

---

# AIMHI: Protecting Sensitive Data through Federated Co-Training

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Federated learning offers collaborative training among distributed sites without  
2 sharing sensitive local information by sharing the sites' model parameters. It is  
3 possible, though, to make non-trivial inferences about sensitive local information  
4 from these model parameters. We propose a novel co-training technique called  
5 AIMHI that uses a public unlabeled dataset to exchange information between  
6 sites by sharing predictions on that dataset. This setting is particularly suitable  
7 to healthcare, where hospitals and clinics hold small labeled datasets with highly  
8 sensitive patient data and large national health databases contain large amounts of  
9 public patient data. We show that the proposed method reaches a model quality  
10 comparable to federated learning while maintaining privacy to high degree.

## 11 1 Introduction

12 Can we collaboratively train models from distributed sensitive datasets while maintaining data privacy  
13 at a level required in healthcare? Federated learning [12] allows distributed sites, e.g., hospitals  
14 or clinics, to collaboratively train a joint model without directly disclosing their sensitive data by  
15 instead periodically sharing model parameters. An attacker or curious observer can, however, make  
16 inferences about local data from model parameters [11] and model updates [21]. Differential privacy  
17 provides a rigorous and measurable privacy guarantee [5] that can be achieved by perturbing model  
18 parameters appropriately citepwei2020federated. This perturbation, however, can reduce model  
19 quality, resulting in a trade-off between privacy and quality. As we show in our experiments, even  
20 with substantial perturbation one can infer membership of a training sample [17] with high probability,  
21 i.e., whether a data point is present in a local dataset from the model parameters shared in federated  
22 learning.

23 We propose to instead use a distributed co-training approach [9], where sites train local models and  
24 exchange predictions on a shared unlabeled dataset, instead of sharing model parameters. By forming  
25 a consensus from the shared predictions, one obtains pseudo-labels for the shared unlabeled dataset  
26 that can be used for local training. Iterating this process improves the consensus, and thereby the  
27 quality of pseudo-labels, effectively nudging local models to come to an agreement. We show in  
28 our experiments that this approach, which we call AIMHI (AI Models for Healthcare Improvement),  
29 achieves the same model quality as federated learning (FedAvg [12]), but protects privacy to a high  
30 level—membership inference is significantly less likely compared to vanilla federated learning and  
31 federated learning with differential privacy. These results indicate that, if a public unlabeled dataset  
32 is available, this approach constitutes a more favorable trade-off between privacy and model quality.

33 The AIMHI approach requires a large shared unlabeled dataset, and is in spirit similar to distributed  
34 distillation [2]. While such unlabeled public datasets are not always available, in healthcare, large  
35 public health databases are quite common: the US NCHS databases, the UK's NHS databases, the

36 UK Biobank [18], the MIMIC-III database [6], or the planned European EHDS contain vast amounts  
 37 of patient data that can be used in many application scenarios.

38 Our contributions are

- 39 (i) a novel distributed co-training approach to collaboratively train models from privacy-  
 40 sensitive distributed data sources, and
- 41 (ii) a preliminary empirical evaluation of model quality and privacy on the CIFAR10 benchmark  
 42 dataset, indicating high model quality and a substantial improvement in privacy.

## 43 2 Preliminaries

44 We assume learning algorithms  $\mathcal{A} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{H}$  that trains a model  $h \in \mathcal{H}$  using a dataset  
 45  $D \subset \mathcal{X} \times \mathcal{Y}$  from an input space  $\mathcal{X}$  and output space  $\mathcal{Y}$ , i.e.,  $h = \mathcal{A}(D)$ . Given a set of  $m \in \mathbb{N}$   
 46 clients with local datasets  $D^1, \dots, D^m \subset \mathcal{X} \times \mathcal{Y}$  drawn iid from a data distribution  $\mathcal{D}$  and a loss  
 47 function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ , the goal is to find a single model  $h^* \in \mathcal{H}$  that minimizes the risk

$$\mathcal{E}_{\mathcal{D}}(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}}(\ell(h(x), y)) .$$

48 In centralized learning, the datasets are pooled as  $D = \bigcup_{i \in [m]} D^i$  and  $\mathcal{A}$  is applied to  $D$ , usually to  
 49 minimize the empirical risk

$$\mathcal{E}_{emp}(h, D) = \sum_{(x,y) \in D} \ell(h(x), y) .$$

50 In federated learning (FL) we assume that the learning algorithm is iterative [cf. Chp. 2.1.4 7], i.e.,  
 51  $\mathcal{A} : \mathcal{X} \times \mathcal{Y} \times \mathcal{H} \rightarrow \mathcal{H}$  that updates a model  $h_{t+1} = \mathcal{A}(D, h_t)$ . In this case, centralized learning  
 52 means applying  $\mathcal{A}$  to  $D$  until convergence. Note that applying  $\mathcal{A}$  on  $D$  can be the application to any  
 53 random subset, e.g., as in mini-batch training, and convergence is measured in terms of low training  
 54 loss, small gradient, or small deviation from previous iterate.

55 In standard federated learning [12],  $\mathcal{A}$  is applied in parallel for  $b \in \mathbb{N}$  rounds on each client locally  
 56 to produce local models  $h^1, \dots, h^m$ . These models are then centralized and aggregated using an  
 57 aggregation operator  $\mathfrak{a} : \mathcal{H}^m \rightarrow \mathcal{H}$ , i.e.,  $\bar{h} = \mathfrak{a}(h^1, \dots, h^m)$ . The aggregated model  $\bar{h}$  is then  
 58 redistributed to local clients which perform another  $b$  rounds of training using  $\bar{h}$  as a starting point.  
 59 This is iterated until convergence of  $\bar{h}$  with the goal to minimize the empirical risk over all local  
 60 datasets[12], i.e.,

$$\mathcal{E}_{emp}(h, D^1, \dots, D^m) = \frac{1}{m} \sum_{k=1}^m \mathcal{E}_{emp}(h, D^k) = \frac{1}{m} \sum_{k=1}^m \sum_{(x,y) \in D^k} \ell(h(x), y) .$$

61 When aggregating by averaging, this method is also known as federated averaging. Next, we describe  
 62 our proposed distributed co-training approach.

## 63 3 AIMHI: Distributed Co-Training

64 We propose a semi-supervised, distributed co-training approach that collaboratively trains models via  
 65 sharing predictions. It uses an unlabeled dataset  $U$ , producing pseudo-labels  $L$  for it by forming a  
 66 consensus of the predictions of all local models. Unlabeled data and pseudo-labels form an additional  
 67 public, shared dataset  $P$  that is combined with local data for training. The details are described in  
 68 Alg. 1: at each client  $i$ , the local model is updated using the local dataset  $D^i$  combined with the current  
 69 pseudo-labeled public dataset  $P$ . The updated model is used to produce improved pseudo-labels  
 70  $L^i$  for the unlabeled data  $U$ , which are sent to a server every  $b$  rounds. At the server, as soon as  
 71 all local prediction  $L^1, \dots, L^m$  are received, a consensus  $L$  is formed and broadcasted back to the  
 72 clients. Forming a consensus is similar to obtaining a prediction from an ensemble [4]. For our  
 73 classification experiments, we use vanilla majority voting [3]. Note that more elaborate consensus  
 74 mechanisms offer a rich design space for improvements. On receiving the new consensus labels  $L$   
 75 from the server, the client updates the public pseudo-labeled dataset  $P$  and performs another iteration  
 76 of local training.

---

**Algorithm 1:** AIMHI

---

**Input:** communication period  $b$ , learning algorithm  $\mathcal{A}$ ,  $m$  clients with local datasets  $D^1, \dots, D^m$ , unlabeled shared dataset  $U$ , total number of rounds  $T$

**Output:** final models  $h_T^1, \dots, h_T^m$

```
1 initialize local models  $h_0^1, \dots, h_0^m$ 
2  $P \leftarrow \emptyset$ 
3 Locally at client  $i$  at time  $t$  do
4    $h_t^i \leftarrow \mathcal{A}(D_i \cup P, h_{t-1}^i)$ 
5   if  $t \% b = b - 1$  then
6      $L^i \leftarrow h_t^i(U)$ 
7     send  $L^i$  to server
8     receive  $L$  from server
9      $P \leftarrow (U, L)$ 
10  end
11 At server at time  $t$  do
12  receive local pseudo-labels  $L^1, \dots, L^m$ 
13   $L \leftarrow \text{consensus}(L^1, \dots, L^m)$ 
14  send  $L$  to all clients
```

---

## 77 4 Empirical Evaluation

78 We perform a preliminary empirical evaluation of the performance of AIMHI in comparison to vanilla  
79 federated learning (i.e., FedAvg) and federated learning with differential privacy on the CIFAR10  
80 image classification dataset [10]. For that, we measure their prediction accuracy on a test set, as well  
81 as their privacy vulnerability.

### 82 4.1 Privacy Vulnerability

83 We measure privacy vulnerability by performing membership inference attacks against AIMHI and FL  
84 according to two different attack scenarios per approach. In both attacks, the attacker creates an attack  
85 model using a model it constructs from its training and test datasets. Similar to previous work [17], we  
86 assume that the training data of the attacker has a similar distribution to the training data of the client.  
87 Once the attacker has its attack model, it uses this model for membership inference. In blackbox  
88 attacks (in which the attacker does not have access to intermediate model parameters), it only uses the  
89 classification scores it receives from the target model (i.e., client’s model) for membership inference.  
90 On the other hand, in whitebox attacks (in which the attacker can observe the intermediate model  
91 parameters), it can use additional information in its attack model. Since the proposed AIMHI does  
92 not reveal intermediate model parameters to any party, it is only subject to blackbox attacks. Vanilla  
93 federated learning on the other hand is subject to whitebox attacks. Each inference attack produces a  
94 membership score of a queried data point, indicating the likelihood of the data point being a member  
95 of the training set. We measure the success of membership inference as ROC AUC of these scores.  
96 The **vulnerability (VUL)** of a method is the ROC AUC of membership attacks over  $K$  runs over  
97 the entire training set (also called attack epochs) according to the attack model and scenario. A  
98 vulnerability of 1.0 means that membership can be inferred with certainty, whereas 0.5 means that  
99 deciding on membership is a random guess.

100 We assume the following attack model: clients are honest and the server may be semi-honest (follow  
101 the protocol execution correctly, but it may try to infer sensitive information about the clients). The  
102 main goal of a semi-honest server is to infer sensitive information about the local training data  
103 of the clients. This is a stronger attacker assumption compared to a semi-honest client since the  
104 server receives the most amount of information from the clients during the protocol, and a potential  
105 semi-honest client can only obtain indirect information about the other clients. We also assume that  
106 parties do not collude. The attack scenarios are<sup>1</sup>:

---

<sup>1</sup>Note that the relaxed scenarios are more in favor of FL than the realistic ones.

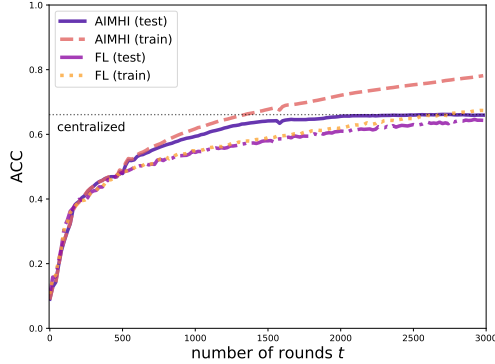


Figure 1: Test accuracy (ACC) over time on CIFAR10 with ACC of average model (FL) and average ACC of local models for AIMHI.

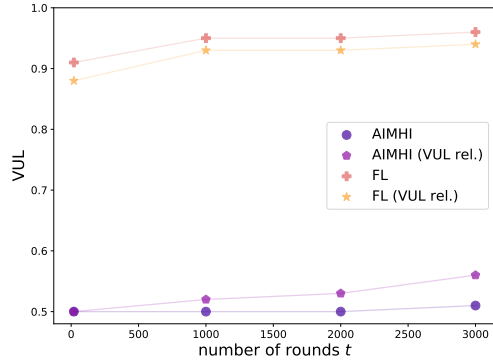


Figure 2: Privacy vulnerability (VUL) over time for AIMHI and FL for realistic and relaxed attack scenarios.

- 107 • *AIMHI realistic*: the attacker can send a (forged) unlabeled dataset to the clients and observe
- 108 their predictions, equivalent to one attack epoch ( $K = 1$ );
- 109 • *FL realistic*: the attacker receives model parameters and can run an arbitrary number of
- 110 attacks—we use  $K = 500$  attack epochs;
- 111 • *AIMHI relaxed*: the attacker can send a (forged) unlabeled dataset in each communication
- 112 round, albeit to different models in each round—we simulate this by assuming pessimistically
- 113 that models are sufficiently similar over all rounds and set  $K = T/b = 150$ ;
- 114 • *FL relaxed*: the attacker cannot copy model parameters on their machine, but can only
- 115 perform an attack epoch in each communication round, thus being able to perform  $K =$
- 116  $T/b = 150$  attack epochs after  $T = 3000$  rounds.

117 To measure vulnerability, we use the ML Privacy Meter tool [13]. This tool allows us to quantify  
 118 the privacy risks associated with machine learning models by performing a range of membership  
 119 inference attacks [14] on machine learning models and measuring attack success. It simulates different  
 120 levels of access and model knowledge for attackers, e.g., limiting attackers to only predictions, or  
 121 loss values, or assuming they have access to the model’s parameters. For our experiments, we assume  
 122 that for AIMHI the attacker only has access to the predictions (i.e., output of the last layer), while for  
 123 FL, the attacker can access all layers.

## 124 4.2 Experimental Setup

125 We use the common CIFAR10 dataset [10] which consists of 50000 training and 10000 test images  
 126 with 10 classes. We use 10000 samples drawn iid from the training images as training set, 40000  
 127 images as unlabeled dataset and the 10000 test images as test set. We use  $m = 5$  clients, each with  
 128 a local training set of  $n = 2000$  samples. Clients use a convolutional neural network (details are  
 129 provided in the Appendix) and communicate every  $b = 20$  rounds for FL and AIMHI. Both FL and  
 130 AIMHI are run for  $T = 3000$  rounds.

131 The experiments are implemented in TensorFlow [1], the code is publicly available<sup>2</sup>. For our  
 132 experiments, we use a simple CNN architecture given in Table 1. As optimizer we use Adam with  
 133 learning rate  $\alpha = 0.001$ , exponential decay rates  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ , and  $\epsilon = 10^{-7}$ .

## 134 4.3 Differential Privacy for Federated Learning

135 A common defense against membership inference attacks is applying appropriate clipping and noise  
 136 before sending models. This guarantees  $\epsilon, \delta$ -differential privacy for local data [20] at the cost of a  
 137 slight-to-moderate loss in model quality. This technique is also proven to defend against backdoor  
 138 and poisoning attacks [19]. We compare AIMHI against federated learning with differential privacy

<sup>2</sup>[https://anonymous.4open.science/r/AIM\\_HI](https://anonymous.4open.science/r/AIM_HI)

Layers	Activation function
Conv2D(32, 3, 3)	relu
MaxPooling2D(2, 2)	
Conv2D(64, 3, 3)	relu
MaxPooling2D((2, 2)	
Flatten layer	
Dense(64)	relu
Dense(10)	softmax

Table 1: Model architecture for AIMHI and federated learning.

	Accuracy	VUL	VUL (relaxed)
<b>AIMHI</b>	<b>0.659</b>	<b>0.51</b>	0.56
<b>FL</b>	0.643	0.96	0.94
<b>DP-FL</b> ( $C = 2.0, \sigma = 0.01$ )	0.425	0.85	0.87

Table 2: Test accuracy (ACC) and privacy vulnerability (VUL) for AIMHI and federated learning, both vanilla federated averaging (FL) and federated averaging with differential privacy (DP-FL) on  $m = 5$  clients with local training set size  $|D^1| = \dots = |D^m| = 8000$  and an unlabeled dataset of size  $|U| = 10000$ .

139 through noise (DP-FL). For that, each client  $i$  first clips their parameters  $w^i$  to

$$w_c^i = \frac{w^i}{\max\left\{1, \frac{\|w^i\|}{C}\right\}}$$

140 for a constant  $C > 0$ , and then add Gaussian noise  $\tilde{w}^i = w_c^i + \epsilon$  with  $\epsilon \sim \mathcal{N}(0, \sigma)$  for  $\sigma \geq 0$ . The  
141 level of privacy depends on the choice of  $C$  and  $\sigma$  [20].

#### 142 4.4 Results

143 The results presented in Table 2 show that AIMHI and FL achieve comparable test accuracy after  
144  $T = 3000$  rounds<sup>3</sup>. Note that centralized training on the 10000 training samples achieves a test  
145 accuracy of 0.661, so FL and AIMHI both achieve virtually optimal model quality. At the same  
146 time, FL is vulnerable to membership attacks, both in the realistic (VUL) and the relaxed (VUL(rel.))  
147 attack scenarios with a vulnerability of over 0.9. Note that since the attacker can run an arbitrary  
148 number of attacks, the privacy vulnerability is in principle 1.0, we instead measure the practical  
149 vulnerability under a large, but limited number of attacks. AIMHI on the other hand preserves privacy  
150 ( $VUL = 0.51$ ) in the realistic attack scenario, and still has very little vulnerability (0.56) under the  
151 relaxed scenario. Adding noise for differential privacy reduces vulnerability considerably to around  
152 0.85, but at the cost of accuracy which drops to around 0.42. We investigate the convergence in terms  
153 of test accuracy in Figure 1: Both AIMHI and FL converge quickly, with AIMHI converging slightly  
154 faster. Looking at the development of privacy vulnerability over time, we observe that vulnerability  
155 for FL is already high after the first communication round (i.e., after  $t = 20$  rounds), as shown in  
156 Figure 2, and increases to 0.94 after  $t = 1000$  rounds. This holds also for the relaxed scenario,  
157 although with slightly less vulnerability. For AIMHI realistic instead the vulnerability remains low  
158 (0.5 to 0.51). For AIMHI relaxed we observe an increase in vulnerability, as expected, since we  
159 assume in this scenario that the attacker performs a number of attacks proportional to the number of  
160 communication rounds.

<sup>3</sup>For FL, we report the test accuracy of the average model; for AIMHI, we report the average of test accuracies for local models—we observe that at  $T = 3000$  their variance is 0.

## 161 5 Discussion and Conclusion

162 We presented a novel distributed co-training approach for privacy-preserving federated learning that  
163 protects sensitive local datasets, where vanilla federated learning is susceptible to privacy attacks, at  
164 the same time achieving similar model quality as FL. Our initial experiments support the hypothesis  
165 that distributed co-training can be competitive with FL, given a large unlabeled dataset, while  
166 preserving data privacy to a much higher degree.

167 Co-training requires a shared unlabeled dataset, which is not available in all application scenarios. In  
168 healthcare, however, it is not uncommon to have large quantities of unlabeled data points available. A  
169 limitation of AIMHI is that local datasets must be sufficiently large to create useful local models [cf.  
170 8]—otherwise, the poor quality pseudo-labeled dataset will not improve local training. Choosing  
171 more elaborate consensus methods [15, 16] is an interesting direction for future work that can improve  
172 performance even with small local datasets. An important advantage of co-training is that local  
173 models do not have to have the same architecture, as FL requires. In fact, local models can be  
174 arbitrary. This includes interpretable models, like decision trees or rule ensembles, for which no FL  
175 method exists so far. Exploring distributed co-training for such model classes could open a whole  
176 new avenue for federated learning.

## 177 References

- 178 [1] Martín Abadi, others Dean, Tucker, Yu, and TensorFlow: Large-scale machine learning on  
179 heterogeneous systems, 2015. Software available from tensorflow.org. 4
- 180 [2] Ilai Bistriz, Ariana Mann, and Nicholas Bambos. Distributed distillation for on-device learning.  
181 *Advances in Neural Information Processing Systems*, 33:22593–22604, 2020. 1
- 182 [3] Gavin Brown and Ludmila I Kuncheva. “good” and “bad” diversity in majority vote ensembles.  
183 In *International workshop on multiple classifier systems*, pages 124–133. Springer, 2010. 2
- 184 [4] Thomas G Dietterich. Ensemble methods in machine learning. In *International workshop on*  
185 *multiple classifier systems*, pages 1–15. Springer, 2000. 2
- 186 [5] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Founda-*  
187 *tions and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014. 1
- 188 [6] Alistair Johnson, Tom Pollard, and Roger Mark. Mimic-iii clinical database (version 1.4).  
189 *PhysioNet*, 10:C2XW26, 2016. 2
- 190 [7] Michael Kamp. *Black-Box Parallelization for Machine Learning*. PhD thesis, Rheinische  
191 Friedrich-Wilhelms-Universität Bonn, Universitäts-und Landesbibliothek Bonn, 2019. 2
- 192 [8] Michael Kamp, Jonas Fischer, and Jilles Vreeken. Federated learning from small datasets. *arXiv*  
193 *preprint arXiv:2110.03469*, 2021. 6
- 194 [9] Balaji Krishnapuram, David Williams, Ya Xue, Lawrence Carin, Mário Figueiredo, and Alexan-  
195 der Hartemink. On semi-supervised classification. *Advances in neural information processing*  
196 *systems*, 17, 2004. 1
- 197 [10] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report,  
198 University of Toronto, Toronto, 2009. 3, 4
- 199 [11] Chuan Ma, Jun Li, Ming Ding, Howard H Yang, Feng Shu, Tony QS Quek, and H Vincent Poor.  
200 On safeguarding privacy and security in the framework of federated learning. *IEEE network*, 34  
201 (4):242–248, 2020. 1
- 202 [12] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.  
203 Communication-efficient learning of deep networks from decentralized data. In *Artificial*  
204 *Intelligence and Statistics*, pages 1273–1282, 2017. 1, 2
- 205 [13] Sasi Kumar Murakonda and Reza Shokri. MI privacy meter: Aiding regulatory compliance by  
206 quantifying the privacy risks of machine learning. *arXiv preprint arXiv:2007.09339*, 2020. 4

- 207 [14] Milad Nasr, Reza Shokri, and Amir Houmansadr. Comprehensive privacy analysis of deep  
208 learning: Passive and active white-box inference attacks against centralized and federated  
209 learning. In *2019 IEEE symposium on security and privacy (SP)*, pages 739–753. IEEE, 2019. 4
- 210 [15] Nicolas Papernot, Shuang Song, Ilya Mironov, Ananth Raghunathan, Kunal Talwar, and Ulfar  
211 Erlingsson. Scalable private learning with pate. In *International Conference on Learning*  
212 *Representations*, 2018. 6
- 213 [16] Constantin Marc Seibold, Simon Reiß, Jens Kleesiek, and Rainer Stiefelhagen. Reference-  
214 guided pseudo-label generation for medical semantic segmentation. In *Proceedings of the AAAI*  
215 *Conference on Artificial Intelligence*, volume 36, pages 2171–2179, 2022. 6
- 216 [17] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference  
217 attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy*  
218 *(SP)*, pages 3–18. IEEE, 2017. 1, 3
- 219 [18] Cathie Sudlow, John Gallacher, Naomi Allen, Valerie Beral, Paul Burton, John Danesh, Paul  
220 Downey, Paul Elliott, Jane Green, Martin Landray, et al. Uk biobank: an open access resource  
221 for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS*  
222 *medicine*, 12(3):e1001779, 2015. 2
- 223 [19] Ziteng Sun, Peter Kairouz, Ananda Theertha Suresh, and H Brendan McMahan. Can you really  
224 backdoor federated learning? *arXiv preprint arXiv:1911.07963*, 2019. 4
- 225 [20] Kang Wei, Jun Li, Ming Ding, Chuan Ma, Howard H Yang, Farhad Farokhi, Shi Jin, Tony QS  
226 Quek, and H Vincent Poor. Federated learning with differential privacy: Algorithms and  
227 performance analysis. *IEEE Transactions on Information Forensics and Security*, 15:3454–  
228 3469, 2020. 4, 5
- 229 [21] Ligeng Zhu and Song Han. Deep leakage from gradients. In *Federated learning*, pages 17–31.  
230 Springer, 2020. 1