tribhuvanesh Orekondy²

Mario Fritz¹

¹ CISP Helmholz Center for Information Security ² Max Planck Institute for Informatics
Saarland Informatics Campus, Germany

Abstract

The wide-spread availability of rich data has fueled the growth of machine learning applications in numerous domains. However, growth in domains with highly-sensitive data (e.g., medical) is largely hindered as the private nature of data prohibits it from being shared. To this end, we propose *Gradient-sanitized Wasserstein Generative Adversarial Networks* (GS-WGAN), which allows releasing a sanitized form of the sensitive data with rigorous privacy guarantees. In contrast to prior work, our approach is able to distort gradient information more precisely, and thereby enabling training deeper models which generate more informative samples. Moreover, our formulation naturally allows for training GANs in both centralized and federated (i.e., decentralized) data scenarios. Through extensive experiments, we find our approach consistently outperforms state-of-the-art approaches across multiple metrics (e.g., sample quality) and datasets.

1 Introduction

Releasing statistical and sensory data to a broad community has contributed towards advances in numerous machine learning (ML) techniques e.g., object recognition (ImageNet [10]), language modeling (RCV [22]), recommendation systems (Netflix ratings [6]). However, in many sensitive domains (e.g., medical, financial), similar advances are often held back as the private nature of collected data prohibits release in its original form. Privacy-preserving data publishing [5, 12, 15] provides a reasonable solution, where only a sanitized form of the original data (with rigorous privacy guarantees) is publicly released.

Traditionally, sanitization is performed in a differentially private (DP) framework [11]. The sanitization method employed is often hand-crafted for the given input data [26, 28, 42] and the specific data-dependent task the sanitized data is intended for (e.g., answering linear queries) [7, 13, 19, 34]. As a result, such sanitization techniques greatly restrict the expressiveness of the released data distribution and fail to generalize to novel tasks unanticipated by the publisher. Instead, recent privacy-preserving techniques [5, 40, 41, 43] build on top of successes in generative adversarial network (GANs) [16] literature, to generate synthetic data faithful to the original input distribution. Specifically, GANs are trained using a privacy-preserving algorithm (e.g., using DP-SGD [1]) and demonstrate promising results in modeling a variety of real-world high-dimensional data distributions. Common to most privacy-preserving training algorithms for neural network models is *manipulating* the gradient information generated during backpropagation. Manipulation most commonly involves clipping the gradients (to bound sensitivity) and adding calibrated random noise (to introduce stochasticity). Although recent techniques that employ such an approach demonstrate reasonable success, they are largely limited to shallow networks and fail to sufficiently capture the sample quality of the original data.

In this paper, towards the goal of a generative model capable of synthesizing high-quality samples in a privacy-preserving manner, we propose a differentially private GAN. We first identify that in

such a data-publishing scenario, only a subset of the trained model (specifically the generator) and its parameters need to be publicly-released. This insight allows us to surgically manipulate the gradient information during training, and thereby allowing more meaningful gradient updates. By coupling the approach with a Wasserstein [2] objective with gradient-penalty term [17], we further improve the amount of gradient information flow during training. The new objective additionally allows us to precisely estimate the gradient norms and analytically determine the sensitivity values. As an added benefit, we find our approach bypasses an intensive and fragile hyper-parameter search for DP-specific hyperparameters (particularly clipping values).

Contributions. (i) A novel gradient-sanitized Wasserstein GAN (GS-WGAN), which is capable of generating high-dimensional data with DP guarantee; (ii) Our approach naturally extends to both centralized and decentralized datasets. In the case of decentralized scenarios, our work can provide user-level DP guarantee [24] under an untrusted server; (iii) Extensive evaluations on various datasets demonstrate that our method significantly improves the sample quality of privacy-preserving data over state-of-the-art approaches.

2 Related Work

We review several generative models in the area of differential privacy, as well as their relations to our work.

DP-SGD GAN. Training GANs via DP-SGD [1, 5, 36, 40, 43] has proven effective in generating high-dimensional sanitized data. However, DP-SGD relies on carefully tuning of the clipping bound of gradient norm, i.e., the sensitivity value. Specifically, the optimal clipping bound varies greatly with the model architecture and the training dynamics, making the implementation of DP-SGD difficult. Unlike DP-SGD, our framework provides precise estimation of the sensitivity value, avoiding the intensive search of hyper-parameters.

PATE. Private Aggregation of Teacher Ensembles (PATE) is recently adapted to generative models and two main approaches were studied: PATE-GAN [41] and G-PATE [27]. PATE-GAN trained multiple teacher discriminators on disjoint data partitions together with a student discriminator. In contrast, we consider a simplified model without a student discriminator.

G-PATE [27] is similar to our work in the sense that, both works trained the discriminator non-privately while only training the generator with DP guarantee, and both sanitized gradients that the generator received from the discriminator. However, G-PATE suffers from two main limitations: (i) gradients need to be discretized by using manually selected bins in order to suit for the PATE framework and (ii) high-dimensional gradients in the PATE framework bring high privacy costs and thus dimension reduction techniques are required. Our framework can effectively avoid these two limitations and achieve better sample quality due to the novel gradient sanitation, see our experiments.

Fed-Avg GAN [3]. While many works focus on centralized setting, the decentralized case has rarely been studied. To address this, Federated Average GAN (Fed-Avg GAN) proposed to adapt GAN training by using the DP-Fed-Avg [29] algorithm, providing user-level DP guarantee under trusted server. In comparison with Fed-Avg GAN that merely works on decentralized data, our work can tackle both centralized and decentralized data using a single framework. Note that Fed-Avg sanitized parameter gradients of the discriminator in a similar way to DP-SGD, it also suffers from the difficulty of turning hyper-parameters.

3 Background

DP provides rigorous privacy guarantees for algorithms while allowing for quantitative privacy analysis. We below present several definitions and theorems that will be used in this work.

Definition 3.1. (Differential Privacy (DP) [11]) A randomized mechanism \mathcal{M} with range \mathcal{R} is (ϵ, δ) -DP, if

$$\Pr[\mathcal{M}(S) \in \mathcal{O}] \leq e^\epsilon \cdot \Pr[\mathcal{M}(S') \in \mathcal{O}] + \delta \quad (1)$$

holds for any subset of outputs $\mathcal{O} \subseteq \mathcal{R}$ and for any adjacent datasets S and S' , where S and S' differ from each other with only one training example. \mathcal{M} is the GAN training algorithm in our case, ϵ corresponds to the upper bound of privacy loss, and δ is the probability of breaching DP constraints.

Definition 3.2. (Rényi Differential Privacy (RDP) [31]) A randomized mechanism \mathcal{M} is (λ, ε) -RDP with order λ , if

$$D_\lambda(\mathcal{M}(S) \parallel \mathcal{M}(S')) = \frac{1}{\lambda - 1} \log \mathbb{E}_{x \sim \mathcal{M}(S)} \left[\left(\frac{\Pr[\mathcal{M}(S) = x]}{\Pr[\mathcal{M}(S') = x]} \right)^{\lambda-1} \right] \leq \varepsilon \quad (2)$$

holds for any adjacent datasets S and S' , where $D_\lambda(P \parallel Q) = \frac{1}{\lambda-1} \log \mathbb{E}_{x \sim Q} [(P(x)/Q(x))^\lambda]$ denotes the Rényi divergence. Moreover, a (λ, ε) -RDP mechanism \mathcal{M} is also $(\varepsilon + \frac{\log 1/\delta}{\lambda-1}, \delta)$ -DP.

In contrast to DP, RDP provides convenient composition properties to accumulate privacy cost over a sequence of mechanisms.

Theorem 3.1. (Composition) For a sequence of mechanisms $\mathcal{M}_1, \dots, \mathcal{M}_k$ s.t. \mathcal{M}_i is (λ, ε_i) -RDP $\forall i$, the composition $\mathcal{M}_1 \circ \dots \circ \mathcal{M}_k$ is $(\lambda, \sum_i \varepsilon_i)$ -RDP.

Definition 3.3. (Gaussian Mechanism [14, 31]) Let $f : X \rightarrow \mathbb{R}^d$ be an arbitrary d -dimensional function with sensitivity being

$$\Delta_2 f = \max_{S, S'} \|f(S) - f(S')\|_2 \quad (3)$$

over all adjacent datasets S and S' . The Gaussian Mechanism \mathcal{M}_σ , parameterized by σ , adds noise into the output, i.e.,

$$\mathcal{M}_\sigma(x) = f(x) + \mathcal{N}(0, \sigma^2 I). \quad (4)$$

\mathcal{M} is $(\lambda, \frac{\lambda \Delta_2^2 f^2}{2\sigma^2})$ -RDP.

Theorem 3.2. (Post-processing [14]) If \mathcal{M} satisfies (ε, δ) -DP, $F \circ \mathcal{M}$ will satisfy (ε, δ) -DP for any function F with \circ denoting the composition operator.

4 Proposed Method

Generative Adversarial Networks (GANs) [16]. Our approach models the underlying (private) data distribution using a generative neural network, building on top of recent successes of GANs. GANs (see Fig. 1(a)) formulate the task of sample generation as a zero-sum two-player game, between two neural network models: discriminator D and generator G . The discriminator D is rewarded for correctly classifying whether a given sample is ‘real’ (i.e., from the input data distribution) or ‘fake’ (generated by the generator). In contrast, the task of the generator G is (given some random noise z) to generate samples which fool the discriminator (i.e., causes misclassifications). After training the models in an adversarial manner, the discriminator is discarded and the generator is used as a proxy to draw samples from the original distribution.

Differentially Private GANs. Releasing the generator as a substitute for the original training data distribution entails privacy risks. Consequently, along the lines of recent work [5, 36, 40, 43], our goal is instead to train the GAN in a privacy-preserving manner, such that any privacy leakage upon disclosing the generator is bounded. A simple approach towards the goal is replacing the typical training procedure (SGD) with a differentially private variant (DP-SGD [1]) and thereby limiting the contribution of a particular training example in the final trained model. DP-SGD enforces the desired privacy requirement by (i) clipping the gradients \mathbf{g}_t to have an L_2 -norm of C at each training step; and (ii) sampling random noise and adding it to the gradients, before performing descent on the trained parameters $\boldsymbol{\theta}$:

$$\mathbf{g}^{(t)} := \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_D, \boldsymbol{\theta}_G) \quad (\text{gradient}) \quad (5)$$

$$\hat{\mathbf{g}}^{(t)} := \mathcal{M}_{\sigma, C}(\mathbf{g}^{(t)}) = \text{clip}(\mathbf{g}^{(t)}, C) + \mathcal{N}(0, \sigma^2 C^2 I) \quad (\text{sanitization mechanism}) \quad (6)$$

$$\boldsymbol{\theta}^{(t+1)} := \boldsymbol{\theta}^{(t)} - \eta \cdot \hat{\mathbf{g}}^{(t)} \quad (\text{gradient descent step}) \quad (7)$$

While such an approach provides rigorous privacy guarantees, there are multiple shortcomings: (i) the sanitization mechanism $\mathcal{M}_{\sigma, C}$, primarily due to clipping, significantly destroys the original gradient information, and thereby affects utility; and (ii) finding a reasonable clipping value C in the mechanism to balance utility with privacy is especially challenging. In particular, as the gradient norms exhibit a heavy-tailed distribution, choosing a clipping value requires an exhaustive

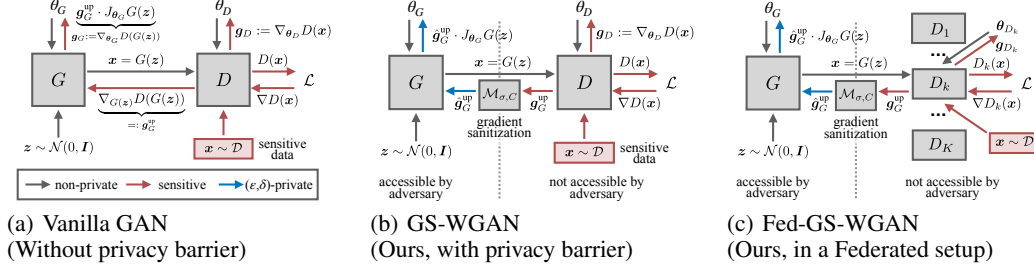


Figure 1: Approach outline. Our gradient sanitization scheme ensures DP training of the generator.

search. Moreover, since the clipping value is extremely sensitive to many other hyperparameters (e.g., learning rate, architecture), it requires persistent re-tuning. Now, we discuss how we address these shortcomings within our differentially private GAN approach.

Selectively applying Sanitization Mechanism. We begin by exploiting the fact that after training the GAN, only the generator G is released. Consequently, we can perform gradient steps by selectively applying the sanitization mechanism only to the corresponding subset of parameters θ_G :

$$\theta_D^{(t+1)} := \theta_D^{(t)} - \eta_D \cdot g_D^{(t)} \quad (\hat{g}_D^{(t)} = g_D^{(t)}; \text{Discriminator}) \quad (8)$$

$$\theta_G^{(t+1)} := \theta_G^{(t)} - \eta_G \cdot \hat{g}_G^{(t)} \quad (\hat{g}_G^{(t)} = \mathcal{M}_{\sigma, C}(g_G^{(t)}); \text{Generator}) \quad (9)$$

Apart from reducing the number of parameters sanitized, this also provides a benefit of more reliably training a discriminator. In addition, we exploit the chain rule to further narrow the scope of the sanitization mechanism:

$$g_G = \nabla_{\theta_G} \mathcal{L}_G(\theta_G) = \nabla_{G(z; \theta_G)} \mathcal{L}_G(\theta_G) \cdot J_{\theta_G} G(z; \theta_G) \quad (10)$$

$$\hat{g}_G = \mathcal{M}_{\sigma, C}(\underbrace{\nabla_{G(z)} \mathcal{L}_G(\theta_G)}_{g_G^{\text{upstream}}}) \cdot \underbrace{J_{\theta_G} G(z; \theta_G)}_{J_G^{\text{local}}} \quad (11)$$

The above becomes easier to intuit by considering a typical loss function $\mathcal{L}_G(\theta_G) = -D(G(z; \theta_G))$. As illustrated in Fig. 1(b), Eq. 11 can then be considered as placing the privacy barrier for gradient information backpropagating from the discriminator back to the generator, by applying the sanitization mechanism on g_G^{upstream} . Note that the second term (J_G^{local}) is the local generator jacobian computed independent of training data, and hence does not require sanitization. Consequently, using a more precise application of the sanitization mechanism on the gradient information, our goal here is to maximally preserve the true gradient direction during training.

Bounding sensitivity using Wasserstein distance. To bound the sensitivity of the optimizer on individual training examples, a key step in sanitization mechanisms is to *clip* (Eq. 6) the gradient vector g (Eq. 5) before updating parameters (Eq. 7). Clipping is typically performed in L_2 norm, by replacing the gradient vector g by $g/\max(1, \|g\|_2/C)$ to ensure $\|g\|_2 \leq C$. However, clipping significantly destroys gradient information, as reasonable choices of C (e.g., 4 [1]) are significantly lower than the gradient-norms observed (12 ± 10 in our case) when training neural networks using standard loss functions. We propose to alleviate the issue by leveraging a more suitable loss function, which generates bounded gradients (with norms close to 1) by construction. Specifically, we use as our loss the Wasserstein-1 metric [2], which measures the statistical distance between the real and generated data distributions. Here, the training process can be interpreted as minimizing integral probability metrics (IPMs) $\sup_{f \in \mathcal{F}} |\int_M f dP - \int_M f dQ|$ between real (P) and generated (Q) data distributions, where $\mathcal{F} = \{f : \|f\|_L \leq 1\}$ (i.e., the discriminator function f is 1-Lipschitz continuous). It follows that the optimal discriminator has a gradient norm being 1 almost everywhere under P and Q .

We incorporate the norm constraint into our training objective in the form of a gradient penalty term [17]:

$$\mathcal{L}_D = -\mathbb{E}_{x \sim P}[D(x)] + \mathbb{E}_{\tilde{x} \sim Q}[D(\tilde{x})] + \lambda \mathbb{E}[(\|\nabla D(\alpha x + (1 - \alpha)\tilde{x})\|_2 - 1)^2] \quad (12)$$

$$\mathcal{L}_G = -\mathbb{E}_{z \sim P_z}[D(G(z))] \quad (13)$$

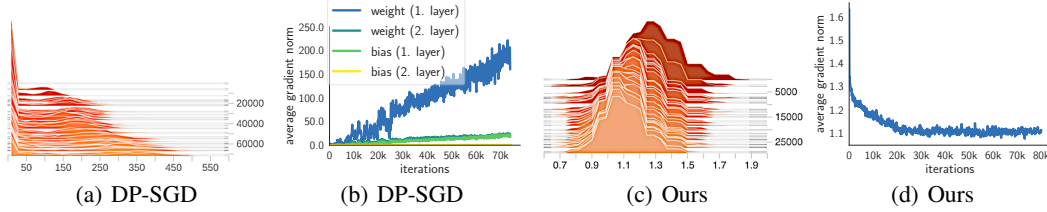


Figure 2: Gradient norm (before clipping) dynamics during the GAN training process. In the experiment, the clipping bound is chosen to be 1 and 1.1 in 2(c) and 2(a) respectively.

where \mathcal{L}_D and \mathcal{L}_G represent training objectives for the discriminator and the generator, respectively. λ is the hyper-parameter for weighting the gradient penalty term and P_z denotes the prior distribution for the latent code variable z . The variable $\alpha \sim \mathcal{U}[0, 1]$, uniformly sampled from $[0, 1]$, regulates the interpolation between real and generated samples.

As a natural consequence of the Wasserstein objective, reducing the norms of our target gradient g_G^{upstream} (Equation 11) during training is integrated in our training objective (last term in Equation 12). Consequently, we observe significantly lower gradient norms during training (see Fig. 2(c)-2(d)) compared to training using a standard GAN loss (see Fig. 2(a)-2(b)). As a result, bounding the sensitivity (i.e., gradient norms) is now largely delegated to our training procedure and clipping using the sanitization mechanism destroys significantly less information. Additionally, by choosing a clipping threshold of $C=1$ (i.e., $\|g\|_2 \leq 1$), we obtain a fixed and bounded sensitivity, and eliminate the need for intensive hyper-parameter search for the optimal clipping threshold. Following this normalization solution, a data-independent privacy cost can be determined by the following theorem, whose proof is provided in Appendix.

Theorem 4.1. Each generator update step satisfies $(\lambda, 2B\lambda/\sigma^2)$ -RDP where B is the batch size.

Privacy Amplification by Subsampling. A well-known approach for increasing privacy of a mechanism is to apply the mechanism to a random subsample of the database, rather than on the entire dataset [4, 25, 38]. Intuitively, subsampling decreases the chances of leaking information about a particular individual since nothing about that individual can be leaked once the individual is not included in the subsample. In order to further reduce the privacy cost, we subsample the whole dataset into different subsets and train multiple discriminators independently on each subset. At each training step, the generator randomly queries one discriminator while the selected discriminator updates its parameters on the generated data and its associated subsampled dataset.

Extending to Federated Learning. In addition to improving the privacy guarantee, performing subsampling in our setup also naturally accommodates training a generative model on decentralized datasets (with a discriminator trained on each disjoint data subset). Recently, Augenstein et al. [3] identified such techniques are extremely relevant when training models in a federated setup [30], i.e., when the training data is private and distributed among edge devices. We outline our method to train a differentially private GAN in a federated setup in Figure 1(c) and remark some subtle differences between our approach and Fed-Avg GAN [3] here: (i) the discriminators are retained at each client in our framework while they are shared between the server and client in Fed-Avg GAN; (ii) the gradients are sanitized at each client before sending to the server, with which we provide DP guarantee even under an untrusted server. In contrast, the unprocessed information is accumulated at the server before being sanitized in Fed-Avg GAN; and (iii) The gradients w.r.t. the samples are transferred in GS-WGAN, while Fed-Avg GAN transfers the gradients w.r.t. discriminator network parameters.

5 Experiment

5.1 Experiment Setup

To validate the applicability of our method to high-dimensional data, we conduct experiments on image datasets. In line with previous works, we use MNIST [21] and Fashion-MNIST [39] dataset. We model the joint distribution of images and the corresponding labels, i.e., the label is supplied to both the generator and the discriminator, and the image is generated conditioned on the input. During both training and inference, we use a uniform prior distribution for generating labels, which

| | | IS \uparrow | FID \downarrow | MLP \uparrow Acc | CNN \uparrow Acc | Avg \uparrow Acc | Calibrated \uparrow Acc |
|---------------|---------------------|---------------|------------------|-----------------------|-----------------------|-----------------------|------------------------------|
| MNIST | Real | 9.80 | 1.02 | 0.98 | 0.99 | 0.88 | 100 % |
| | G-PATE ¹ | 3.85 | 177.16 | 0.25 | 0.51 | 0.34 | 40% |
| | DP-SGD GAN | 4.76 | 179.16 | 0.60 | 0.63 | 0.52 | 59% |
| | DP-Merf | 2.91 | 247.53 | 0.63 | 0.63 | 0.57 | 66% |
| | DP-Merf AE | 3.06 | 161.11 | 0.54 | 0.68 | 0.42 | 47% |
| | Ours | 9.23 | 61.34 | 0.79 | 0.80 | 0.60 | 69% |
| Fashion-MNIST | Real | 8.98 | 1.49 | 0.88 | 0.91 | 0.79 | 100% |
| | G-PATE | 3.35 | 205.78 | 0.30 | 0.50 | 0.40 | 54% |
| | DP-SGD GAN | 3.55 | 243.80 | 0.50 | 0.46 | 0.43 | 53% |
| | DP-Merf | 2.32 | 267.78 | 0.56 | 0.62 | 0.51 | 65% |
| | DP-Merf AE | 3.68 | 213.59 | 0.56 | 0.62 | 0.45 | 55% |
| | Ours | 5.32 | 131.34 | 0.65 | 0.65 | 0.53 | 67% |

Table 1: Quantitative Results on MNIST and Fashion-MNIST ($\epsilon = 10$, $\delta = 10^{-5}$)

is independent of the training dataset and thus does not incur additional privacy cost (in contrast, [36] needs to assume the labels are non-private).

Evaluation Metrics. We evaluate along two fronts: *privacy* (determined by ϵ) and *utility*. For utility, we consider two metrics: (a) *sample quality*: realism of the samples produced – evaluated by **Inception Score (IS)** [23, 35] and **Frechet Inception Distance (FID)** [20] (standard in GAN literature); and (b) *usefulness for downstream tasks*: we train downstream classifiers on 60k privately-generated data points and evaluate the prediction accuracy on real test set. We consider Multi-layer Perceptrons (MLP), Convolutional Neural Networks (CNN) and 11 scikit-learn [32] classifiers (e.g., SVMs, Random Forest). We include the following metrics in the main paper: **MLP Acc** (MLP accuracy), **CNN Acc** (CNN accuracy), **Avg Acc** (Averaged accuracy of all classification models), **Calibrated Acc** (Averaged accuracy of all classification models normalized by the accuracy when trained on real data). The detailed results are presented in Appendix.

Architecture and Warm-start. We highlight two strategies adopted for improving the sample quality as well as reducing the privacy cost: (i) *Better model architecture*: While previous works are limited to shallow networks and thereby bottle-necking generated sample quality, our framework allows stable training with a complex model architecture (DCGAN [33] architecture for the discriminator, ResNet architecture (adapted from BigGAN [8]) for the generator) to help improve the sample quality; and (ii) *Discriminator warm-starting*: To bootstrap the training process, we pre-train discriminators along with a non-private generator for a few steps, and we subsequently train the private generator by using the warm-starting values of the discriminators. Note that our framework allows pre-training on the original private dataset without compromising privacy (in contrast, [43] needs to use external public datasets).

5.2 Comparison with Baselines

Baselines. We consider the following state-of-the-art methods designed for DP high-dimensional data generation: **DP-Merf** and **DP-Merf AE** [18], **DP-SGD GAN** [36, 40, 43], and **G-PATE** [27]. While PATE-GAN [41] demonstrates promising results on low-dimensional data, we currently do not consider it as we were unable to extend it to our image datasets (more details in appendix) for a fair comparison. For DP-Merf, DP-Merf AE, and G-PATE, we use the source code provided by the authors. For DP-SGD GAN, we adopt the implementation of [36], which is the only work that provides executable code with privacy analysis. For a fair comparison, we evaluate all methods with a privacy budget of $(\epsilon, \delta) = (10, 10^{-5})$ (consistently used in previous works) over 60K generated samples.

Results. We present the qualitative results in Figure 3 and the quantitative results in Table 1. In terms of sample quality, we find (Table 1, columns IS and FID) our method consistently provides significant improvements over baselines. For instance, considering inception scores, we find a relative improvement of 94% (9.23 vs. 4.76 of DP-SGD GAN) on MNIST and 45% on Fashion-MNIST (5.32 vs. 3.68 of DP-Merf AE).

¹ PATE provides data-dependent ϵ , i.e., publishing ϵ value will introduce privacy cost. Thus, G-PATE is not directly comparable to other methods and is excluded from our analysis study (section 5.3).






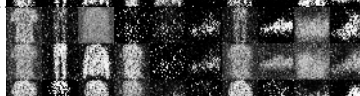
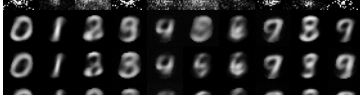
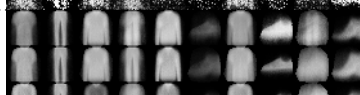

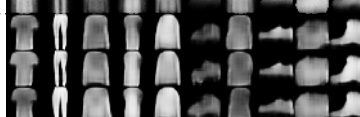
| Method | MNIST | Fashion-MNIST |
|------------|---|--|
| G-PATE |  |  |
| DP-SGD GAN |  |  |
| DP-Merf |  |  |
| DP-Merf AE |  |  |
| Ours |  |  |

Figure 3: Generated samples with $(\epsilon, \delta) = (10, 10^{-5})$

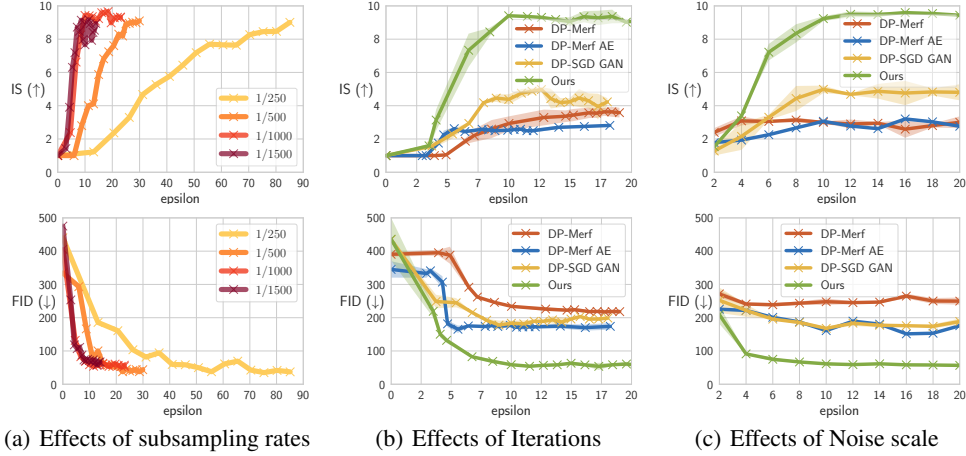


Figure 4: Privacy-utility trade-off on MNIST with $\delta = 10^{-5}$. (Top row: IS. Bottom row: FID.)

Furthermore, our method also generates samples that better capture the statistical properties of the original data and are thereby making aiding performances of downstream tasks. For instance, our approach increases performance of a downstream MLP classifier (Table 1, column MLP Acc) by 25% (0.79 vs. 0.63 of DP-Merf) on MNIST and 16% (0.65 vs. 0.56 of DP-Merf) on Fashion-MNIST. In a word, our approach demonstrates significant improvements across multiple metrics and high-dimensional image datasets.

5.3 Influence of Hyperparameters

The privacy/utility performances of our approach is primarily determined by three factors: (i) **sub-sampling rates** γ , (ii) number of training **iterations**, and (iii) **noise scale** σ . We now investigate how these factors influence privacy cost ϵ and utility (sample quality measured by IS and FID), and additionally compare with baselines:

(i) **Subsampling rates**: We evaluate the sample quality of our method considering multiple choices of subsampling rates ($\gamma \in [1/250, 1/500, 1/1000, 1/1500]$) over the training iterations. The results are presented in Figure 4(a), where the x -axis corresponds to the ϵ value evaluated at different iterations. We observe that the sub-sampling rate should be sufficiently small for achieving a reasonable sample quality while providing a strong privacy guarantee. A value of 1/1000 yields relatively good privacy-utility trade-off, while further decreasing the sub-sampling rate does not

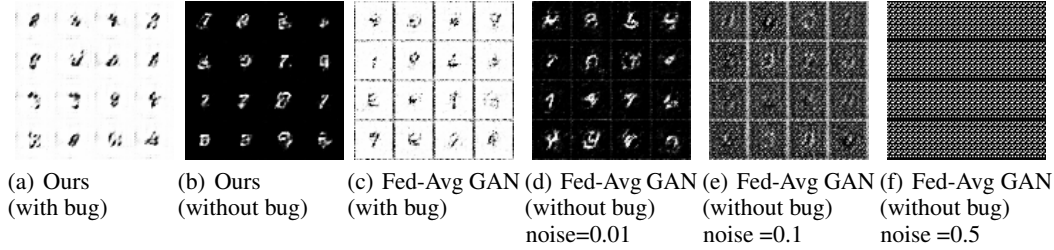


Figure 5: Qualitative Results on Federated EMNIST.

necessarily improve the results. *(ii) Iterations:* We evaluate all methods during the course of training, where more iterations lead to higher utilities, but at the expense of accumulating a higher privacy cost ϵ . From Figure 4(b), we find our approach yields better sample qualities with fewer iterations (and hence lower ϵ). Specifically, across the range of iterations, we find IS increases by 10-90%, while the FID decreases by 20-60% compared to baselines. *(iii) Noise scale:* We calibrate the noise scale of each method to certain privacy budget ϵ and show the resulting privacy-utility curves in Figure 4(c). Similar to the previous case, our method achieves a consistent improvement in both metrics spanning a broad range of noise scale (privacy budget ϵ).

5.4 Federated Setting Evaluation

Our approach allows to perform privacy-preserving training of a GAN in federated setup, where sensitive user dataset is partitioned across K clients (e.g., edge devices). Such a training scheme is useful to privately inspect data for debugging. For evaluation, we consider a real-world debugging task introduced in [3]: to detect the erroneous flipping of pixel intensities, which occurs in a fraction of client devices. Two GAN models are trained: one on client data that are suspected to be erroneously flipped (with bug) and one on the client data that are believed to be normal (without bug). The samples generated by these two GAN models should exhibit different appearance such that the bug can be detected by inspecting the generated samples.

We conduct experiments on the Federated EMNIST dataset [9] and compare our GS-WGAN with **Fed-Avg GAN** [3]. As shown in Figure 5(a) and 5(b), the presence of bug is clearly identifiable by inspecting the samples generated by our model. Moreover, as shown in Table 2, our GS-WGAN yields better sample quality ($0.28\times$ smaller FID) with a significantly lower privacy cost ($10^4\times$ smaller ϵ) compared to Fed-Avg GAN. Furthermore, our method shows better robustness against large injected noise. This is illustrated in Figure 5(e) and 5(f): a noise scale larger than 0.1 inevitably leads to failure in training Fed-Avg GAN, whereas our method can tolerate 10 times larger noise scale. In addition, we show in the last column of Table 2 the amortized communication cost (CT) required for performing one update step on the generator. Specifically, this corresponds to the total number of transferred bytes (including both server-to-client and client-to-server) averaged over all participating clients. Our GS-WGAN allows each client to retain its discriminator locally and only the gradients w.r.t. generated samples are communicated (which is significantly more compact than gradients w.r.t. model parameters, as done by Fed-Avg GAN). We observe that GS-WGAN achieves a magnitude of 10^2 gain in reducing the communication cost.

| | IS \uparrow | FID \downarrow | epsilon \downarrow | CT (byte) \downarrow |
|-------------|---------------|------------------|----------------------|-------------------------|
| Fed Avg GAN | 10.88 | 218.24 | 9.99×10^6 | $\sim 3.94 \times 10^7$ |
| Ours | 11.25 | 60.76 | 5.99×10^2 | $\sim 1.50 \times 10^5$ |

Table 2: Quantitative Results on Federated EMNIST ($\delta = 1.15 \times 10^{-3}$)

6 Conclusion

In this paper, we presented a differentially-private approach *GS-WGAN* to sanitize sensitive high-dimensional datasets with provable privacy guarantees while simultaneously preserving informativeness of the sanitized samples. Our primary insight is that privacy-preserving training (which sacrifices utility) can be selectively applied only to the generator (which is publicly released) while the discriminator (which is discarded post-training) can be trained optimally. Additionally, introducing a Wasserstein training objective allowed us to exploit the Lipschitz property of the discriminator and led to tighter estimates of sensitivity values to better estimate privacy. Our extensive evaluation

presents encouraging results: sensitive datasets can be effectively distilled to sanitized forms which nonetheless preserves informativeness of the data and allows training downstream models.

References

- [1] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (CCS)*, 2016.
- [2] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In *International Conference on Machine Learning (ICML)*, 2017.
- [3] S. Augenstein, H. B. McMahan, D. Ramage, S. Ramaswamy, P. Kairouz, M. Chen, R. Mathews, and B. A. y Arcas. Generative models for effective ML on private, decentralized datasets. In *International Conference on Learning Representations (ICLR)*, 2020.
- [4] B. Balle, G. Barthe, and M. Gaboardi. Privacy amplification by subsampling: Tight analyses via couplings and divergences. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [5] B. K. Beaulieu-Jones, Z. S. Wu, C. Williams, and C. S. Greene. Privacy-preserving generative deep neural networks support clinical data sharing. *bioRxiv*. DOI, 10, 2017.
- [6] J. Bennett, S. Lanning, et al. The netflix prize. In *Proceedings of KDD cup and workshop*, volume 2007. Citeseer, 2007.
- [7] A. Blum, K. Ligett, and A. Roth. A learning theory approach to noninteractive database privacy. *Journal of the ACM (JACM)*, 60(2), 2013.
- [8] A. Brock, J. Donahue, and K. Simonyan. Large scale gan training for high fidelity natural image synthesis. In *International Conference on Learning Representations (ICLR)*, 2018.
- [9] S. Caldas, P. Wu, T. Li, J. Konečný, H. B. McMahan, V. Smith, and A. Talwalkar. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097*, 2018.
- [10] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [11] C. Dwork. Differential privacy: A survey of results. In *International Conference on Theory and Applications of Models of Computation (TAMC)*. Springer, 2008.
- [12] C. Dwork, M. Naor, O. Reingold, G. N. Rothblum, and S. Vadhan. On the complexity of differentially private data release: efficient algorithms and hardness results. In *Proceedings of the forty-first annual ACM symposium on Theory of computing (STOC)*, 2009.
- [13] C. Dwork, G. N. Rothblum, and S. Vadhan. Boosting and differential privacy. In *IEEE 51st Annual Symposium on Foundations of Computer Science (FOCS)*, 2010.
- [14] C. Dwork, A. Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4), 2014.
- [15] B. C. Fung, K. Wang, R. Chen, and P. S. Yu. Privacy-preserving data publishing: A survey of recent developments. *ACM Computing Surveys (Csur)*, 42(4), 2010.
- [16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems (NeurIPS)*, 2014.
- [17] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [18] F. Harder, K. Adamczewski, and M. Park. Differentially private mean embeddings with random features (dp-merf) for simple & practical synthetic data generation. *arXiv preprint arXiv:2002.11603*, 2020.

- [19] M. Hardt and G. N. Rothblum. A multiplicative weights mechanism for privacy-preserving data analysis. In *IEEE 51st Annual Symposium on Foundations of Computer Science (FOCS)*, 2010.
- [20] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [21] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 1998.
- [22] D. D. Lewis, Y. Yang, T. G. Rose, and F. Li. Rcv1: A new benchmark collection for text categorization research. *Journal of Machine Learning Research (JMLR)*, 5(Apr), 2004.
- [23] C. Li, H. Liu, C. Chen, Y. Pu, L. Chen, R. Henao, and L. Carin. Alice: Towards understanding adversarial learning for joint distribution matching. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [24] J. Li, M. Khodak, S. Caldas, and A. Talwalkar. Differentially private meta-learning. In *International Conference on Learning Representations (ICLR)*, 2020.
- [25] N. Li, W. Qardaji, and D. Su. On sampling, anonymization, and differential privacy or, k-anonymization meets differential privacy. In *Proceedings of the 7th ACM Symposium on Information, Computer and Communications Security*, 2012.
- [26] N. Li, M. Lyu, D. Su, and W. Yang. Differential privacy: From theory to practice. *Synthesis Lectures on Information Security, Privacy, & Trust*, 8(4), 2016.
- [27] Y. Long, S. Lin, Z. Yang, C. A. Gunter, and B. Li. Scalable differentially private generative student model via pate. *arXiv preprint arXiv:1906.09338*, 2019.
- [28] R. Mckenna, D. Sheldon, and G. Miklau. Graphical-model based estimation and inference for differential privacy. In *International Conference on Machine Learning (ICML)*, 2019.
- [29] B. McMahan, D. Ramage, K. Talwar, and L. Zhang. Learning differentially private recurrent language models. In *International Conference on Learning Representations (ICLR)*, 2018.
- [30] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, et al. Communication-efficient learning of deep networks from decentralized data. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017.
- [31] I. Mironov. Rényi differential privacy. In *IEEE 30th Computer Security Foundations Symposium (CSF)*, 2017.
- [32] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research (JMLR)*, 12, 2011.
- [33] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In Y. Bengio and Y. LeCun, editors, *International Conference on Learning Representations (ICLR)*, 2016.
- [34] A. Roth and T. Roughgarden. Interactive privacy via the median mechanism. In *Proceedings of the forty-second ACM symposium on Theory of computing (STOC)*, 2010.
- [35] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- [36] R. Torkzadehmahani, P. Kairouz, and B. Paten. Dp-cgan: Differentially private synthetic data and label generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019.
- [37] T. Van Erven and P. Harremoës. Rényi divergence and kullback-leibler divergence. *IEEE Transactions on Information Theory*, 60(7), 2014.

- [38] Y.-X. Wang, B. Balle, and S. P. Kasiviswanathan. Subsampled renyi differential privacy and analytical moments accountant. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2019.
- [39] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017.
- [40] L. Xie, K. Lin, S. Wang, F. Wang, and J. Zhou. Differentially private generative adversarial network. *arXiv preprint arXiv:1802.06739*, 2018.
- [41] J. Yoon, J. Jordon, and M. van der Schaar. PATE-GAN: Generating synthetic data with differential privacy guarantees. In *International Conference on Learning Representations (ICLR)*, 2019.
- [42] J. Zhang, G. Cormode, C. M. Procopiuc, D. Srivastava, and X. Xiao. Privbayes: Private data release via bayesian networks. *ACM Transactions on Database Systems (TODS)*, 42(4), 2017.
- [43] X. Zhang, S. Ji, and T. Wang. Differentially private releasing via deep generative model (technical report). *arXiv preprint arXiv:1801.01594*, 2018.

Appendix

The appendix contains: the privacy analysis (§A), the algorithm pseudocode (§B), the details of experiment setup (§C), and additional results (§D). The source code will be released at the time of publication.

A Privacy Analysis

The privacy cost (ϵ) computation including: (i) bounding the privacy loss for our gradient sanitization mechanism using RDP; (ii) applying analytical moments accountant of subsampled RDP [38] for a tighter upper bound on the RDP parameters; (iii) tracking the overall privacy cost: multiplying the RDP orders by the number of training iterations and converting the resulting RDP orders to an (ϵ, δ) pair (Definition 3.2 [31]). We below present the theoretical results.

Theorem 4.1. Each generator update step satisfies $(\lambda, 2B\lambda/\sigma^2)$ -RDP where B is the batch size.

Proof. Let $f = \text{clip}(g_G^{\text{upstream}}, C)$, i.e., the clipped gradient before being sanitized. The sensitivity can be derived via the triangle inequality:

$$\Delta_2 f = \max_{S, S'} \|f(S) - f(S')\|_2 \leq 2C \quad (14)$$

with $C = 1$ in our case. Hence, we have $\mathcal{M}_{\sigma, C}$ is $(\lambda, 2\lambda/\sigma^2)$ -RDP.

Each generator update step (which operates on a batch of data) can be expressed as

$$\hat{g}_G = \frac{1}{B} \sum_{i=1}^B \mathcal{M}_{\sigma, C}(\nabla_{G(z_i)} \mathcal{L}_G(\theta_G)) \cdot J_{\theta_G} G(z_i; \theta_G) \quad (15)$$

This can be seen as a composition of B Gaussian mechanism. Concretely, we want to bound the Rényi divergence $D_\lambda(\hat{g}_G(S) \parallel \hat{g}_G(S'))$ with S, S' denoting the neighbouring datasets. We use the following properties of Rényi divergence [37]:

(i) Data-processing inequality : $D_\lambda(P_Y \parallel Q_Y) \leq D_\lambda(P_X \parallel Q_X)$ if the transition probabilities $A(Y|X)$ in the Markov chain $X \rightarrow Y$ is fixed.

(ii) Additivity : For arbitrary distributions P_1, \dots, P_N and Q_1, \dots, Q_N let $P^N = P_1 \times \dots \times P_N$ and $Q^N = Q_1 \times \dots \times Q_N$. Then $D_\lambda(P^N \parallel Q^N) = \sum_{n=1}^N D_\lambda(P_n \parallel Q_n)$

Let u and v denote the output distribution of the sanitization mechanism $\mathcal{M}_{\sigma, C}$ when applied on S and S' respectively, and h the post-processing function (i.e., multiplication with the local Jacobian). We have,

$$D_\lambda(\hat{g}_G(S), \hat{g}_G(S')) \leq D_\lambda(h_1(u_1) * \dots * h_B(u_B) \parallel h_1(v_1) * \dots * h_B(v_B)) \quad (16)$$

$$\leq D_\lambda((h_1(u_1), \dots, h_B(u_B)) \parallel (h_1(v_1), \dots, h_B(v_B))) \quad (17)$$

$$= \sum_b D_\lambda(h_b(u_b) \parallel h_b(v_b)) \quad (18)$$

$$\leq \sum_b D_\lambda(u_b \parallel v_b) \quad (19)$$

$$\leq B \cdot \max_b D_\lambda(u_b \parallel v_b) \quad (20)$$

$$\leq B \cdot 2\lambda/\sigma^2 \quad (21)$$

where (3)(4)(6) are based on the data-processing theorem; (5) follows from the additivity; and the last equation follows from the $(\lambda, 2\lambda/\sigma^2)$ -RDP of $\mathcal{M}_{\sigma, C}$. \square

Theorem A.1. (RDP for Subsampled Mechanisms [38]) Given a dataset containing n datapoints with domain \mathcal{X} and a randomized mechanism \mathcal{M} that takes an input from \mathcal{X}^m for $m \leq n$, let the randomized algorithm $\mathcal{M} \circ \text{subsample}$ be defined as: (i) **subsample**: subsample without replacement m datapoints of the dataset (with subsampling rate $\gamma = m/n$); (ii) apply \mathcal{M} : a randomized algorithm

taking the subsampled dataset as the input. For all integers $\lambda \geq 2$, if \mathcal{M} is $(\lambda, \epsilon(\lambda))$ -RDP, then $\mathcal{M} \circ \text{subsample}$ is $(\lambda, \epsilon'(\lambda))$ -RDP where

$$\begin{aligned} \epsilon'(\lambda) \leq & \frac{1}{\lambda-1} \log \left(1 + \gamma^2 \binom{\lambda}{2} \min \left\{ 4(e^{\epsilon(2)} - 1), e^{\epsilon(2)} \min \{ 2, (e^{\epsilon(\infty)} - 1)^2 \} \right\} \right. \\ & \left. + \sum_{j=3}^{\lambda} \gamma^j \binom{\lambda}{j} e^{(j-1)\epsilon(j)} \min \{ 2, (e^{\epsilon(\infty)} - 1)^j \} \right) \end{aligned}$$

In practice, we adopt the official implementation of [38]¹ for computing the accumulated privacy cost (i.e., tracking the RDP orders and converting RDP to (ϵ, δ) -DP).

B Algorithm

We present the pseudocode of our proposed method in Algorithm 1 (Centralized setup) and Algorithm 2 (Federated setup).

Algorithm 1: Centralized GS-WGAN Training

Input: Dataset S , subsampling rate γ , noise scale σ , warm-start iterations T_w , training iterations T , learning rates η_D and η_G , the number of discriminator iterations per generator iteration n_{dis} , batch size B

Output: Differentially Private generator G with parameters θ_G , total privacy cost ϵ

- 1 Subsample (without replacement) the dataset S into subsets $\{S_k\}_{k=1}^K$ with rate γ ($K = 1/\gamma$);
- 2 **for** k **in** $\{1, \dots, K\}$ **in parallel do**
- 3 Initialize non-private generator θ_G^k , discriminator θ_D^k **for** $step$ **in** $\{1, \dots, T_w\}$ **do**
- 4 **for** t **in** $\{1, \dots, n_{dis}\}$ **do**
- 5 Sample batch $\{x_i\}_{i=1}^B \subseteq S_k$;
- 6 Sample batch $\{z_i\}_{i=1}^B$ with $z_i \sim P_z$;
- 7 $\theta_D^k \leftarrow \theta_D^k - \eta_D \cdot \frac{1}{B} \sum_i \nabla_{\theta_D^k} \mathcal{L}_D(\theta_D^k; x_i, G(z_i; \theta_G^k))$;
- 8 **end**
- 9 $\theta_G^k \leftarrow \theta_G^k - \eta_G \cdot \frac{1}{B} \sum_i \nabla_{\theta_G^k} \mathcal{L}_G(\theta_G^k; G(z_i; \theta_G^k), \theta_D^k)$;
- 10 **end**
- 11 Initialize private generator θ_G ;
- 12 **for** $step$ **in** $\{1, \dots, T\}$ **do**
- 13 Sample subset index $k \sim \mathcal{U}[1, K]$;
- 14 **for** t **in** $\{1, \dots, n_{dis}\}$ **do**
- 15 Sample batch $\{x_i\}_{i=1}^B \subseteq S_k$;
- 16 Sample batch $\{z_i\}_{i=1}^B$ with $z_i \sim P_z$;
- 17 $\theta_D^k \leftarrow \theta_D^k - \eta_D \cdot \frac{1}{B} \sum_i \nabla_{\theta_D^k} \mathcal{L}_D(\theta_D^k; x_i, G(z_i; \theta_G))$;
- 18 **end**
- 19 $\theta_G \leftarrow \theta_G - \eta_G \cdot \frac{1}{B} \sum_i \mathcal{M}_{\sigma, C}(\theta_G; G(z_i; \theta_G), \theta_D^k) \cdot \mathbf{J}_{\theta_G} G(z_i; \theta_G)$;
- 20 Accumulate privacy cost ϵ ;
- 21 **end**
- 22 **end**
- 23 **return** Generator $G(\cdot; \theta_G)$, privacy cost ϵ

¹ <https://github.com/yuxiangw/autodp>

Algorithm 2: Federated (Decentralized) GS-WGAN Training

Input: Client index set $\{1, \dots, K\}$, noise scale σ , warm-start iterations T_w , training iterations T , learning rates η_D and η_G , the number of discriminator iterations per generator iteration n_{dis} , batch size B

Output: Differentially Private generator G with parameters θ_G , total privacy cost ε

```
1 for each client  $k$  in  $\{1, \dots, K\}$  in parallel do
2   | ClientWarmStart( $k$ )
3 end
4 Initialize private generator  $\theta_G$ ;
5 for  $step$  in  $\{1, \dots, T\}$  do
6   | Sample client index  $k \sim \mathcal{U}[1, K]$ ;
7   | for  $t$  in  $\{1, \dots, n_{dis}\}$  do
8     | Sample batch  $\{z_i\}_{i=1}^B$  with  $z_i \sim P_z$ ;
9     |  $\{\hat{g}_i^{up}\}_{i=1}^B \leftarrow \text{ClientUpdate}(k, G(z_i; \theta_G))$ 
10    | end
11    |  $\theta_G \leftarrow \theta_G - \eta_G \cdot \frac{1}{B} \sum_i \hat{g}_i^{up} \cdot \mathbf{J}_{\theta_G} G(z_i; \theta_G)$ ;
12    | Accumulate privacy cost  $\varepsilon$ ;
13  end
14 return Generator  $G(\cdot; \theta_G)$ , privacy cost  $\varepsilon$ 
```

```
16 Procedure ClientWarmStart( $k$ )
17   Get local dataset  $S_k$ ;
18   Initialize local generator  $\theta_G^k$ , discriminator  $\theta_D^k$ ;
19   for  $step$  in  $\{1, \dots, T_w\}$  do
20     | for  $t$  in  $\{1, \dots, n_{dis}\}$  do
21       | Sample batch  $\{x_i\}_{i=1}^B \subseteq S_k$ ;
22       | Sample batch  $\{z_i\}_{i=1}^B$  with  $z_i \sim P_z$ ;
23       |  $\theta_D^k \leftarrow \theta_D^k - \eta_D \cdot \frac{1}{B} \sum_i \nabla_{\theta_D^k} \mathcal{L}_D(\theta_D^k; x_i, G(z_i; \theta_G^k))$ ;
24     | end
25     |  $\theta_G^k \leftarrow \theta_G^k - \eta_G \cdot \frac{1}{B} \sum_i \nabla_{\theta_G^k} \mathcal{L}_G(\theta_G^k; G(z_i; \theta_G^k), \theta_D^k)$ ;
26   end
```

```
27 Procedure ClientUpdate( $k, G(z_i; \theta_G)$ )
28   Get local dataset  $S_k$ , local discriminator  $D(\cdot; \theta_D^k)$ ;
29   Sample batch  $\{x_i\}_{i=1}^B \subseteq S_k$ ;
30    $\theta_D^k \leftarrow \theta_D^k - \eta_D \cdot \frac{1}{B} \sum_i \nabla_{\theta_D^k} \mathcal{L}_D(\theta_D^k; x_i, G(z_i; \theta_G))$ ;
31   return  $\mathcal{M}_{\sigma, C}(\theta_G; G(z_i; \theta_G), \theta_D^k)$ 
```

C Experiment Setup

C.1 Hyperparameters

We adopt the hyperparameters setting in [17] for the GAN training, and list below the hyperparameters relevant for privacy computation.

Centralized Setting. We use by default a subsampling rate of $\gamma=1/1000$, noise scale $\sigma=1.07$, pretraining (warm-start) for 2K iterations and subsequently training for 20K iterations.

Federated Setting. We use by default a noise scale $\sigma=1.07$, pretraining (warm-start) for 2K iterations and subsequently training for 30K iterations.

C.2 Datasets

Centralized Setting. MNIST [21] and Fashion-MNIST [39] datasets contain 60K training images and 10K testing images. Each image has dimension 28×28 and belongs to one of the 10 classes.

Federated Setting. **Federated EMNIST** [9] dataset contains 28×28 gray-scale images of handwritten letters and numbers, grouped by user. The entire dataset contains 3400 users with 671,585 training examples and 77,483 testing examples. Following [3], the users are filtered by the prediction accuracy of a 36-class (10 numeric digits + 26 letters) CNN classifier. For evaluating the sample quality, we train GAN models on the users’ data which yields classification accuracy $\geq 93.9\%$ (866 users); For simulating the debugging task, we randomly choose 50% of the users and pre-process their data by flipping the pixel intensities. To mimic the real-world situation where the server is blind to the erroneous pre-processing, users with low classification accuracy $\leq 88.2\%$ are selected (2136 users) as they are suspected to be affected by erroneous flipping (with bug). Note that only a fraction of them is indeed affected by the bug (1720 with bug, 416 without bug). This has the realistic property that the client data is non-IID and poses additional difficulties in the GAN training.

C.3 Evaluation Metrics

In line with previous literature, we use Inception Score (IS) [23, 35] and Frechet Inception Distance (FID) [20] for measuring sample quality, and classification accuracy for evaluating the usefulness of generated samples. We present below a detailed explanation of the evaluation metrics we adopted in the experiments.

Inception Score (IS). Formally, the IS is defined as follows,

$$\text{IS} = \exp \left(\mathbb{E}_{\mathbf{x} \sim G(\mathbf{z})} D_{KL}(P(y|\mathbf{x}) \| P(y)) \right)$$

which corresponds to exponential of the KL divergence between the conditional class $P(y|\mathbf{x})$ and the marginal class distribution $P(y)$, where both $P(y|\mathbf{x})$ and $P(y)$ are measured by the output distribution of a pre-trained classifier when passing the generated samples as input. Intuitively, the IS should exhibit a high value if $P(y|\mathbf{x})$ has low entropy (i.e., the generated images are sharp and contain clear objects) and $P(y)$ is of high entropy (i.e., the generated samples have a high diversity covering all the different classes). In our experiments, we use pre-trained classifiers on the real datasets (with test accuracy equals to 99.25%, 93.75%, 92.16% on the MNIST, Fashion-MNIST and Federated EMNIST dataset respectively)² for computing the IS.

Frechet Inception Distance (FID). The FID is formularized as follows,

$$\text{FID} = \|\mu_r - \mu_g\|^2 + \text{tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2})$$

where $\mathbf{x}_r \sim \mathcal{N}(\mu_r, \Sigma_r)$ and $\mathbf{x}_g \sim \mathcal{N}(\mu_g, \Sigma_g)$ are the 2048-dimensional activations of the Inception-v3 pool3 layer for real and generated samples respectively. A lower FID value indicates a smaller discrepancy between the real and generated samples, which corresponds to a better sample quality and diversity. Following previous works³, we rescale the images and convert them to RGB by repeating the grayscale channel three times before inputting them to the Inception network.

Classification Accuracy. We consider the following classification models in our experiments: Multi-layer Perceptron (MLP), Convolutional Neural Network (CNN), AdaBoost (adaboost), Bagging (bagging), Bernoulli Naive Bayes (bernoulli nb), Decision tree (decision tree), Gaussian Naive Bayes (gaussian nb), Gradient Boosting (gbm), Linear Discriminant Analysis (lda), Linear Support Vector Machine (linear svc), Logistic Regression (logistic reg), Random Forest (random forest), and XGBoost (xgboost). For implementing the CNN model, we use two hidden layers (with dropout) each containing 32 and 64 kernels and apply ReLU as the activation function. For implementing the MLP, we use one hidden layer with 100 neurons and set ReLU as the activation function. All the other classification models are implemented using the default hyperparameters supplied by the scikit-learn [32] package.

C.4 Baseline Methods

We present more details about the implementation of the baseline methods. In particular, we provide the default value of the privacy hyperparameters below.

² https://github.com/ChunyuanyanLI/MNIST_Inception_Score

³ https://github.com/google/compare_gan

DP-Merf (AE)⁴ We use as default a batch size=500 ($\gamma=1/120$), noise scale $\sigma=0.588$, training iteration=600 (epoch=5) for implementing DP-Merf, and batch size=500, noise scale $\sigma=0.686$, training iteration=2040 (epoch=17) for implementing DP-Merf AE.

DP-SGD GAN⁵ We set the default hyper-parameters as follows: gradient clipping bound $C=1.1$, noise scale $\sigma=2.1$, batch size=600, training iterations=30K.

G-PATE We use 2000 teacher discriminators with batch size of 30 and set noise scales $\sigma_1=600$ and $\sigma_2=100$, consensus threshold $T=0.5$. A random projection with projection dimension=10 is applied.

PATE-GAN⁶ When extending PATE-GAN to high-dimensional image datasets, we observe that after a few iterations, the generated samples are classified as fake by all teacher discriminators and the learning signals (gradients) for student discriminator and the generator vanish. Consequently, the training stuck at the early stage where the losses remain unchanged and no progress can be observed. While this issue is well resolved by careful design of the prior distribution, as reported in the original paper, we find that this technique has a limited effect when applied to the high-dimensional image dataset. In addition, we make the following attempts to address this issue: (i) changing the network initialization (ii) increasing (or decreasing) the network capacity of the student discriminator, the teacher discriminators, and the generator (iii) increasing the number of iterations for updating the student discriminator and/or the generator. Despite some progress in preserving the gradients for larger iterations, none of the above attempts successfully eliminate the issue, as the training inevitably gets stuck within 1K iterations.

D Additional Results

Effects of gradient clipping. We show in Figure A1 the gradient norm distribution before and after gradient clipping. The clipping bound is set to be 1.1 for DP-SGD and 1 for our method. In contrast to DP-SGD, the clipping operation distorts less information in our framework, witnessed by a much smaller difference in the average gradient norm before and after the clipping. Moreover, the gradients used in our method exhibit much less variance both before and after the clipping compared with DP-SGD.

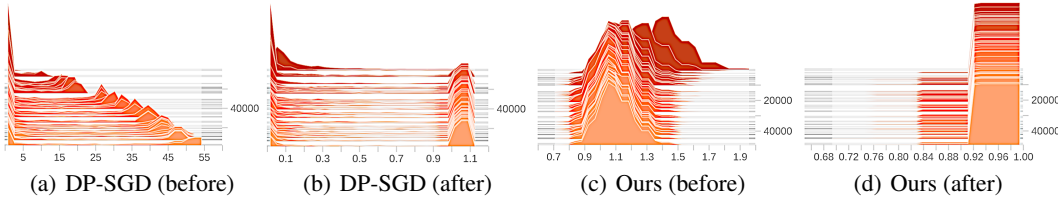


Figure A1: Effects of gradient clipping.

Comparison to Baselines. We provide the detailed quantitative results in Table A1 and A2, which are supplementary to Table 1 in the main paper. We show in parentheses the calibrated accuracy, i.e., the absolute accuracy of each classifier trained on generated data divided by the accuracy when trained on real data. The results are averaged over five runs.

Privacy-utility Curves. We show in Figure A2 the privacy-utility curves of different methods when applied to the Fashion-MNIST dataset. We evaluate over three runs and show the corresponding mean and standard deviation. Similar to the results shown in Figure 4 in the main paper, our method achieves a consistent improvement over prior methods across a broad range of privacy budget ϵ .

⁴ <https://github.com/frhrdr/Differentially-Private-Mean-Embeddings-with-Random-Features-for-Synthetic-Data-Generation>

⁵ <https://github.com/reihaneh-torkzadehmahani/DP-CGAN>

⁶ <https://bitbucket.org/mvdschaar/mlforhealthlabpub/src/2534877d99c8fdf19cbade16057990171e249ef3/alg/pategan/>

| | Real | GAN (non-private) | G-PATE | DP-SGD GAN | DP-Merf | DP-Merf AE | Ours |
|---------------|------|-------------------|-------------|------------|-------------|------------|-------------|
| MLP | 0.98 | 0.84 (85%) | 0.25 (26%) | 0.60 (61%) | 0.63 (64%) | 0.54 (55%) | 0.79 (81%) |
| CNN | 0.99 | 0.84 (85%) | 0.51 (52%) | 0.64 (65%) | 0.63 (64%) | 0.68 (69%) | 0.80 (81%) |
| adaboost | 0.73 | 0.28 (39%) | 0.11 (16%) | 0.32 (44%) | 0.38 (52%) | 0.21 (29%) | 0.21 (29%) |
| bagging | 0.93 | 0.46 (49%) | 0.36 (38%) | 0.44 (47%) | 0.43 (46%) | 0.33 (35%) | 0.45 (48%) |
| bernoulli nb | 0.84 | 0.80 (95%) | 0.71 (84%) | 0.62 (74%) | 0.76 (90%) | 0.50 (60%) | 0.77 (92%) |
| decision tree | 0.88 | 0.40 (45%) | 0.13 (14%) | 0.36 (41%) | 0.29 (33%) | 0.27 (31%) | 0.35 (40%) |
| gaussian nb | 0.56 | 0.71 (126%) | 0.61 (110%) | 0.37 (66%) | 0.57 (102%) | 0.17 (30%) | 0.64 (114%) |
| gbm | 0.91 | 0.50 (55%) | 0.11 (12%) | 0.45 (49%) | 0.36 (40%) | 0.20 (22%) | 0.39 (43%) |
| lda | 0.88 | 0.84 (95%) | 0.60 (68%) | 0.59 (67%) | 0.72 (82%) | 0.55 (63%) | 0.78 (89%) |
| linear svc | 0.92 | 0.81 (88%) | 0.24 (26%) | 0.56 (61%) | 0.58 (63%) | 0.43 (47%) | 0.76 (83%) |
| logistic reg | 0.93 | 0.83 (90%) | 0.26 (28%) | 0.60 (65%) | 0.66 (71%) | 0.55 (59%) | 0.79 (85%) |
| random forest | 0.97 | 0.39 (41%) | 0.33 (34%) | 0.63 (65%) | 0.66 (68%) | 0.45 (46%) | 0.52 (54%) |
| xgboost | 0.91 | 0.44 (49%) | 0.15 (16%) | 0.60 (66%) | 0.70 (77%) | 0.54 (59%) | 0.50 (55%) |
| Average | 0.88 | 0.63 (71%) | 0.34 (40%) | 0.52 (59%) | 0.57 (66%) | 0.42 (47%) | 0.60 (69%) |

Table A1: Classification accuracy on MNIST ($\epsilon = 10$, $\delta = 10^{-5}$).

| | Real | GAN (non-private) | G-PATE | DP-SGD GAN | DP-Merf | DP-Merf AE | Ours |
|---------------|------|-------------------|------------|------------|-------------|------------|------------|
| MLP | 0.88 | 0.77 (88%) | 0.30 (34%) | 0.50 (57%) | 0.56 (64%) | 0.56 (64%) | 0.65 (74%) |
| CNN | 0.91 | 0.73 (80%) | 0.50 (54%) | 0.46 (51%) | 0.54 (59%) | 0.62 (68%) | 0.64 (70%) |
| adaboost | 0.56 | 0.41 (74%) | 0.42 (75%) | 0.21 (38%) | 0.33 (59%) | 0.26 (46%) | 0.25 (45%) |
| bagging | 0.84 | 0.57 (68%) | 0.38 (45%) | 0.32 (38%) | 0.40 (47%) | 0.45 (54%) | 0.47 (56%) |
| bernoulli nb | 0.65 | 0.59 (91%) | 0.57 (88%) | 0.50 (77%) | 0.62 (95%) | 0.54 (83%) | 0.55 (85%) |
| decision tree | 0.79 | 0.53 (67%) | 0.24 (30%) | 0.33 (42%) | 0.25 (32%) | 0.36 (46%) | 0.40 (51%) |
| gaussian nb | 0.59 | 0.55 (93%) | 0.57 (97%) | 0.28 (47%) | 0.59 (100%) | 0.12 (20%) | 0.48 (81%) |
| gbm | 0.83 | 0.44 (53%) | 0.25 (30%) | 0.38 (46%) | 0.27 (33%) | 0.30 (36%) | 0.38 (46%) |
| lda | 0.80 | 0.77 (96%) | 0.55 (69%) | 0.55 (69%) | 0.67 (84%) | 0.65 (81%) | 0.67 (84%) |
| linear svc | 0.84 | 0.77 (91%) | 0.30 (36%) | 0.39 (46%) | 0.46 (55%) | 0.40 (48%) | 0.65 (77%) |
| logistic reg | 0.84 | 0.76 (90%) | 0.35 (42%) | 0.51 (61%) | 0.59 (70%) | 0.50 (60%) | 0.68 (81%) |
| random forest | 0.88 | 0.69 (78%) | 0.33 (37%) | 0.51 (58%) | 0.61 (69%) | 0.55 (63%) | 0.54 (61%) |
| xgboost | 0.83 | 0.65 (78%) | 0.49 (59%) | 0.52 (63%) | 0.62 (75%) | 0.55 (66%) | 0.47 (57%) |
| Average | 0.79 | 0.61 (77%) | 0.40 (54%) | 0.42 (53%) | 0.50 (65%) | 0.45 (56%) | 0.53 (67%) |

Table A2: Classification accuracy on Fashion-MNIST ($\epsilon = 10$, $\delta = 10^{-5}$).

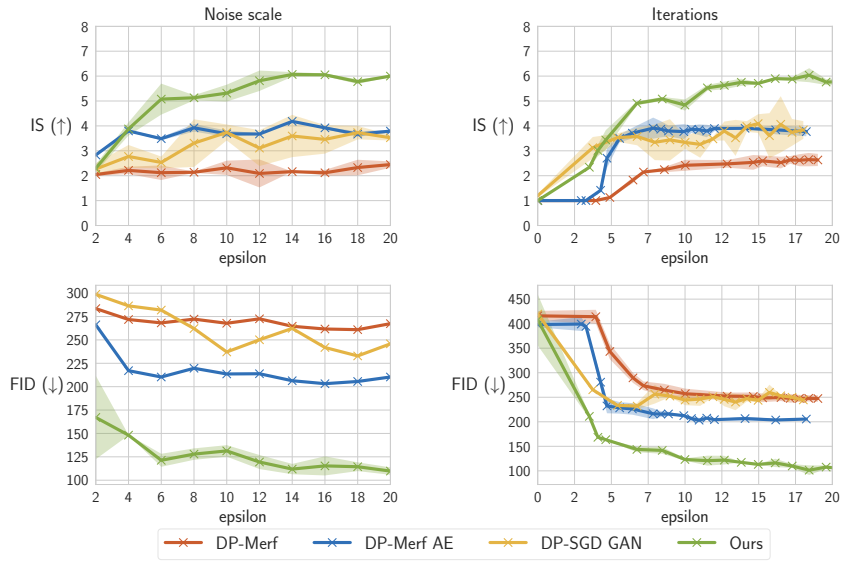


Figure A2: Privacy-utility trade-off on Fashion-MNIST with $\delta = 10^{-5}$. (Left: Effects of noise scale. Right: Effects of Iterations.)