

---

# Shielding Federated Learning: Aligned Dual Gradient Pruning Against Gradient Leakage

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Federated learning (FL) is a distributed learning framework that claims to protect  
2 user privacy. However, gradient inversion attacks (GIAs) reveal severe privacy  
3 threats to FL, which can recover the users' training data from outsourced gradients.  
4 Existing defense methods adopt different techniques, e.g., differential privacy,  
5 cryptography, and gradient perturbation, to against the GIAs. Nevertheless, all  
6 current state-of-the-art defense methods suffer from a trade-off between privacy,  
7 utility, and efficiency in FL. To address the weaknesses of existing solutions, we  
8 propose a novel defense method, Aligned Dual Gradient Pruning (ADGP), based  
9 on gradient sparsification, which can improve communication efficiency while  
10 preserving the utility and privacy of the federated training. Specifically, ADGP  
11 slightly changes gradient sparsification with a stronger privacy guarantee. Through  
12 primary gradient parameter selection strategies during training, ADGP can also  
13 significantly improve communication efficiency with a theoretical analysis of its  
14 convergence and generalization. Our extensive experiments show that ADGP can  
15 effectively defend against the most powerful GIAs and significantly reduce the  
16 communication overhead without sacrificing the model's utility.

## 17 1 Introduction

18 Federated learning (FL) [1] is a distributed learning framework, where multiple users train and send  
19 their gradients of the local models to the server without sharing their local data [1, 2, 3]. FL claims  
20 to protect the users' training data since the users do not need to share local data with the server  
21 directly. However, recent studies reveal that gradients can be used to recover the original training  
22 data information via gradient inversion attacks (GIAs) [4, 5]. To against GIAs, a large number of  
23 studies have been proposed, where they leverage the advanced privacy protection techniques, such as  
24 differential privacy (DP) [6], cryptography [7, 8, 9] and gradient perturbation [10, 11, 12]. However,  
25 none of the existing defense methods could take care of all privacy, utility, and efficiency difficulties  
26 in FL. For example, DP and cryptography-based methods could effectively defend GIAs, but sacrifice  
27 either the utility or efficiency respectively [6, 7, 8, 9]. In order to achieve better utility and efficiency  
28 in FL, perturbation-based methods design various gradient perturbations [10, 11, 12], but all existing  
29 perturbation-based methods could only defend one or two kinds of GIAs in practice.

30 For example, recent perturbation-based defense methods (*i.e.*, Precode [12], Soteria [10], and  
31 ATS [11]) can effectively defend against optimization-based GIAs [5, 13, 14, 15], but fail to work  
32 against the active GIAs [16, 17, 18]. On the contrary, the classic Top- $k$  based gradient sparsification  
33 method [19, 20] is generally considered as a bad privacy protection solution on optimization-based  
34 GIAs, but in fact performs much better than recent defense methods under the active attack from  
35 our experiments as shown in Table 2. The new findings inspire us to seek for a more practical  
36 perturbation-based defense against both optimization-based and active GIAs.

Table 1: Comparison of our method with existing privacy-preserving FL methods. Note:  $\checkmark$  represents the scheme has a high guarantee for the property, while  $\times$  represents otherwise.

Defense	Privacy					Utility	Efficiency
	Analytical attack	Optimization attack		Active server attack			
	R-GAP [22]	DLG [4]	IVG [5]	Curious [16]	Rob [21]		
Precode [12]	$\checkmark$	$\checkmark$	$\checkmark$	$\times$	$\times$	$\times$	$\times$
ATS [11]	$\times$	$\checkmark$	$\times$	$\times$	$\times$	$\checkmark$	$\times$
Soteria [10]	$\checkmark$	$\checkmark$	$\checkmark$	$\times$	$\times$	$\checkmark$	$\times$
DP [23]	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\times$	$\times$
Top- $k$ [19]	$\checkmark$	$\checkmark$	$\times$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
ADGP (Ours)	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$

37 In this paper, we propose a new gradient pruning based method, **Aligned Dual Gradient Pruning**  
 38 (ADGP). Specifically, ADGP consists of two components: dual gradient pruning (DGP) and gradient  
 39 location bounding. Dual gradient pruning is a novel gradient sparsification technique, which removes  
 40 top- $k_1$  largest gradient parameters and the bottom- $k_2$  smallest gradient parameters from the local  
 41 model. DGP leads to a strong privacy protection against both optimization-based and active GIAs.  
 42 To further reduce the expensive download costs caused by the asymmetric gradient selection among  
 43 different users, we propose gradient location bounding strategy to make the aggregated global gradient  
 44 stay in the same sparsified region. In summary, ADGP achieves a better utility and privacy trade-off,  
 45 increases FL system efficiency, and is robust against active attacks.

46 Furthermore, we give the theoretical analysis of ADGP, which proves the reconstruction error is  
 47 proportional to gradient distance. So removing larger gradient parameters could enlarge the gradient  
 48 distance resulting in a low reconstruction error. However, removing larger gradients will significantly  
 49 impact the model’s utility. Thus, to improve the sparsification ratio, which is essential to robustness  
 50 against active attack [21, 16], we also remove the model parameters with smaller gradients. In such a  
 51 way, our method could significantly mitigate GIAs without affecting the model’s utility.

52 We conduct extensive experiments over multiple datasets and models to evaluate our method. The  
 53 quantitative and visualized results show that our design can effectively make recovered images  
 54 recognizable under different attacks, and reduce nearly half of the communication costs. Our  
 55 contributions are as follows: 1) We revisit gradient sparsification to show its potential on mitigating  
 56 GIA; 2) We propose an improved gradient pruning strategy to provide sufficient privacy guarantee  
 57 while balancing the model accuracy and the system efficiency; 3) We conduct extensive experiments  
 58 to show that our design outperforms perturbation-based defense methods *w.r.t* privacy protection,  
 59 model accuracy, and system efficiency.

## 60 2 Related work

61 Federated learning [1, 3] is considered to be a privacy-preserving framework for distributed machine  
 62 learning as the training data is not directly outsourced. However, the emerging of GIAs [4, 5, 21, 16,  
 63 22, 24, 25, 26, 27] shatters this conception. It has been proved that the attacker (*e.g.*, a curious server)  
 64 can easily recover the private training data from gradient information to a great extent. The privacy  
 65 guarantee of federated learning urgently needs to be strengthened.

66 **Cryptographic Defense.** Traditionally, there are two approaches to construct privacy-preserving  
 67 federated learning: using DP to disturb gradients [6, 23, 28, 29, 30] or using cryptographic tools to  
 68 perform secure aggregation [7, 8, 9, 31, 32]. DP [6] is a popular and effective privacy protection  
 69 mechanism by adding random noise to the raw data, but it is well known that the noises intro-  
 70 duced by DP can greatly degrade the model accuracy when meaningful privacy is enforced [33].  
 71 Cryptographic-based secure aggregation can guarantee both privacy and accuracy simultaneously, but  
 72 incurs expensive computation and communication costs [34]. Using the shuffle model [35, 36] can  
 73 only provide anonymity. Moreover, it totally changes the system model of FL since an additional  
 74 semi-trusted third party is introduced to work cooperatively with the server.

75 **Gradient Perturbation Defense.** Recently, researchers have begun to explore the possibility of  
 76 constructing new gradient perturbation mechanisms to better balancing privacy and accuracy. Sun *et*  
 77 *al.* [10] proposed Soteria, a scheme that perturbs the representation of training data by pruning the

78 gradients of a single fully connected layer. Gao *et al.* [11] proposed ATS, a training data augmentation  
 79 policy by transforming original sensitive images into alternative inputs, to reduce the visibility of  
 80 reconstructed images. Scheliga *et al.* [12] presented Precode to extend the model architecture by using  
 81 variational bottleneck (VB) [37] to prevent attackers from obtaining optimal solutions to reconstructed  
 82 data. These defenses work well against GIAs in the semi-honest setting [4, 5, 38, 13], but fail to  
 83 protect privacy when an active server modifies the model to launch GIAs [21, 16]. Moreover, these  
 84 works suffer from high computation costs or huge communication burden.

85 **Gradient Sparsification Defense.** From an independent research domain, gradient sparsification  
 86 has been commonly used for saving communication bandwidth. The most common sparsification  
 87 strategy is Top- $k$  selection, which selects top  $k$  gradient parameters with the largest absolute values  
 88 [19, 20]. It has been widely proved that gradient sparsification provides very limited privacy protection  
 89 ability [4, 10, 11, 12, 39] unless a high pruning ratio (*e.g.*, removing 99% of the gradients) is used  
 90 at the cost of 10% accuracy drop [39]. However, we emphasize that this is misunderstood as they  
 91 only consider the Top- $k$  sparsification strategy that has never received an in-depth investigation in the  
 92 field of security. It is originally designed for improving system efficiency, thus a direct application  
 93 inherently suffers from many weaknesses. As shown in Section 4, a slight modification can unleash  
 94 the potential of gradient sparsification to provide a strong privacy guarantee.

### 95 3 Threat Model and Attacks

96 In this work, we consider a strong threat scenario, where an active server, after receiving gradients  
 97 from users, tries to reconstruct the local training data and is motivated to modify model parameters  
 98 in each iteration to strengthen the attack performance. Note that the server also wants to obtain a  
 99 high-quality global model with high accuracy. More specifically, we consider the following three  
 100 kinds of GIAs:

101 **Analytical attack.** Analytical attack exploits the structure of the gradients to recover the input,  
 102 such as using gradient bias terms [40]. Recently proposed R-GAP attack [22] exploits the recursive  
 103 relationship between gradient layers to solve the input. An effective analytical attack depends on the  
 104 specific structure and parameters of gradients.

105 **Optimization attack.** Optimization attack is first proposed using L-BFGS optimizer to solve  
 106  $\min \|\frac{\partial l(\mathbf{x}, \mathbf{y})}{\partial \mathbf{W}} - \frac{\partial l(\mathbf{x}^*, \mathbf{y}^*)}{\partial \mathbf{W}}\|_2^2$  and gets dummy data  $\mathbf{x}^*$  and dummy label  $\mathbf{y}^*$ , where  $\mathbf{y}$  is the label of  
 107  $\mathbf{x}$  [4]. The state-of-the-art optimization attack method IVG [5] uses Adam to optimize the cosine  
 108 distance and has been widely used to evaluate defense works [10, 11, 12].

109 Despite different optimizers can be used to achieve better attack quality [5, 13, 14, 15], the existing  
 110 attacks are all measured by the distance between the generated gradients  $\nabla \mathbf{W}^*$  and the original  
 111 gradients  $\nabla \mathbf{W}$ . We therefore propose a general definition for optimization attack to better evaluate  
 112 its performance. As shown in Definition 1, a smaller  $\varepsilon$  indicates a stronger optimization attack.

113 **Definition 1.** An optimization attack is a  $(\varepsilon, \delta)$ -attack, if it satisfies:

$$\mathbb{P}(\mathbb{E}(\mathcal{D}(\nabla \mathbf{W}, \nabla \mathbf{W}^*)) \leq \varepsilon) \geq 1 - \delta. \quad (1)$$

114 where  $\mathbb{P}$  represents the probability,  $\mathbb{E}$  represents the expectation,  $\mathcal{D}$  is the distance function commonly  
 115 instantiated with Euclidean or cosine distance.

116 **Active server attack.** In this kind of attack, the server can actively modify the global model to  
 117 realize a better attack result rather than honestly executing the protocols [16, 17, 18]. Recently  
 118 proposed Rob attack [21] adds imprint modules to the model and uses the difference between the  
 119 gradient parameters in adjacent rows of the imprint module to recover the data, achieving the best  
 120 attack effect in the literature.

## 121 4 Aligned Dual Gradient Pruning

### 122 4.1 Analysis of Gradient Sparsification

123 We owe the failure of common Top- $k$  gradient sparsification methods to two reasons: 1) the distance  
 124 between the Top- $k$  sparsified (*i.e.*, perturbed) gradient  $\mathbf{g}$  and the real gradient  $\nabla \mathbf{W}$  is small; and 2)  
 125 large gradient parameters in  $\nabla \mathbf{W}$  also reveal label information about user data.

126 To explain the first reason, we investigate the relationship between the reconstruction error of user  
 127 data and distance of perturbed gradient v.s. real gradient, as shown in Proposition 1.

128 **Proposition 1.** For any given input  $\mathbf{x}$  and shared model  $\mathbf{W}$ , the distance between the recovered data  
 129  $\mathbf{x}'$  and the real data  $\mathbf{x}$  is bounded by:

$$\|\mathbf{x} - \mathbf{x}'\|_2 \geq \frac{\|\nabla \mathbf{W} - \mathbf{g}\|_2}{\|\partial \varphi(\mathbf{x}, \mathbf{W}) / \partial \mathbf{x}\|_2}, \quad (2)$$

130 where  $\varphi$  is the mapping from  $\mathbf{x}$  to  $\nabla \mathbf{W}$ , i.e., the reconstruction quality is limited by  $\|\nabla \mathbf{W} - \mathbf{g}\|_2$ .

131 The proof of the above proposition is moved to the supplementary due to space limit (the same  
 132 hereinafter). From this Proposition, it is clear that the reconstruction error is proportional to the  
 133 gradient distance  $\|\nabla \mathbf{W} - \mathbf{g}\|_2$ , i.e., effective defense methods should enlarge the gradient distance  
 134 as much as possible. However, for the Top- $k$  based gradient sparsification [19, 20], the  $k$  largest  
 135 parameters are retained, making the gradient distance small by nature.

136 To explain the second reason, we consider a  $L$ -layer perceptron model trained with cross-entropy  
 137 loss for classification. Let a column vector  $\mathbf{r} = [r_1, r_2, \dots, r_n]$  be the logits (the output of the  
 138  $L$ -th linear layer) that input to the softmax layer, the confidence score probability vector is thus  
 139  $\left[ \frac{e^{r_1}}{\sum_j e^{r_j}}, \frac{e^{r_2}}{\sum_j e^{r_j}}, \dots, \frac{e^{r_n}}{\sum_j e^{r_j}} \right]$  and the succinct form of the cross-entropy loss becomes  $\ell(\mathbf{x}, y) =$   
 140  $-\log\left(\frac{e^{r_y}}{\sum_j e^{r_j}}\right)$ . Focus on the  $L$ -th layer  $\mathbf{W}^L \mathbf{x} + \mathbf{b}^L = \mathbf{r}$ , it is easy to find

$$\frac{\partial \ell(\mathbf{x}, y)}{\partial b_i} = \frac{\partial \ell(\mathbf{x}, y)}{\partial r_i} \cdot \frac{\partial r_i}{\partial b_i} = \frac{\partial \ell(\mathbf{x}, y)}{\partial r_i} = \frac{e^{r_i}}{\sum_j e^{r_j}} - \mathbb{I}_{i=y}, \quad (3)$$

141 and

$$\nabla \mathbf{W}^L = \frac{\partial \ell(\mathbf{x}, y)}{\partial \mathbf{r}} \cdot \mathbf{x}^T = \left[ \frac{\partial \ell(\mathbf{x}, y)}{\partial r_1}, \dots, \frac{\partial \ell(\mathbf{x}, y)}{\partial r_n} \right] \cdot \mathbf{x}^T. \quad (4)$$

142 For a given  $\mathbf{x}$  (and so  $\mathbf{x}^T$  is fixed), the magnitude of certain elements of the gradient matrix  $\nabla \mathbf{W}^L$   
 143 (i.e., the  $i$ -th row) is particularly large if  $i$  is the true label of the training data  $\mathbf{x}$  due to reason that  
 144  $\left| \frac{\partial \ell(\mathbf{x}, y)}{\partial r_i} \right| = \sum_{j \neq i} \left| \frac{\partial \ell(\mathbf{x}, y)}{\partial r_j} \right|$ .

145 To summarize, due to the above two reasons, we conclude that common Top- $k$  gradient sparsification  
 146 cannot provide sufficient protection for user data against passive optimization attacks. From another  
 147 point of view, a sufficient gradient sparsification ratio also plays an important role in defending against  
 148 active server attacks. As mentioned in Section 3, active attackers can exploit the correspondence of  
 149 partial gradient parameters to recover the real data. So, the gradient sparsity rate will directly destroy  
 150 the relationship among gradient parameters constructed by the active attacker. Intuitively, the higher  
 151 the sparsity rate, the more severe the impact. As will be validated in Section 6, a higher sparsity rate  
 152 can prevent the attacker from obtaining useful gradient information.

## 153 4.2 Dual Gradient Pruning

154 Generally speaking, large gradient parameters of local model need to be removed to make the  
 155 gradient difference larger, but the difference should also be appropriately bounded to maintain high  
 156 model accuracy. Moreover, it is also necessary to delete small gradient parameters to achieve a  
 157 high sparsification ratio. With these observations, we propose dual gradients pruning (DGP), a new  
 158 parameter selection strategy for gradient sparsification.

159 The users first sort the absolute values of all  $\text{Size}(\nabla \mathbf{W})$  local gradient parameters in the descending  
 160 order. Let  $\mathcal{T}_{k_1}(\nabla \mathbf{W})$  represent the set of top- $k_1$  elements of  $\nabla \mathbf{W}$ ,  $\mathcal{B}_{k_2}(\nabla \mathbf{W})$  represent the set of  
 161 its bottom- $k_2$  elements. Then the users remove  $\mathcal{T}_{k_1}(\nabla \mathbf{W})$  and  $\mathcal{B}_{k_2}(\nabla \mathbf{W})$  from  $\nabla \mathbf{W}$  for gradient  
 162 sparsification. Note that we set  $p = k_2/k_1$  as a hyperparameter to regulate the trade-off between  
 163 privacy and accuracy. Clearly, even with a fixed value  $p$ , different user will have different sets of  
 164  $\mathcal{T}_{k_1}(\cdot)$  and  $\mathcal{B}_{k_2}(\cdot)$  because their respective local models could be different from each other.

165 We emphasize that although such dual gradients pruning strategy is very simple, it can significantly  
 166 mitigate GIAs without affecting the model accuracy. A rigorous security proof is shown in Section 5,  
 167 and experimental results can be found in Section 6.

---

**Algorithm 1:** Aligned Dual Gradient Pruning (ADGP)

---

**Input** : Original gradient matrix  $\nabla\mathbf{W}$ , location binary matrix  $\mathcal{I}$ , values of  $k_1$  and  $k$   
**Output** : Sparsified gradient matrix  $\mathbf{g} = \{\mathbf{g}^1, \dots, \mathbf{g}^L\}$   
**for**  $l \leftarrow 1$  to  $L$  **do**  
    Remove parameters in  $\mathcal{T}_{k_1}(\nabla\mathbf{W}^l)$  from  $\nabla\mathbf{W}^l$   
    Keep parameters in  $\nabla\mathbf{W}^l$  when location is in  $\mathcal{I}$ , and discard all other parameters  
    Upload  $\mathbf{g}^l = \mathcal{T}_k(\nabla\mathbf{W}^l)$  to the server

---

---

**Algorithm 2:** A Complete Illustration of our Defense

---

**Input** : Initial global model  $\mathbf{W}^0$ , value  $k$  and  $k_1$ , total rounds  $T$ , total users  $N$   
**Output** : Shared global model  $\mathbf{W}^T$   
Set  $\mathbf{e}^0 = 0$   
**for**  $t \leftarrow 0$  to  $T - 1$  **do**  
    Randomly select a user to broadcast the location matrix  $\mathcal{I}^t$  of its parameter set  $\mathcal{T}_{2k}$   
    **for**  $i \leftarrow 1$  to  $N$  **do**  
         $\mathbf{P}_i^t = \nabla\mathbf{W}_i^t + \mathbf{e}_i^t$   
         $\mathbf{g}_i^t = \text{ADGP}(k_1, k, \mathbf{P}_i^t, \mathcal{I}^t)$   
         $\mathbf{e}_i^{t+1} = \mathbf{P}_i^t - \mathbf{g}_i^t$   
    Server side aggregation:  
     $\mathbf{W}^{t+1} = \mathbf{W}^t - \eta \frac{\sum_{i=1}^N \mathbf{g}_i^t}{N}$

---

### 168 4.3 A Complete Illustration of Our Method

169 Although dual gradient pruning provides a sufficient privacy guarantee as well as reduces upload cost  
170 of users, users' download costs could still be expensive. This is because different users have different  
171 sets of  $\mathcal{T}_{k_1}(\cdot)$  and  $\mathcal{B}_{k_2}(\cdot)$  when sparsifying their own local gradients, which will ultimately make the  
172 global gradient become dense after aggregation.

173 We thus propose aligned dual gradient pruning (ADGP), an improved scheme to align the selected  
174 gradients across different users. Similar to DGP, for best privacy, each user will still firstly identify  
175 his top- $k_1$  gradients location set  $\mathcal{T}_{k_1}$ . Different from DGP, ADGP also wants to save users' download  
176 cost by ensuring that all users' uploaded sparsified gradients reside in the same location set. This is  
177 achieved by randomly selecting a user, who identifies a top- $2k$  ( $k_1 < k$ ) location set  $\mathcal{T}_{2k}$  (represented  
178 with a binary matrix  $\mathcal{I}$ ) and broadcasts  $\mathcal{I}$  to all other users. Note that  $\mathcal{T}_{k_1} \subset \mathcal{T}_{2k}$  is not necessary  
179 true. Upon receiving  $\mathcal{I}$ , each user first discards gradient parameters in  $\mathcal{T}_{k_1}$  and then only transmits  
180 the  $k$  largest gradient parameters whose locations belong to  $\mathcal{I}$ . After aggregation at the server side,  
181 users only need to download global gradients' parameters associated with  $\mathcal{I}$ . A detailed illustration  
182 of ADGP is shown in Algorithm 1.

183 For ADGP pruning, in each FL iteration round, all gradient parameters whose locations are outside of  
184  $\mathcal{I}$  will not participate the current round global model aggregation. In the extreme case,  $\mathcal{I}$  can remain  
185 static for all iteration rounds and the local accumulated error (accumulated unused local gradient  
186 parameters) becomes large, thus hindering global model convergence. To reduce this negative impact  
187 and increase convergence speed, we design an error feedback mechanism. In particular, at the iteration  
188 round  $t$ , after user  $i$  obtaining his local gradient  $\nabla\mathbf{W}_i^t$ , he will combine  $\nabla\mathbf{W}_i^t$  with an error term  
189 accumulated in the previous  $(t - 1)$  rounds before performing the ADGP sparsification pruning. A  
190 complete illustration of our method is shown in Algorithm 2.

## 191 5 Theoretical Analysis

192 This section presents the security analysis with regard to passive attacks (i.e., analytical and opti-  
193 mization attacks presented in Section 3), as well as the generalization and convergence analyses  
194 of the proposed ADGP algorithm. Following the literature studies in [41, 42], for a given  $L$ -layer  
195 centralized model, we model the first  $(L - 1)$  layers as a robust feature extractor of any input sample.

196 Thus, the function of this model is characterized by  $f(x|\mathbf{W}) = \mathbf{W}x + \mathbf{b}$ , and the optimization  
 197 objective is the loss  $\ell(\mathbf{x}, y)$  (such as cross-entropy or L2 loss). To facilitate analyses and following  
 198 literature studies [19, 41, 43, 44], the assumptions about the smoothness of DGP, ADGP and  $l$ , as  
 199 well as the variance of the stochastic gradient are employed.

**Assumption 1.** *The pruning mechanisms DGP( $k_1, k_2, \nabla \mathbf{W}^t$ ) and ADGP( $k_1, k, \nabla \mathbf{W}^t, \mathcal{I}^t$ ) are both bi-Lipschitz, so the following conditions hold:*

$$\begin{aligned} \|\nabla \mathbf{W} - \text{DGP}(k_1, k_2, \nabla \mathbf{W})\|_2^2 &= \|\text{DGP}(0, 0, \nabla \mathbf{W}) - \text{DGP}(k_1, k_2, \nabla \mathbf{W})\|_2^2 \geq \gamma_1 \|\nabla \mathbf{W}\|_2^2, \\ \|\nabla \mathbf{W} - \text{ADGP}(k_1, k, \nabla \mathbf{W}^t, \mathcal{I}^t)\|_2^2 &\leq \gamma_2 \|\nabla \mathbf{W}\|_2^2, \end{aligned}$$

200 where  $\gamma_1$  is a constant determined by  $k_1$  and  $k_2$ , and  $\gamma_2$  is a constant determined by  $k_1$  and  $k$ .

201 **Assumption 2.** *The objective function  $l : R^d \rightarrow R$  has a low bound  $l^*$  and it is Lipschitz-smooth, i.e.,*  
 202 *for any  $x_1, x_2$ ,  $\|\nabla l(x_1) - \nabla l(x_2)\|_2 \leq K\|x_1 - x_2\|_2$  and  $l(x_1) \leq l(x_2) + \langle \nabla l(x_2), x_1 - x_2 \rangle +$   
 203  $\frac{K}{2}\|x_1 - x_2\|_2^2$ .*

204 **Assumption 3.** *The full gradient  $\nabla l(\mathbf{W}^t)$  is bounded, i.e.,  $\|\nabla l(\mathbf{W}^t)\|_2^2 \leq G^2$ , and the federated*  
 205 *stochastic gradient  $\nabla \mathbf{W}_i^t$  ( $i = [1, N]$ ) is the unbiased estimation of the full gradient  $\nabla l(\mathbf{W}^t)$ , i.e.,*  
 206  $\mathbb{E}(\nabla \mathbf{W}_i^t) = \nabla l(\mathbf{W}^t)$ . *Moreover, the variance between  $\nabla \mathbf{W}_i^t$  and  $\nabla l(\mathbf{W}^t)$  is bounded:  $\mathbb{E}\|\nabla \mathbf{W}_i^t -$   
 207  $\nabla l(\mathbf{W}^t)\|_2^2 \leq \sigma^2$ .*

208 **Security Analysis.** It is noted that, for the same sparsification ratio, user's uploaded gradient  
 209 parameters from ADGP is generally smaller than that from DGP. Indeed, the uploaded gradient  
 210 parameters from both methods are the same only when  $\mathcal{T}_{k_1} \subset \mathcal{T}_{2k}$  holds. From this observation and  
 211 referring to Proposition 1, DGP is the security lower bound of our design for privacy protection.  
 212 So, our focus is the security analysis of DGP. As shown in the theorem below, we prove that DGP  
 213 achieves a stronger privacy protection in the sense of Definition 1.

214 **Theorem 1.** *For any  $(\varepsilon, \delta)$  optimization attack, under the presence of DGP, it will be degenerated*  
 215 *to  $(\varepsilon + \sqrt{\gamma_2}\|\nabla \mathbf{W}\|_2, \delta)$ -attack if  $\mathcal{D}$  is measured by Euclidean distance, and degenerated to  $(1 -$   
 216  $\sqrt{\gamma_1}(1 - \varepsilon), \delta)$ -attack if  $\mathcal{D}$  is measured by cosine distance.*

217 The Theorem is based on Assumption 1 about DGP. It reveals that, with the same successful chance  
 218  $1 - \delta$ , DGP weakens the attacker's capability to optimize a better estimation of the true  $\nabla \mathbf{W}$ .

219 **Generalization and Convergence Analyses.** The generalization analysis aims to quantify how the  
 220 trained model performs on the test data, and it is achieved by analyzing the how ADGP affects the  
 221 properties of the optima reached (without gradient pruning) [41, 42]. Assisted with Assumption 3  
 222 and Assumption 1 about ADGP gradient pruning, the following Lemma can be obtained.

223 **Lemma 1.** *Let  $\mathbf{e}^t = \sum_{i=1}^N \mathbf{e}_i^t/N$  be the averaged accumulated error among all users at iteration  $t$ ,*  
 224 *the expectation of the norm of  $\mathbf{e}^t$  is bounded, i.e.,*

$$\mathbb{E}\|\mathbf{e}^t\|_2^2 \leq \frac{\gamma_2}{2} \left( \frac{2 + \gamma_2}{1 - \gamma_2} \right)^2 (G^2 + \sigma^2). \quad (5)$$

225 Note that the difference between the averaged pruned gradient  $\mathbf{g}^t = \sum_{i=1}^N \mathbf{g}_i^t/N$  and the averaged  
 226 Fed-SGD gradient  $\nabla \mathbf{W}^t = \sum_{i=1}^N \nabla \mathbf{W}_i^t/N$  is simply  $\|\sum_{i=0}^{T-1} (\nabla \mathbf{W}^t - \mathbf{g}^t)\|_2^2 = \|\mathbf{e}^T\|_2^2$ . So  
 227 the lemma above indicates that the accumulated gradient difference between our algorithm and  
 228 Fed-SGD is bounded. That said, the optima reached by ADGP and the optima reached by Fed-  
 229 SGD will eventually be the same if the algorithm converge. Armed with Lemma 1 and based on  
 230 Assumptions 1, 2 and 3, we demonstrate the convergence of the our algorithm.

231 **Theorem 2.** *The averaged norm of the full gradient  $\nabla l(\mathbf{W}^t)$  derived from centralized training is*  
 232 *correlated with the our algorithm as follows:*

$$\frac{\sum_{t=0}^{T-1} \mathbb{E}\|\nabla l(\mathbf{W}^t)\|_2^2}{T} \leq 4 \frac{l^0 - l^*}{\eta T} + 4\eta^2 K^2 \frac{\gamma_2}{2} \left( \frac{2 + \gamma_2}{1 - \gamma_2} \right)^2 (G^2 + \sigma^2) + 2K\eta(G^2 + \frac{\sigma^2}{N}), \quad (6)$$

233 where  $l^0$  is the initialization of the objective  $l$ , and  $\eta$  is the learning rate.

234 The immediate implication of Theorem 2 is that, with an appropriate learning rate  $\eta$ , ADGP converges  
 235 similar to Fed-SGD (slower by a negligible term  $\mathcal{O}(\frac{1}{T})$ ), as shown in Corollary 1.

Table 2: Evaluation of defense performance under three attacks.

Attack	Metric	Baseline	Precode	DP	Soteria	ATS-I	ATS-II	Top- $k$	Ours
IVG	PSNR ( $\downarrow$ )	34.8805	9.6441	<b>6.9554</b>	9.2447	16.6894	31.3200	14.1338	7.6192
	LPIPS ( $\uparrow$ )	0.0016	0.4473	<b>0.5504</b>	0.3774	0.1621	0.0015	0.2754	0.4829
	SSIM ( $\downarrow$ )	0.9273	0.4793	<b>0.2451</b>	0.4173	0.6851	0.9189	0.5336	0.2923
R-GAP	PSNR ( $\downarrow$ )	36.7656	-	<b>5.0691</b>	5.1817	10.8442	42.0900	5.1017	5.1196
	LPIPS ( $\uparrow$ )	0.0007	-	0.3621	0.3532	0.2094	1.8e-05	0.4817	<b>0.4863</b>
	SSIM ( $\downarrow$ )	0.9307	-	0.2483	0.2124	0.3962	0.9121	0.2027	0.1928
Rob	PSNR ( $\downarrow$ )	102.8838	109.6553	<b>8.7491</b>	102.8838	9.6166	115.9886	13.0685	13.0804
	LPIPS ( $\uparrow$ )	0.0960	0.1488	<b>1.3434</b>	0.0960	0.6410	0.0486	0.8920	0.9184
	SSIM ( $\downarrow$ )	0.8969	0.8440	0.2064	0.8969	0.2545	0.9490	0.0428	<b>0.0229</b>
Final model accuracy		<b>93.4400</b>	93.1699	86.8900	93.2300	93.3900	93.3900	93.2099	93.1700

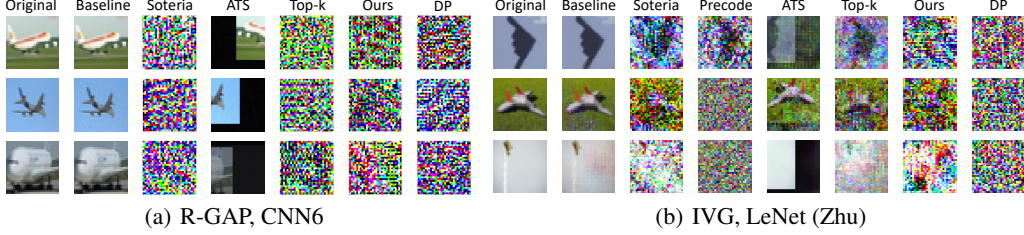


Figure 1: Visualization of the reconstructed data under R-GAP and IVG attacks.

236 **Corollary 1.** Let  $\eta = \sqrt{\frac{l^0 - l^*}{KT(G^2 + \sigma^2/N)}}$ , we have

$$\frac{\sum_{t=0}^{T-1} \mathbb{E} \|\nabla l(\mathbf{W}^t)\|_2^2}{T} \leq 6 \sqrt{\frac{l^0 - l^*}{KT(G^2 + \sigma^2/N)}} + \mathcal{O}\left(\frac{1}{T}\right). \quad (7)$$

## 237 6 Experiments: Privacy-Accuracy-Efficiency Tradeoff

### 238 6.1 Experimental Setup

239 **Datasets and models.** We conduct experiments on CIFAR10 with LeNet (Zhu) [5], CIFAR10 [45]  
 240 with CNN6 and CIFAR100 [45] with LeNet (Zhu) and ResNet18 respectively. We run these experi-  
 241 ments in a pytorch environment by using a single RTX 2080 Ti GPU and 2.10GHz CPU.

242 **Evaluation Metrics.** We quantify the privacy effect of defenses, follow [39, 46], we visualize the  
 243 reconstructed data and use learned perceptual image patch similarity (LPIPS), peak signal-to-noise  
 244 ratio (PSNR), structural similarity (SSIM) to measure the quality of the recovered data. A better  
 245 defense scheme should has a larger LPIPS ( $\uparrow$ ), smaller peak signal-to-noise ratio (PSNR) ( $\downarrow$ ) and  
 246 structural similarity (SSIM) ( $\downarrow$ ).

247 **Attack methods.** We evaluate our defense against IVG attack [5], R-GAP attack [22], and Rob  
 248 attack [21], which represent three kinds of state-of-the-art GIAs, as illustrated in Section 3. All  
 249 these attacks are implemented strictly following the original settings, *i.e.*, IVG is evaluated on  
 250 CIFAR10 with LeNet (Zhu), R-GAP is evaluated on CIFAR10 with CNN6, Rob attack is evaluated  
 251 on CIFAR100 with LeNet (Zhu). More settings for attacks are shown in the supplementary.

252 **Defense methods.** We compare our method with five state-of-the-art defenses: Soteria [10],  
 253 ATS [11], Precode [12], Differential Privacy (DP) [2], and Top- $k$  based gradient sparsification<sup>1</sup> [19].  
 254 Besides, we set Fed-SGD [3] as the baseline that adopts no defenses. Following the DP design  
 255 in [2], we use the Gaussian differential privacy mechanism with  $\epsilon = 10.7$ ,  $\delta = 10^{-5}$ , which is the  
 256 suggested best setting for the privacy-accuracy trade-off and can make most models converge. When  
 257 quantifying the defense performance of ATS, we not only evaluate the similarity between the raw  
 258 images and the recovered data (ATS-I), but also evaluate the similarity between the disturbed training

<sup>1</sup>Hereinafter, we abuse the notion of  $k$  to denote the send rate  $(k/\text{Size}(\nabla \mathbf{W})) \times 100\%$  since it will not cause ambiguity. And the sparse ratio is 1- $k$ . The smaller the ratio  $k$  is, the better communication efficiency.

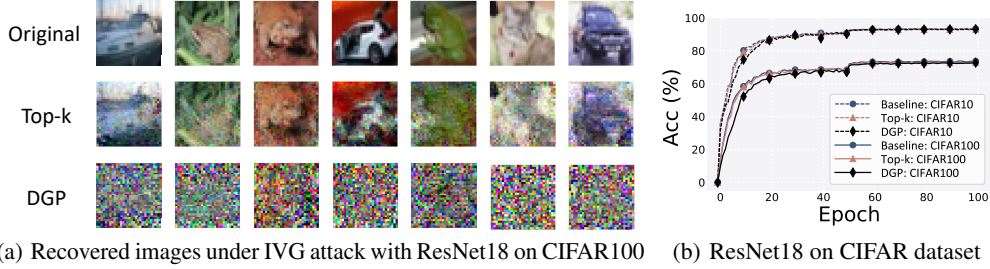


Figure 2: A detailed comparison between Top- $k$  and our DGP on privacy and model accuracy.

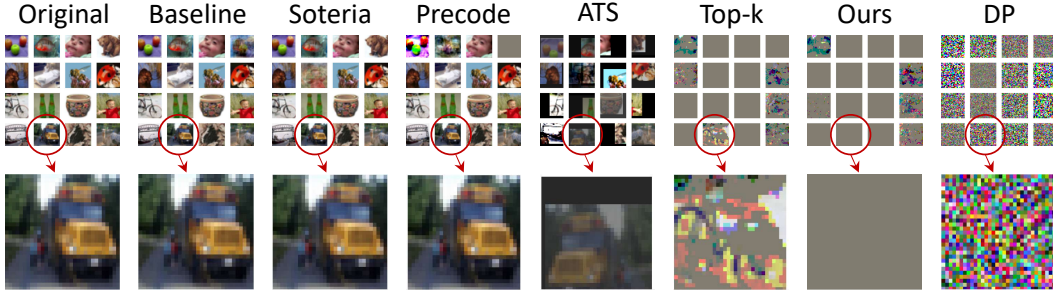


Figure 3: Visualization of reconstructed images under Rob attack with batchsize=16.

259 images (*i.e.*, the real inputs) and the recovered data (ATS-II). We set the send rate  $k = 0.2$  and the  
 260 regulation hyperparameter  $p = 15$ . The supplementary gives more experiments under different  $p$  and  
 261  $k$ . The rest defense schemes remain the original settings.

## 262 6.2 Defense Performance Evaluation

263 Table 2 shows the defense performance with PSNR, SSIM, and LPIPS under three attacks. The  
 264 results show that ATS, Soteria, Precode perform poorly under Rob attack, while Top- $k$  is vulnerable  
 265 to IVG attack although it performs better under Rob attack. In most cases, our scheme performs  
 266 comparably with DP and outperforms all the other defenses. More evaluation results under Rob  
 267 attack are presented in our supplementary.

268 We also visualize the reconstructed images in order to perceptually demonstrate the defense per-  
 269 formance. Figure 1(a) shows the the recovered images against R-GAP and IVG attacks. We can  
 270 see that all the existing defenses can well defend against R-GAP attack except ATS because it does  
 271 not damage the gradient structure, proving that a slight perturbation on gradients can mitigate the  
 272 analytical attacks easily. We are not able to provide the result of Precode because its VB operation  
 273 destroys the model structure thus analytical attack R-GAP cannot be implemented. In Figure 1(b),  
 274 recovered images under IVG attack are presented. We can find that the attacker can still recover the  
 275 outline of inputs with ATS and Top- $k$  defenses. DP, Soteria, Precode, and our scheme can still make  
 276 the recovered images unrecognizable. Figure 3 evaluates the defenses against Rob attack. It shows  
 277 that ATS, Soteria, and Precode fail to work and most inputs can be reconstructed.

278 In Rob attack, the attacker uses the gradient of the imprint module to reconstruct the training data.  
 279 Our method, Top- $k$ , and DP can effectively defend against Rob attack because the gradients of  
 280 all layers are sparsed or perturbed, including those of the malicious imprint modules. However,  
 281 we emphasize that the main weakness of Top- $k$  is its vulnerability to optimization attacks (*e.g.*,  
 282 IVG), as widely demonstrated in the literature [4, 10, 11, 39, 12]. We thus further evaluate Top- $k$   
 283 and our scheme under IVG attack with ResNet18 on CIFAR datasets. We set  $k_1/\text{Size}(\nabla\mathbf{W}) =$   
 284  $0.05$ ,  $k_2/\text{Size}(\nabla\mathbf{W}) = 0.75$ .

## 285 6.3 Model Accuracy Evaluation

286 To evaluate model performance, we train ResNet18, LeNet (Zhu), VGG13\_bn [47] on CIFAR10,  
 287 CIFAR100 with ten users, respectively. We set epoch=100, the learning rate  $\eta=0.1$  if epoch  $\leq 50$ ,



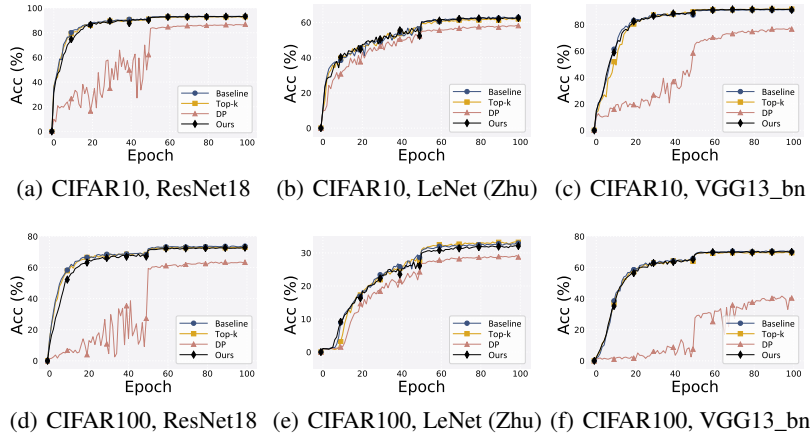


Figure 4: Evaluation of model accuracy with different datasets and model architectures.

Table 3: Commu. cost in one iteration (MB)

Method	Baseline	Soteria	Precode	ATS	DP	Top-k	Ours
Resnet18	85.2506	85.2268	88.1644	85.2506	85.2506	43.7979	27.3067
VGG13	71.8385	71.7318	74.8433	71.8385	71.8385	34.9625	22.9697
LeNet	0.1207	0.0764	1.8624	0.1207	0.1207	0.0493	0.0424

288  $\eta=0.01$  if epoch  $>50$ , and  $\eta=0.05$  if epoch  $>70$ . We show in Table 2 the accuracy of ResNet18 over  
 289 CIFAR10 under different defenses, and here we only compare our scheme with the baseline Fed-SGD,  
 290 DP, and Top- $k$  since they perform best for privacy protection. Because [41] showed that the error  
 291 feedback is beneficial to improve the model accuracy, even without using gradient sparsification. To  
 292 give a fair comparison, we set the error feedback mechanism as the basic setting for all the defenses.  
 293 The experimental results in Figure 4 show that we achieve similar model performance with the  
 294 baseline, while DP, as expected, significantly damage the model accuracy.

## 295 6.4 Efficiency Evaluation

296 To clearly demonstrate the system efficiency, we evaluate the communication cost, which is obtained  
 297 by computing the total overheads of sending updated gradients and receiving aggregated gradients.  
 298 For ease of presentation, we only show the results for one iteration. As shown in Table 3, our scheme  
 299 reduces more than half of the communication costs compared with existing defenses, and our gradient  
 300 sparsification incurs negligible computation burden. The specific computation cost evaluation is  
 301 presented in the supplementary.

## 302 7 Conclusions, Limitations, and Broader Impact

303 Our work firstly reveals the risks of privacy-preserving methods that only perturb the gradients of  
 304 some layers. Through a comprehensive analysis of gradient inversion attacks, we show that it is  
 305 necessary to perturb or sparse the gradients of each layer for privacy preservation. And considering  
 306 the challenge of high communication cost in federated learning, we propose aligned dual gradient  
 307 sparsification method to achieve the trade-off between privacy protection, model performance, and  
 308 efficient communication, and give sufficient theoretical support. We hope that our newly proposed  
 309 gradient sparsification method can shed new light on addressing privacy leakage concern as well as  
 310 saving communication bandwidth.

311 In terms of limitations, the success of our scheme relies on selecting a reliable user to broadcast its  
 312 gradient locations. Randomly selecting users may encounter malicious users that destroy the entire  
 313 system. Our design is delegated to protecting privacy and has no negative societal impacts in practice.

314 **References**

- 315 [1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient  
316 learning of deep networks from decentralized data,” in *Proceedings of the 20th International  
317 Conference on Artificial Intelligence and Statistics (AISTATS’17)*, 2017, pp. 1273–1282.
- 318 [2] M. Naseri, J. Hayes, and E. De Cristofaro, “Local and central differential privacy for robustness  
319 and privacy in federated learning,” in *Proceedings of the 29th Network and Distributed System  
320 Security Symposium (NDSS’22)*, 2022.
- 321 [3] Q. Li, Z. Wen, Z. Wu, S. Hu, N. Wang, Y. Li, X. Liu, and B. He, “A survey on federated  
322 learning systems: vision, hype and reality for data privacy and protection,” *IEEE Transactions  
323 on Knowledge and Data Engineering*, 2021.
- 324 [4] L. Zhu, Z. Liu, and S. Han, “Deep leakage from gradients,” in *Advances in Neural Information  
325 Processing Systems (NeurIPS’19)*, 2019, pp. 14 747–14 756.
- 326 [5] J. Geiping, H. Bauermeister, H. Dröge, and M. Moeller, “Inverting gradients-how easy is it to  
327 break privacy in federated learning?” in *Advances in Neural Information Processing Systems  
328 (NeurIPS’20)*, 2020, pp. 16 937–16 947.
- 329 [6] C. Dwork, A. Roth *et al.*, “The algorithmic foundations of differential privacy,” *Found. Trends  
330 Theor. Comput. Sci.*, vol. 9, no. 3-4, pp. 211–407, 2014.
- 331 [7] R. Gilad-Bachrach, K. Laine, K. Lauter, P. Rindal, and M. Rosulek, “Secure data exchange:  
332 A marketplace in the cloud,” in *Proceedings of the 2019 ACM SIGSAC Conference on Cloud  
333 Computing Security Workshop (CCSW’19)*, 2019, pp. 117–128.
- 334 [8] S. Hardy, W. Henecka, H. Ivey-Law, R. Nock, G. Patrini, G. Smith, and B. Thorne, “Private fed-  
335 erated learning on vertically partitioned data via entity resolution and additively homomorphic  
336 encryption,” *arXiv preprint arXiv:1711.10677*, 2017.
- 337 [9] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage,  
338 A. Segal, and K. Seth, “Practical secure aggregation for privacy-preserving machine learning,” in  
339 *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security  
340 (CCS’17)*, 2017, pp. 1175–1191.
- 341 [10] J. Sun, A. Li, B. Wang, H. Yang, H. Li, and Y. Chen, “Soteria: Provable defense against privacy  
342 leakage in federated learning from representation perspective,” in *Proceedings of the 2021  
343 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR’21)*, 2021, pp.  
344 9311–9319.
- 345 [11] W. Gao, S. Guo, T. Zhang, H. Qiu, Y. Wen, and Y. Liu, “Privacy-preserving collaborative learn-  
346 ing with automatic transformation search,” in *Proceedings of the 2021 IEEE/CVF Conference  
347 on Computer Vision and Pattern Recognition (CVPR’21)*, 2021, pp. 114–123.
- 348 [12] D. Scheliga, P. Mäder, and M. Seeland, “Precode-a generic model extension to prevent deep  
349 gradient leakage,” in *Proceedings of the 2022 IEEE/CVF Winter Conference on Applications of  
350 Computer Vision (WACV’22)*, 2022, pp. 1849–1858.
- 351 [13] W. Wei, L. Liu, M. Loper, K.-H. Chow, M. E. Gursoy, S. Truex, and Y. Wu, “A framework for  
352 evaluating gradient leakage attacks in federated learning,” *arXiv preprint arXiv:2004.10397*,  
353 2020.
- 354 [14] Y. Wang, J. Deng, D. Guo, C. Wang, X. Meng, H. Liu, C. Ding, and S. Rajasekaran, “Sapag: A  
355 self-adaptive privacy attack from gradients,” *arXiv preprint arXiv:2009.06228*, 2020.
- 356 [15] M. Balunović, D. I. Dimitrov, R. Staab, and M. Vechev, “Bayesian framework for gradient  
357 leakage,” *arXiv preprint arXiv:2111.04706*, 2021.
- 358 [16] F. Boenisch, A. Dziedzic, R. Schuster, A. S. Shamsabadi, I. Shumailov, and N. Paper-  
359 not, “When the curious abandon honesty: Federated learning is not private,” *arXiv preprint  
360 arXiv:2112.02918*, 2021.
- 361 [17] X. Pan, M. Zhang, Y. Yan, J. Zhu, and M. Yang, “Exploring the security boundary of data  
362 reconstruction via neuron exclusivity analysis.”
- 363 [18] Y. Wen, J. Geiping, L. Fowl, M. Goldblum, and T. Goldstein, “Fishing for user data in large-  
364 batch federated learning via gradient magnification,” *arXiv preprint arXiv:2202.00580*, 2022.

- 365 [19] D. Alistarh, T. Hoeffler, M. Johansson, N. Konstantinov, S. Khirirat, and C. Renggli, “The  
366 convergence of sparsified gradient methods,” in *Advances in Neural Information Processing  
367 Systems (NeurIPS’18)*, 2018, pp. 5977–5987.
- 368 [20] Y. Lin, S. Han, H. Mao, Y. Wang, and W. J. Dally, “Deep gradient compression: Reducing the  
369 communication bandwidth for distributed training,” *arXiv preprint arXiv:1712.01887*, 2017.
- 370 [21] L. Fowl, J. Geiping, W. Czaja, M. Goldblum, and T. Goldstein, “Robbing the fed: Directly obtain-  
371 ing private data in federated learning with modified models,” *arXiv preprint arXiv:2110.13057*,  
372 2021.
- 373 [22] J. Zhu and M. Blaschko, “R-gap: Recursive gradient attack on privacy,” in *Proceedings of the  
374 9th International Conference on Learning Representations (ICLR’21)*, 2021.
- 375 [23] R. C. Geyer, T. Klein, and M. Nabi, “Differentially private federated learning: A client level  
376 perspective,” *arXiv preprint arXiv:1712.07557*, 2017.
- 377 [24] B. Zhao, K. R. Mopuri, and H. Bilen, “idlg: Improved deep leakage from gradients,” *arXiv  
378 preprint arXiv:2001.02610*, 2020.
- 379 [25] J. Qian and L. K. Hansen, “What can we learn from gradients?” *arXiv preprint  
380 arXiv:2010.15718*, 2020.
- 381 [26] H. Yin, A. Mallya, A. Vahdat, J. M. Alvarez, J. Kautz, and P. Molchanov, “See through gradients:  
382 Image batch recovery via gradinversion,” in *Proceedings of the 2021 IEEE/CVF Conference on  
383 Computer Vision and Pattern Recognition (CVPR’21)*, 2021, pp. 16 337–16 346.
- 384 [27] L. Fan, K. W. Ng, C. Ju, T. Zhang, C. Liu, C. S. Chan, and Q. Yang, “Rethinking privacy  
385 preserving deep learning: How to evaluate and thwart privacy attacks,” in *Federated Learning*.  
386 Springer, 2020, vol. 12500, pp. 32–50.
- 387 [28] X. Chen, Z. S. Wu, and M. Hong, “Understanding gradient clipping in private sgd: A geometric  
388 perspective,” in *Advances in Neural Information Processing Systems (NeurIPS’20)*, 2020, pp.  
389 13 773–13 782.
- 390 [29] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, “Deep  
391 learning with differential privacy,” in *Proceedings of the 2016 ACM SIGSAC conference on  
392 computer and communications security (CCS’16)*, 2016, pp. 308–318.
- 393 [30] L. Yu, L. Liu, C. Pu, M. E. Guroy, and S. Truex, “Differentially private model publishing for  
394 deep learning,” in *Proceedings of the 2019 IEEE Symposium on Security and Privacy (SP’19)*,  
395 2019, pp. 332–349.
- 396 [31] P. Mohassel and Y. Zhang, “Secureml: A system for scalable privacy-preserving machine  
397 learning,” in *Proceedings of the 2017 IEEE Symposium on Security and Privacy (SP’17)*, 2017,  
398 pp. 19–38.
- 399 [32] G. Danner and M. Jelasity, “Fully distributed privacy preserving mini-batch gradient descent  
400 learning,” in *Proceedings of the 15th International conference on distributed applications and  
401 interoperable systems (IFIP’15)*, 2015, pp. 30–44.
- 402 [33] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. Quek, and H. V. Poor,  
403 “Federated learning with differential privacy: Algorithms and performance analysis,” *IEEE  
404 Transactions on Information Forensics and Security*, vol. 15, pp. 3454–3469, 2020.
- 405 [34] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz,  
406 Z. Charles, G. Cormode, R. Cummings *et al.*, “Advances and open problems in federated  
407 learning,” *Foundations and Trends® in Machine Learning*, vol. 14, no. 1–2, pp. 1–210, 2021.
- 408 [35] R. Liu, Y. Cao, H. Chen, R. Guo, and M. Yoshikawa, “Flame: Differentially private federated  
409 learning in the shuffle model,” *arXiv preprint arXiv:2009.08063*, 2020.
- 410 [36] L. Sun, J. Qian, and X. Chen, “Ldp-fl: Practical private aggregation in federated learning with  
411 local differential privacy,” in *Proceedings of the Thirtieth International Joint Conference on  
412 Artificial Intelligence (IJCAI’21)*, 2021, pp. 1571–1578.
- 413 [37] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy, “Deep variational information bottleneck,”  
414 *arXiv preprint arXiv:1612.00410*, 2016.
- 415 [38] Z. Wang, M. Song, Z. Zhang, Y. Song, Q. Wang, and H. Qi, “Beyond inferring class repre-  
416 sentatives: User-level privacy leakage from federated learning,” in *Proceedings of 2019 IEEE  
417 Conference on Computer Communications (INFOCOM’19)*, 2019, pp. 2512–2520.

- 418 [39] Y. Huang, S. Gupta, Z. Song, K. Li, and S. Arora, “Evaluating gradient inversion attacks  
419 and defenses in federated learning,” in *Advances in Neural Information Processing Systems*  
420 (*NeurIPS’21*), 2021, pp. 7232–7241.
- 421 [40] L. T. Phong, Y. Aono, T. Hayashi, L. Wang, and S. Moriai, “Privacy-preserving deep learning:  
422 Revisited and enhanced,” in *Proceedings of the 2017 International Conference on Applications*  
423 *and Techniques in Information Security (ATIS’17)*, 2017, pp. 100–110.
- 424 [41] S. P. Karimireddy, Q. Rebjock, S. Stich, and M. Jaggi, “Error feedback fixes signsgd and  
425 other gradient compression schemes,” in *Proceedings of the 36th International Conference on*  
426 *Machine Learning (ICML’19)*, 2019, pp. 3252–3261.
- 427 [42] A. C. Wilson, R. Roelofs, M. Stern, N. Srebro, and B. Recht, “The marginal value of adaptive  
428 gradient methods in machine learning,” in *Advances in Neural Information Processing Systems*  
429 (*NeurIPS’17*), 2017, pp. 4148–4158.
- 430 [43] C.-Y. Chen, J. Ni, S. Lu, X. Cui, P.-Y. Chen, X. Sun, N. Wang, S. Venkataramani, V. V.  
431 Srinivasan, W. Zhang *et al.*, “Scalecom: Scalable sparsified gradient compression for  
432 communication-efficient distributed training,” in *Advances in Neural Information Process-*  
433 *ing Systems (NeurIPS’20)*, 2020, pp. 13 551–13 563.
- 434 [44] X. Dai, X. Yan, K. Zhou, H. Yang, K. K. Ng, J. Cheng, and Y. Fan, “Hyper-sphere quantization:  
435 Communication-efficient sgd for federated learning,” *arXiv preprint arXiv:1911.04655*, 2019.
- 436 [45] A. Krizhevsky, G. Hinton *et al.*, “Learning multiple layers of features from tiny images,” 2009.
- 437 [46] J. Jeon, K. Lee, S. Oh, J. Ok *et al.*, “Gradient inversion with generative image prior,” in *Advances*  
438 *in Neural Information Processing Systems (NeurIPS’21)*, 2021, pp. 29 898–29 908.
- 439 [47] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image  
440 recognition,” *arXiv preprint arXiv:1409.1556*, 2014.

## 441 Checklist

- 442 1. For all authors...
- 443 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
444 contributions and scope? [Yes]
- 445 (b) Did you describe the limitations of your work? [Yes]
- 446 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 447 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
448 them? [Yes]
- 449 2. If you are including theoretical results...
- 450 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 451 (b) Did you include complete proofs of all theoretical results? [Yes]
- 452 3. If you ran experiments...
- 453 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
454 mental results (either in the supplemental material or as a URL)? [Yes]
- 455 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
456 were chosen)? [Yes] See Sec. 6
- 457 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
458 ments multiple times)? [Yes] see supplementary
- 459 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
460 of GPUs, internal cluster, or cloud provider)? [Yes]
- 461 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 462 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 463 (b) Did you mention the license of the assets? [No]
- 464 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 465 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
466 using/curating? [Yes] all the license of the data, basically the ones used for research  
467 are openly accessible

- 468 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
469 information or offensive content? [No]
- 470 5. If you used crowdsourcing or conducted research with human subjects...
- 471 (a) Did you include the full text of instructions given to participants and screenshots, if  
472 applicable? [N/A]
- 473 (b) Did you describe any potential participant risks, with links to Institutional Review  
474 Board (IRB) approvals, if applicable? [N/A]
- 475 (c) Did you include the estimated hourly wage paid to participants and the total amount  
476 spent on participant compensation? [N/A]