
Model Extraction Attacks on Split Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Federated learning (FL) is a popular collaborative learning scheme involving
2 multiple clients and a server. FL focuses on client’s data privacy but exposes
3 interfaces for Model Extraction (ME) attacks. As FL periodically collects and
4 shares model parameters, a malicious client can download the latest model and thus
5 steal model Intellectual Property (IP). Split Federated Learning (SFL), a recent
6 variant of FL, splits the model into two, giving one part of the model (client-
7 side model) to clients, and the remaining part (server-side model) to the server.
8 While SFL was primarily designed to facilitate training on resource-constrained
9 devices, it prevents some ME attacks by blocking prediction queries. In this work,
10 we expose the vulnerability of SFL and show how ME attacks can be launched
11 by malicious clients querying the gradient information from server-side. We
12 propose five ME attacks that differ in the gradient usage in data crafting, generating,
13 gradient matching and soft-label crafting as well as in the attacker data availability
14 assumptions. We show that the proposed ME attacks work exceptionally well for
15 SFL. For instance, when the server-side model has five layers, our proposed ME
16 attack can achieve over 90% accuracy with less than 2% accuracy degradation with
17 VGG-11 on CIFAR-10.

18 1 Introduction

19 Federated Learning (FL) [McMahan et al., 2017] has become increasingly popular thanks to its ability
20 to protect users’ data privacy and comply with General Data Protection Regulation (GDPR) policy. In
21 FL, clients locally update their model copies, and the FL server collects them by averaging the model
22 parameters, then distributing the averaged model again to its clients. Such a setting only allows model
23 parameters to be shared with the server, and direct data sharing is avoided. One drawback of FL is
24 its clients need to train the entire model locally, which is usually challenging for resource-limited
25 edge devices. As the countermeasure, Split Federated Learning (SFL) scheme [Thapa et al., 2020] is
26 proposed as a variant of FL. In SFL, a neural network is split into a client-side model and a server-side
27 model, where the client-side model is shared among multiple clients and processed locally on their
28 devices. During training, clients offload the intermediate activations to server, where the heavy-duty
29 computation is performed at the heavy-duty server and the computed gradients are sent back to clients.
30 SFL follows the same model averaging routine as FL to synchronize the model. SFL avoids collecting
31 clients’ raw data and also reduces the computational overhead at the client-end.

32 In addition to the computation advantage introduced by SFL, it can also provide model IP protection
33 which is absent in FL. The high training cost of high-performance NN makes the NN model a valuable
34 Intellectual Property (IP). Unlike FL which handles the entire NN over to the clients, SFL preserves
35 the server-side model which prevents potential IP theft (Fig. 1 (a)). Moreover, according to our
36 investigation, SFL shows resistance to Model Extraction (ME) attack [Tramèr et al., 2016, Jagielski
37 et al., 2020]. In an ME attack, the model IP can be acquired by querying a publicly accessible

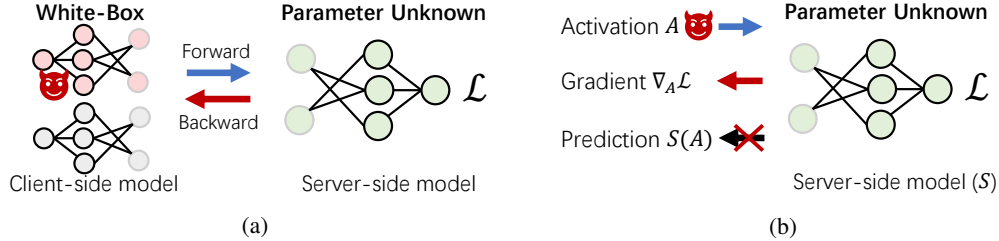


Figure 1: Model Extraction (ME) attack in SFL. (a) The attacker knows the client-side model parameters but does not know the server-side model parameters. (b) Server-side model allows gradient query access but does not allow prediction queries.

38 prediction API of the model. Prior ME attacks cannot be applied to SFL as the protocol does not
 39 allow prediction query access, as illustrated in Fig. 1 (b).

40 As prior ME attacks fail in SFL, we propose SFL-specific ME attacks in this work. Assuming
 41 the client-side model is a white box, the attacker posing as a participant client can get gradient
 42 information of inputs sent to the server. We propose five variants of ME attacks on SFL enumerated as
 43 Craft-ME, GAN-ME, GM-ME, Train-ME, and SoftTrain-ME. These ME attack variants extensively
 44 cover different gradient usage, including data crafting, data generating, gradient matching and
 45 soft label crafting, as well as assumptions on data such as no data (including randomly generated
 46 noise data) [Truong et al., 2021], only auxiliary data (out-of-distribution data) and training data (in-
 47 distribution data). We consider both train-from-scratch and fine-tuning SFL applications as gradient
 48 consistency is different. We benchmark the performance of five ME attacks on SFL for these two
 49 cases and show that ME attacks can succeed without any data in fine-tuning applications. We also
 50 show that ME performance strongly correlates with the #layers in the server-side model and fail when
 51 the #layers in the server-side model increase. However, increasing layers in server-side model reduces
 52 layers in client-side model, which compromises clients' privacy. Finally, we conclude that using L1
 53 regularization during training can improve SFL schemes' resistance to ME attacks. In summary, we
 54 make the following contributions:

- 55 • We define the unique threat model in SFL and propose five ME attacks. These attacks
 56 differ in the usage of gradients (data crafting, data generating, gradient matching and soft
 57 label crafting) and data assumptions (no data, only auxiliary data and training data). To our
 58 knowledge, this is the first work that studies ME attacks for SFL.
- 59 • We study the performance of the proposed attacks and find that even when the attacker has
 60 no data, the model can be extracted with high accuracy. With auxiliary data or with training
 61 data, ME attack performance can be further improved. For a 5-layer-in-server SFL, the
 62 strongest ME attack can derive a surrogate model with over 90% accuracy, and less than 2%
 63 accuracy degradation compared to the original VGG-11 on CIFAR-10 model.
- 64 • We find the ME attack performance decreases when more layers are present in server-side
 65 model. However, such a split model configuration compromises clients' privacy. We also
 66 show resistance to ME attack can be improved by regularizing the client-side model.

67 2 Related Work

68 **Model-split learning schemes.** The key idea for model-split learning schemes is to split the model
 69 so that part of it is processed in the client and the rest is offloaded to the server. This idea was first
 70 proposed in Kang et al. [2017], Teerapittayanon et al. [2017], Liu et al. [2018] for inference tasks.
 71 Gupta and Raskar [2018] then extended this idea for split learning, a collaborative multi-client neural
 72 network training. However, the round-robin design needed clients to learn sequentially and thus has a
 73 huge disadvantage in terms of training time.

74 **Split federated learning (SFL).** In this paper, we consider the SFL scheme where the clients process
 75 their local models in parallel followed by periodic synchronization as in FedAvg [McMahan et al.,
 76 2017]. This is the SFL-V2 scheme that is introduced in Thapa et al. [2020] for its better accuracy
 77 performance. The detailed process is shown in Algorithm 1. At the beginning of each epoch, server
 78 performs the synchronization of client-side model and sends the updated version to all clients. Then,

Algorithm 1 Split Federated Learning

Require: For M clients, instantiate private training data $(\mathbf{X}_i, \mathbf{Y}_i)$ for $1, 2, \dots, M$. Server-side model S has N layers and client-side model C_i has $L - N$ layers.

```
1: initialize  $C_i, S$ 
2: for epoch  $t \leftarrow 1$  to num_epochs do
3:    $C^* = \frac{1}{M} \sum_{i=1}^M C_i$  {Model Synchronization}
4:    $C_i \leftarrow C^*$  for all  $i$ 
5:   for step  $s \leftarrow 1$  to num_batches do
6:     for client  $i \leftarrow 1$  to  $M$  in Parallel do
7:       data batch  $(\mathbf{x}_i, \mathbf{y}_i) \leftarrow (\mathbf{X}_i, \mathbf{Y}_i)$ 
8:        $\mathbf{A}_i = C_i(\mathbf{W}_{C_i}; \mathbf{x}_i)$  {Client forward; send  $\mathbf{A}_i$  to Server}
9:     end for
10:    for client  $i \leftarrow 1$  to  $M$  in Sequential do
11:       $\mathcal{L} = \mathcal{L}_{CE}(S(\mathbf{W}_S; \mathbf{A}_i), \mathbf{y}_i)$  {Server forward}
12:       $\nabla_{\mathbf{A}_i} \mathcal{L} \leftarrow$  back-propagation {Server backward, send  $\nabla_{\mathbf{A}_i} \mathcal{L}$  to Client}
13:      Update  $\mathbf{W}_S$ ;
14:    end for
15:    for client  $i \leftarrow 1$  to  $M$  in Parallel do
16:       $\nabla_{\mathbf{x}_i} \mathcal{L} \leftarrow$  back-propagation {Client backward}
17:      Update  $\mathbf{W}_{C_i}$ ;
18:    end for
19:  end for
20: end for
```

79 clients perform forward propagation locally till layer $L - N$ (the last layer of client-side model),
80 sending the intermediate activation \mathbf{A}_i to the server (line 8). Server accepts the activation and label
81 \mathbf{y}_i sent from clients, and uses them to calculate the loss and initiates the backward process (line 9).
82 The backward process (line 10) consists of several steps: server performs backward propagation on
83 the loss, updates server-side model and sends back gradient $\nabla_{\mathbf{A}_i} \mathcal{L}$ to clients. Clients then continue
84 the backward propagation on their client-side model copies and perform model updates accordingly.
85 While FedGKT [He et al., 2020] leaks prediction logits to clients, SFL does not, making it a promising
86 candidate against ME attacks.

87 **Model extraction (ME) attack.** ME attack targets model prediction service APIs and retrieves
88 confidence score or prediction label for given inputs. The vulnerability of a DNN model to ME is first
89 shown in Tramèr et al. [2016]. Jagielski et al. [2020] shows that high fidelity model extraction can
90 be achieved with fewer queries and the surrogate model can be used for launching more successful
91 adversarial attacks. A successful ME attack not only breaches the model IP, but also makes the model
92 more vulnerable to attacks. ME attack also supports transferable adversarial attacks [Goodfellow
93 et al., 2014], mainly targeted ones [Madry et al., 2017] against the victim model. It can also be used
94 to perform bit-flip attacks [Rakin et al., 2019]; with hardware expertise, a few bit flips on model
95 parameters can degrade ResNet-18 model accuracy to below 1%.

96 **Data privacy in SFL.** Similar to FL, SFL scheme also has data privacy concerns. The most serious
97 one is its vulnerability to MI attacks. In model-based MI attack [Fredrikson et al., 2015], the attacker
98 trains an inverted version of client-side model and can directly reconstruct raw inputs from the
99 intermediate activation. Recent works [Vepakomma et al., 2020, Li et al., 2022] point out this
100 vulnerability and provide practical ways to mitigate MI. However, defenses only work well if the
101 client-side model has enough number of layers.

102 3 Threat Model

103 3.1 Attacker assumptions.

104 **Objectives.** According to Jagielski et al. [2020], there are three model extraction (ME) attack
105 objectives: i) functional equivalence, ii) high accuracy, and iii) high fidelity. However, achieving

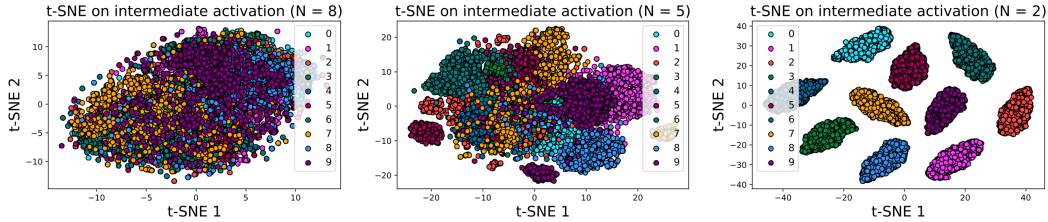


Figure 2: tSNE analysis of intermediate activation with different N (8, 5, 2 from left to right) on a VGG11 model on CIFAR-10 dataset. N denotes number of layers in server-side model in SFL.

106 functional equivalence is difficult in practical applications. Hence, most of the existing practical
 107 ME attacks on NN models focus on achieving high accuracy and fidelity. To achieve the accuracy
 108 goal, the attacker wants to obtain a model that maximizes the prediction correctness. To achieve the
 109 fidelity goal, the attacker wants to derive a model with a similar decision boundary as the victim
 110 model before launching adversarial attacks [Biggio et al., 2013].

111 **Capabilities.** We assume the attacker acts as a client in a multi-client SFL scheme against the
 112 server (model owner). Consider the SFL scheme outlined in Fig. 1. We assume the entire model has a
 113 total of L layers (or layer-like blocks, i.e. BasicBlock in ResNet) out of which the server processes N
 114 layers. The attacker holds *white-box assumption* on the client-side model (consists of $L - N$ layers),
 115 that is, it knows the exact model architecture and parameters for those layers. The attacker holds
 116 a *grey-box assumption* on the N -layer server-side model, that is, it knows its architecture and loss
 117 function while the model parameters are unknown. Also, we assume *server blocks the prediction*
 118 *queries* thus neither logits nor prediction labels are accessible by clients during training, but server
 119 *allows gradient queries* to let client-side models be updated. Based on a client’s activation $\mathbf{A} = C(\mathbf{x})$
 120 and its label \mathbf{y} , gradient information $\nabla_{\mathbf{A}} \mathcal{L}$ is computed and sent back to clients. In addition, we
 121 assume the attacker can perform gradient queries on any input \mathbf{x} , including malicious ones.

122 3.2 Analysis

123 **Partial model extraction problem.** Since the attacker already has a white-box assumption of the
 124 client-side model, the attacker only needs to extract the server-side model to reveal the entire model.
 125 This results in an easier problem setting. We observe client-side model heavily regularizes the feature
 126 space of its output (input of server-side model), making ME attacks easier to succeed, especially
 127 when N is small. As shown in Fig. 2, on a VGG-11 [Simonyan and Zisserman, 2014] model on
 128 CIFAR-10 dataset, as N becomes smaller, tSNE embeddings of intermediate features with different
 129 labels are easier to distinguish. For $N = 2$, ME attack is as simple as separating different clusters
 130 with a linear layer.

131 To show that the extracting part of the model with ME is easier, we study existing ME attacks on
 132 the server-side model, **by assuming that prediction access is allowed**. Specifically, we investigate
 133 CopyCat CNN [Correia-Silva et al., 2018], Knockoff-random [Orekondy et al., 2019] and data-free
 134 ME [Truong et al., 2021]. As shown in Fig. 3 (a), with auxiliary data (CIFAR-100) and enough query
 135 budget, both attacks derive a surrogate model with very high accuracy even for a large N setting.
 136 Moreover, attacker with no data can also succeed with data-free ME as shown in Fig. 3 (a). When the
 137 query budget is equal to 2 million, the data-free ME can extract the model with high accuracy even
 138 when N is equal to 5.

139 **Consistency of gradient query.** For fine-tuning applications [Park et al., 2021], attackers get
 140 consistent gradient information from gradient query, as server-side model parameters are frozen
 141 or updated with a very small learning rate. Gradient consistency is very beneficial for ME attacks.
 142 However, for a training-from-scratch usage, queries to SFL model obtain inconsistent gradient
 143 information as the server-side model drastically changes during training. As shown in Fig. 3 (c), for
 144 the same query input, the gradient is drastically different in different epochs.

145 4 Proposed Model Extraction Attack

146 Previously, Milli et al. [2019] demonstrated that using gradients to reveal one-layer linear transforma-
 147 tion is trivial. Given $f(\mathbf{x}) = \mathbf{W}^T \mathbf{x}$, one can directly infer \mathbf{W} from a single gradient query given that

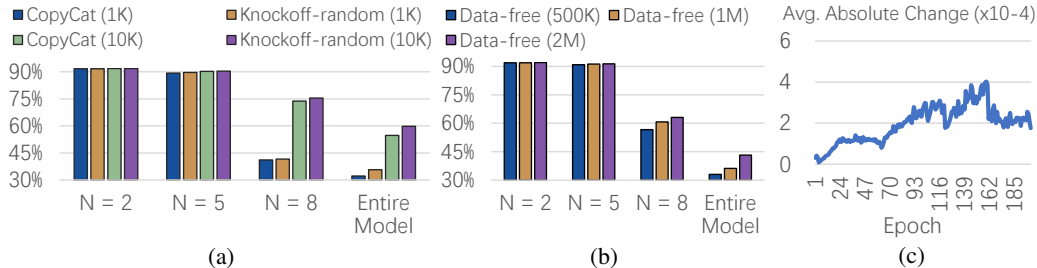


Figure 3: Analyses of SFL threat model. Existing ME attack performance on VGG-11 CIFAR-10 model with different N , assuming prediction query access is allowed: (a) ME attacks using CIFAR-100 as auxiliary dataset; (b) Data-free ME attack that demonstrates ME attack on part of the model is much easier than ME attack on the entire model; (c) Inconsistent gradient problem in training-from-scratch SFL. The y-axis denotes the change in gradient (lower means more consistent) for the same inputs in different epochs.

148 $\mathbf{W}^T = \nabla_{\mathbf{x}} f(\mathbf{x})$. However, using gradient only can go no further than one layer. Milli et al. [2019]
 149 shows that to recover a two-layer ReLU network of the form $f(\mathbf{x}) = \sum_{n=1}^h g(\mathbf{x})_i \mathbf{W}_i \mathbf{A}_i^T \mathbf{x}$, where
 150 $g(\mathbf{x}) = \mathbb{1}\{\mathbf{A}\mathbf{x} > 0\}$, \mathbf{A} is of $\mathbb{R}^{h \times d}$ and \mathbf{W} is of \mathbb{R}^h , using input gradient can recover the absolute
 151 value of normal vectors $|\mathbf{W}_i \mathbf{A}_i|$ for $i \in [h]$. In order to get the sign information of $\mathbf{W}_i \mathbf{A}_i$, prediction
 152 query is required which is not supported by SFL’s threat model.

153 So in this paper, we investigate approximate ME attacks that differ in the data assumptions and
 154 gradient usage. We propose five ME attacks as shown in Table 1. We also include a naive baseline
 155 scheme that train the surrogate model from scratch without access to either the client-side model or
 156 gradients. Despite the differences, the five proposed attack methods follow the same strategy, that is,
 157 they all train a randomly initialized surrogate server-side model from scratch. We provide a detailed
 158 illustration of proposed five attacks in Appendix A.2.

Table 1: Model Extraction Attack Methods in SFL

Method	Data Assumption	Prediction Query	Gradient Usage	Client-side model Usage
Craft-ME	None	No	Data Crafting	Initialization
GAN-ME	None	No	Data Generator	Initialization
GM-ME	Natural Auxiliary	No	Gradient Matching	Initialization
Train-ME	Limited Training	No	None	Initialization
SoftTrain-ME	Limited Training	No	Soft Label Crafting	Initialization
Naive Baseline	Limited Training	No	None	None

159 4.1 ME attacks without training data

160 We first consider two cases where the attacker does not have training data: one is that attacker does
 161 not use any data or use randomly generated noise data. This case is motivated by Truong et al. [2021]
 162 which shows that when an attacker does not have a *similar enough* dataset, it is better not to use it at
 163 all. The other case is when the attacker has an auxiliary dataset with different labels from the victim’s
 164 training data (for example, CIFAR-10 and CIFAR-100).

165 **Crafting model extraction (Craft-ME).** Inspired by Han et al. [2018], where data-label pairs
 166 (referred to as instances) with small-loss are shown to present useful guidance for knowledge
 167 distillation, we propose a simple method to craft small-loss instances using gradient queries and use
 168 them to train the surrogate model. We initialize random input \mathbf{x}_r for every class label c , and use the
 169 gradient $\nabla_{\mathbf{x}_r} \mathcal{L}$ to update \mathbf{x}_r . For each input, updating is repeated for a number of steps. By varying
 170 label c , a collection of small-loss instances is derived during SFL training. Then, a surrogate model is
 171 trained from scratch on these small-loss instances.

172 **GAN-based model extraction (GAN-ME).** Recent work [Truong et al., 2021] proposes a GAN-
 173 based approach for data-free ME. The key idea is to use a generator G to continually feed fake inputs
 174 to the victim model V and surrogate model S , and use confidence score matching to let S approach
 175 V . However, the confidence score matching needs prediction query which is not allowed in our case.

176 Thus, we adapt the GAN-based method to gradient-query-only case and propose a two-step method:
 177 First, a conditional-GAN (c-GAN) model $G(z|c)$ is initialized. The attacker trains the generator
 178 during victim model’s training by generating fake data x_f and label c and performing gradient queries
 179 to update G . After the training is done, generator G is used to supply small-loss instances (x_f, c) to
 180 train the surrogate model (the unknown part). We observe a serious mode collapse problem during
 181 the GAN training. So we utilize the distance-aware training introduced in Yang et al. [2019], to
 182 encourage the c-GAN to generate more diverse small-loss instances. In the new method, the training
 183 of the generator is not based on a min-max game, or on traditional GAN training. Instead, it simply
 184 trains the generator toward minimizing cross-entropy loss. While the generator G fails to generate
 185 natural-looking inputs even upon convergence, it generates abundant small-loss instances for every
 186 label, and divergence loss helps it generate a good variety of outputs. During SFL training, the
 187 generator can adjust itself to the changing server-side model.

188 **Gradient matching model extraction (GM-ME).** Gradient matching (GM) in ME attack has been
 189 investigated in Jagielski et al. [2020], Milli et al. [2019] and is used in combination with prediction
 190 query to improve the extraction performance. Since, in SFL, prediction query is not allowed, we
 191 navigate this strict threat model’s restriction by adopting gradient matching (GM) loss. For a given
 192 label y_i , GM loss has the following form:

$$\mathcal{L}_{GM} = |\nabla_{x_i} \mathcal{L}(S(C(x_i)), y_i) - \nabla_{x_i} \mathcal{L}(V(C(x_i)), y_i)|_2^2 \quad (1)$$

193 where, x_i denote inputs, C denotes client-side model, S and V denotes the surrogate model and
 194 victim model, respectively. For each input, attacker would query gradients with different label y_i to
 195 get as much information as possible. This attack performs extremely well for small N but degrades
 196 significantly for a larger N . Its performance also depends on the domain similarity between the
 197 auxiliary dataset and the victim dataset.

198 4.2 ME attacks with training data.

199 Next we discuss the case when the attacker has a subset of training data, corresponding to the strongest
 200 data assumption .

201 **Training-based model extraction (Train-ME).** For attackers with a subset of the training data,
 202 derivation of an accurate surrogate model can be done using supervised learning (through minimizing
 203 the cross entropy loss on the available data). We call this Train-ME, similar idea is also adopted in Fu
 204 et al. [2022] to extract the entire model of the other party. Train-ME only relies on the white-box
 205 assumption of the client-side model, using it to initialize the surrogate model and does not need to
 206 use the gradient query at all. Surprisingly, it is one of the most effective ME attacks.

207 **Gradient-based soft label training model extraction (Soft-train-ME).** If gradient query is allowed
 208 and a subset of training data is available, the attacker can achieve better ME attack performance
 209 compared to Train-ME. To utilize gradients, a naive idea is to combine the GM loss with cross-entropy
 210 loss in Train-ME. However, our initial investigation shows they are not compatible; the cross-entropy
 211 loss term usually dominates and the GM loss even hurts the performance. An alternative approach is
 212 to use soft label. We build upon the method in Gu et al. [2020] which shows that gradient information
 213 of incorrect labels is beneficial in knowledge distillation, and use it for surrogate model training.
 214 Specifically, for each input x_i , gradients of the ground truth label as well as incorrect labels are
 215 collected (a total N_C , where N_C is the number of classes). For an input x_i with true label c , its soft
 216 label q_i^k of k -th ($k \neq c$) label is computed as follows:

$$q_i^k = (1 - \alpha) * \frac{\cos(e^k, e^c)}{\sum_{m=1, m \neq c}^{N_C} (\cos(e^m, e^c) + 1)} \quad (2)$$

217 where, e^k denotes flattened gradients of label k , q_i^k denotes soft label for the k -th label k and α is a
 218 constant ($\alpha > 0.5$). The derived (x_i, q_i) pair is then used in the surrogate model training in addition
 219 to the true label c (which is the only difference from the Train-ME).

220 **5 Model Extraction Performance**

221 In this section, we demonstrate the performance of the proposed ME attacks and the baseline attack
 222 on SFL schemes. All experiments are conducted on a single RTX-3090 GPU. For the SFL model
 223 training, we set the total number of epochs to 200, and use SGD optimizer with a learning rate of
 224 0.05 and learning rate decay (multiply by factor of 0.2 at epochs 60, 120 and 160). We assume all
 225 clients participate in every epoch with an equal number of training steps. To perform ME attacks,
 226 the attacker uses an SGD optimizer with a learning rate of 0.02 to train the surrogate model and we
 227 report the best accuracy and fidelity. We evaluate accuracy of the surrogate model on the validation
 228 dataset. We use the label agreement as fidelity, defined as the percentage of samples that the surrogate
 229 and victim models agree with over the entire validation dataset, as in Jagielski et al. [2020]. Detailed
 230 settings for each ME attack are described in Appendix A.1.

231 **5.1 ME attack on SFL with fine-tuning based training**

232 We first perform the proposed ME attacks on fine-tuning SFL version with consistent gradient query.
 233 Here we use a pre-trained model and set the number of gradient queries to 100K. On a victim VGG-11
 234 model on CIFAR-10 dataset, whose original accuracy is 91.89%, performance of all five ME attacks
 235 are shown in Table 2. For each of the ME attacks, we vary hyper-parameters and report the one
 236 that achieves the best attack performance; details on different hyper-parameters are included in
 237 Appendix A.3. When $N = 2$, all five ME attacks can achieve near-optimal accuracy and fidelity
 238 performance. For Craft, GAN and GM ME, the accuracy drops to around 80% when N is 5, it sharply
 239 drops to below 40% when N is 8. For Train and SoftTrain ME, accuracy slightly degrades when N is
 240 5, and reduces to around 70% when N is 8. These results show that ME attack performance strongly
 241 correlates with the #layers in server-side model. With an increasing #layers in server-side model, ME
 242 attack performance reduces as the extraction problem becomes harder with more unknown parameters
 243 and more complicated input feature space.

244 **Observation.** Different attacks present different and interesting characteristics. Craft-ME has a
 245 steady attack performance and can succeed even with a tight query budget. GAN-ME needs a large
 246 query budget to train the c-GAN generator towards convergence but can achieve better accuracy
 247 and fidelity than Craft-ME for $N \leq 5$. GM-ME requires an auxiliary dataset that is similar to the
 248 training data. If CIFAR-100 is used to attack CIFAR-10 model, GM-ME achieves almost perfect
 249 extraction for small N . However, its performance degrades if MNIST or SVHN are used as auxiliary
 250 datasets. For a high N , the surrogate model fails to converge on the GM loss, and its extraction
 251 performance suffers from a sharp drop. For attacks with training data such as Train and SoftTrain
 252 MEs, both accuracy and fidelity are much higher than attacks without training data. When $N \geq 6$,
 253 SoftTrain-ME can achieve slightly better accuracy and fidelity than Train-ME.

Table 2: ME attack performance on SFL on fine-tuning and training-from-scratch applications. The victim is a VGG-11 model on CIFAR-10 with 91.89% validation accuracy. For Train, SoftTrain and Naive baseline, for the fine-tuning setting, data assumption is 1K training data (randomly sampled), and for the train-from-scratch setting, the number of clients is 10 and each client has 5K training data.

Metric	N	Fine-tuning						Training-from-scratch					
		Craft	GAN	GM	Train	SoftTrain	Naive	Craft	GAN	GM	Train	SoftTrain	Naive
Accuracy (%)	2	91.64	91.86	92.02	92.05	91.99		85.99	85.99	53.06	90.58	90.31	
	5	83.46	84.93	80.28	90.82	90.48	49.64	35.58	40.03	12.13	89.86	87.02	72.63
	8	35.48	18.82	12.45	70.28	71.32		15.34	17.49	10.88	78.64	56.78	
Fidelity (%)	2	98.23	98.42	99.87	99.29	99.10		92.37	89.59	54.63	99.34	98.87	
	5	86.32	87.49	84.33	94.84	94.67	50.62	41.32	38.72	11.87	95.40	89.83	72.62
	8	36.11	18.62	12.63	71.79	72.45		15.63	17.44	10.67	80.01	57.78	

254 **5.2 ME attack on training-from-scratch SFL**

255 Next, we investigate the proposed ME attack performance in training-from-scratch SFL case. A good
 256 attack-time-window for gradient-based ME attacks is at the end of training when gradients do not vary
 257 as much and the model converges. So for Craft, GAN and GM-ME, we launch the attack at epoch
 258 160 to get more consistent gradients. As the model is updated by multiple clients, the percentage of

malicious clients also affects the ME attack performance. We found that with more benign clients (or fewer malicious inputs), the server-side model returns more consistent gradients to the attacker. Attack performance for three attacks are shown in Table 2 for 10-client SFL training-from-scratch case. Results with other hyperparameter settings are included in Appendix A.4.

Observation. Because of the poisoning effect of malicious inputs sent by the attacker in gradient-based attacks, in all the methods except Train-ME, the victim model’s accuracy suffers from a 2-3% degradation, resulting in a less accurate surrogate model. For gradient-based attacks without training data (Craft, GAN and GM MEs), we notice a sharp ($> 3\%$) performance drop in both accuracy and fidelity compared to the consistent query case. The sharp drop in ME attack performance is caused by *inconsistent gradients*. Take Craft-ME as an example. Crafted inputs that have a small loss at an earlier epoch of the training can have a large loss in the final model. Training surrogate model with large amount of inconsistent information results in poor training accuracy. Compared to Craft-ME, GAN-ME is more robust to inconsistent gradients as the generator can adjust itself to the change of server-side model, resulting in a better extraction performance. However, when N is larger, the generator does not converge well and its performance drops drastically. Last but not the least, GM-ME completely fails with inconsistent gradients, even for small N . This implies that the GM loss is super sensitive to inconsistent gradients and is only effective in consistent query cases. A comparison of SoftTrain-ME and Train-ME shows that SoftTrain’s advantage diminishes because of the poisoning effect. Thus, Train-ME is more effective for an attacker with training data in training-from-scratch SFL.

5.3 ME attack performance for other architectures and datasets

We perform extensive analysis on the attack performance of the proposed ME schemes on other architectures and datasets. We use Train-ME attack, the best performer for both fine-tuning and training-from-scratch cases and use the accuracy as the performance measure. To test ME attack on different datasets, we perform Train-ME attack with 1K training data on VGG-11 model with N set to 5, on MNIST [LeCun, 1998], FMNSIT [Xiao et al., 2017], SVHN [Netzer et al., 2011] and CIFAR-100 datasets [Krizhevsky et al., 2009] (in addition to CIFAR-10). As shown in Fig. 4 (a), for all datasets except CIFAR-100, ME attack achieves accuracy very close to the original. For CIFAR-100, the extracted accuracy is over 20% below the original because of its task difficulty. Additionally, we also test Train-ME performance with 2% and 20% ImageNet training data of Mobilenet-V2 on ImageNet dataset. As shown in Fig. 4 (b), ME attack fails badly due to the complexity of ImageNet dataset [Deng et al., 2009], resulting in a high accuracy gap of 10% when N is set to 2. For different architectures, we choose Resnet-20, Resnet-32 [He et al., 2016] and Mobilenet-V2 [Sandler et al., 2018] on CIFAR-10 dataset (with necessary adaptations). For Resnet and Mobilenet family, we assign last 4 layer-blocks and 1 FC layer to server-side model. As shown in Fig. 4 (c), with the same proportion of layers (5 out of 11) being assigned to server-side model, ME attack is much less effective on Resnet-20 than on VGG-11. A comparison of the performance of Resnet-32 and Mobilenetv2 with similar proportion of layers being assigned to server-side (5 out of 17 and 20, respectively), ME on Resnet-32 is also much worse than on MobilenetV2. This indicates Resnet architecture is more resistant to ME attack.

6 Discussion

Our evaluation showed that ME attack performance drops with increasing #layers in server-side model. Thus, a simple idea to improve resistance to ME attack is to use a larger N . However this implies that the #layers in client-side model would be smaller, thereby undermining clients’ data privacy. Data privacy in SFL can be represented by the Mean Square Error (MSE) performance of an MI attack as outlined in Li et al. [2022]. The implementation detail is included in Appendix B.1. As shown in Fig. 4 (d), extracted accuracy decreases with a larger N , and MSE decreases too. This means MI attack can provide more precise reconstruction which compromises data privacy. We will extensively consider the tradeoff between ME resistance and data privacy in future.

Defense Method. We also investigate defensive methods against ME attacks on SFL. According to the analysis in Fig. 2, extraction is easier because client-side model is leaked to the attacker. Thus, to make ME attacks more difficult, we can restrict the client-side model’s feature extraction capabilities. Towards this goal, we apply L1 regularization with three different strength ($\lambda = 5e-5, 1e-4$ and $2e-4$)

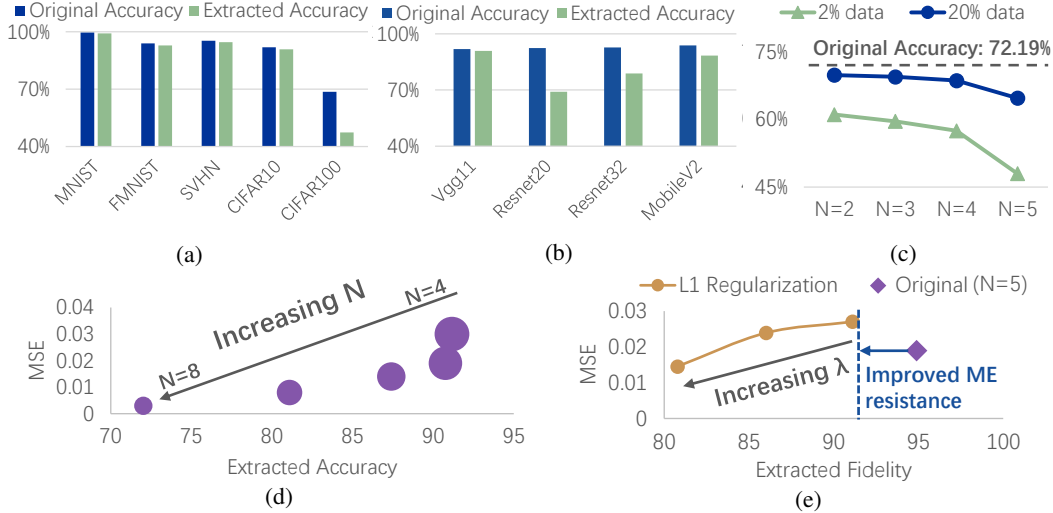


Figure 4: **Top row:** Extensive results on model extraction (ME) attacks. (a) ME attack performance of VGG-11 on other datasets. (b) ME attack performance of MobilenetV2 on ImageNet. (c) ME attack performance of other architectures on CIFAR-10 dataset. **Bottom row:** (d) Tradeoff between ME resistance and data privacy (MSE). (e) L1 regularization as effective defense for ME attacks.

312 on the client-side model to penalize its weight magnitude. As shown in Fig. 4 (e), this simple defense
 313 effectively improves the resistance to ME attack (specifically, Train-ME with 1K data) without hurting
 314 data privacy. While there is some accuracy degradation, as shown in Appendix A.5, this demonstrates
 315 the potential of using regularization to defend ME attacks.

316 7 Ablation Study

317 **ME attack with non-IID data.** We consider the non-IID (independent and identically distributed)
 318 case where the attacker only has data from C classes of CIFAR-10. Results presented in Appendix A.6
 319 show that ME attack performance is still good for $C = 5$ but degrades sharply when $C = 2$.

320 **Adversarial attack based on successful ME attack.** As mentioned before, the goal of ME attack
 321 is to launch more successful adversarial attacks. We perform transfer adversarial attacks using a
 322 surrogate model extracted by the strongest Train-ME attack. As shown in Appendix A.7, SFL with
 323 proper N achieves better resistance to adversarial attacks.

324 **ME attack without architecture information.** In Appendix A.8, we investigate simple variants
 325 (longer, shorter, wider, and thinner) of the original architecture as the surrogate model architecture.
 326 We find that the performance of ME attacks is similar for the different architectures – the exception is
 327 GM-ME which fails for different surrogate architectures.

328 8 Conclusion

329 In this work, we study the model IP protection capability of SFL and its resistance to model extraction
 330 (ME) attack. We propose five viable ME attack methods for the threat model where gradient query
 331 is allowed but prediction query is not allowed. For the case when there are enough number of
 332 layers in server-side model, model IP can be protected well and transfer adversarial attack is not
 333 successful. However, data privacy could be compromised and so this factor needs to be considered
 334 in the development of such schemes. We also point out a possible way of defending such attacks
 335 through regularization and plan to expand on it in the near future.

336 **Broader Impact.** This work points out the vulnerability of Split Federated Learning (SFL), to model
 337 extraction attacks, and should prevent a naive adoption of SFL as a model IP protection method. We
 338 believe that the attacks presented here would initiate research in the development of defense schemes
 339 to mitigate such attacks, help design a more robust SFL and possibly help in the design of neural
 340 network models that are inherently resilient to such attacks.

References

- 341
342 Battista Biggio, Igino Corona, Davide Maiorca, Blaine Nelson, Nedim Šrndić, Pavel Laskov, Giorgio Giacinto,
343 and Fabio Roli. Evasion attacks against machine learning at test time. In *Joint European conference on*
344 *machine learning and knowledge discovery in databases*, pages 387–402. Springer, 2013.
- 345 Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMahan, Virginia
346 Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097*,
347 2018.
- 348 Jacson Rodrigues Correia-Silva, Rodrigo F Berriel, Claudine Badue, Alberto F de Souza, and Thiago Oliveira-
349 Santos. Copycat cnn: Stealing knowledge by persuading confession with random non-labeled data. In *2018*
350 *International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2018.
- 351 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
352 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee,
353 2009.
- 354 Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit confidence informa-
355 tion and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and*
356 *Communications Security*, pages 1322–1333, 2015.
- 357 Chong Fu, Xuhong Zhang, Shouling Ji, Jinyin Chen, Jingzheng Wu, Shanqing Guo, Jun Zhou, Alex X Liu, and
358 Ting Wang. Label inference attacks against vertical federated learning. In *31st USENIX Security Symposium*
359 *(USENIX Security 22)*, Boston, MA, 2022.
- 360 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples.
361 *arXiv preprint arXiv:1412.6572*, 2014.
- 362 Jindong Gu, Zhiliang Wu, and Volker Tresp. Introspective learning by distilling knowledge from online
363 self-explanation. In *Proceedings of the Asian Conference on Computer Vision*, 2020.
- 364 Otkrist Gupta and Ramesh Raskar. Distributed learning of deep neural network over multiple agents. *Journal of*
365 *Network and Computer Applications*, 116:1–8, 2018.
- 366 Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama.
367 Co-teaching: Robust training of deep neural networks with extremely noisy labels. *Advances in neural*
368 *information processing systems*, 31, 2018.
- 369 Chaoyang He, Murali Annavaram, and Salman Avestimehr. Group knowledge transfer: Federated learning of
370 large cnns at the edge. *Advances in Neural Information Processing Systems*, 33:14068–14080, 2020.
- 371 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In
372 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- 373 Matthew Jagielski, Nicholas Carlini, David Berthelot, Alex Kurakin, and Nicolas Papernot. High accuracy and
374 high fidelity extraction of neural networks. In *29th USENIX Security Symposium (USENIX Security 20)*,
375 pages 1345–1362, 2020.
- 376 Yiping Kang, Johann Hauswald, Cao Gao, Austin Rovinski, Trevor Mudge, Jason Mars, and Lingjia Tang.
377 Neurosurgeon: Collaborative intelligence between the cloud and mobile edge. *ACM SIGARCH Computer*
378 *Architecture News*, 45(1):615–629, 2017.
- 379 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 380 Yann LeCun. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- 381 Jingtao Li, Adnan Siraj Rakin, Xing Chen, Zhezhi He, Deliang Fan, and Chaitali Chakrabarti. Rssfl: A
382 resistance transfer framework for defending model inversion attack in split federated learning, 2022.
- 383 Peng Liu, Bozhao Qi, and Suman Banerjee. Edgeeye: An edge service framework for real-time intelligent video
384 analytics. In *Proceedings of the 1st international workshop on edge systems, analytics and networking*, pages
385 1–6, 2018.
- 386 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep
387 learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- 388 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-
389 efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pages
390 1273–1282. PMLR, 2017.

- 391 Smitha Milli, Ludwig Schmidt, Anca D Dragan, and Moritz Hardt. Model reconstruction from model ex-
 392 planations. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 1–9,
 393 2019.
- 394 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in
 395 natural images with unsupervised feature learning. 2011.
- 396 Tribhuvanesh Orekondy, Bernt Schiele, and Mario Fritz. Knockoff nets: Stealing functionality of black-box
 397 models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
 398 4954–4963, 2019.
- 399 Sangjoon Park, Gwanghyun Kim, Jeongsol Kim, Boah Kim, and Jong Chul Ye. Federated split task-agnostic
 400 vision transformer for covid-19 cxr diagnosis. *Advances in Neural Information Processing Systems*, 34, 2021.
- 401 Adnan Siraj Rakin, Zhezhi He, and Deliang Fan. Bit-flip attack: Crushing neural network with progressive bit
 402 search. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 1211–1220,
 403 2019.
- 404 Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2:
 405 Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and
 406 pattern recognition*, pages 4510–4520, 2018.
- 407 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition.
 408 *arXiv preprint arXiv:1409.1556*, 2014.
- 409 Surat Teerapittayanon, Bradley McDanel, and Hsiang-Tsung Kung. Distributed deep neural networks over
 410 the cloud, the edge and end devices. In *2017 IEEE 37th international conference on distributed computing
 411 systems (ICDCS)*, pages 328–339. IEEE, 2017.
- 412 Chandra Thapa, Mahawaga Arachchige Pathum Chamikara, Seyit Camtepe, and Lichao Sun. Splitfed: When
 413 federated learning meets split learning. *arXiv preprint arXiv:2004.12088*, 2020.
- 414 Florian Tramèr, Fan Zhang, Ari Juels, Michael K Reiter, and Thomas Ristenpart. Stealing machine learning
 415 models via prediction apis. In *25th {USENIX} Security Symposium ({USENIX} Security 16)*, pages 601–618,
 416 2016.
- 417 Jean-Baptiste Truong, Pratyush Maini, Robert J Walls, and Nicolas Papernot. Data-free model extraction. In
 418 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4771–4780,
 419 2021.
- 420 Praneeth Vepakomma, Abhishek Singh, Otkrist Gupta, and Ramesh Raskar. Nopeek: Information leakage
 421 reduction to share activations in distributed deep learning. In *2020 International Conference on Data Mining
 422 Workshops (ICDMW)*, pages 933–942. IEEE, 2020.
- 423 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine
 424 learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 425 Dingdong Yang, Seunghoon Hong, Yunseok Jang, Tianchen Zhao, and Honglak Lee. Diversity-sensitive
 426 conditional generative adversarial networks. *arXiv preprint arXiv:1901.09024*, 2019.

427 Checklist

428 The checklist follows the references. Please read the checklist guidelines carefully for information on
 429 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
 430 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
 431 the appropriate section of your paper or providing a brief inline description. For example:

- 432 • Did you include the license to the code and datasets? **[Yes]** See Section x.
- 433 • Did you include the license to the code and datasets? **[No]** The code and the data are
 434 proprietary.
- 435 • Did you include the license to the code and datasets? **[N/A]**

436 Please do not modify the questions and only use the provided macros for your answers. Note that the
 437 Checklist section does not count towards the page limit. In your paper, please delete this instructions
 438 block and only keep the Checklist section heading above along with the questions/answers below.

- 439 1. For all authors...
- 440 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
441 contributions and scope? [Yes]
- 442 (b) Did you describe the limitations of your work? [Yes] We provide the limitations of the
443 proposed attacking methods by experiments described in Section 5.2 and Section 5.3.
- 444 (c) Did you discuss any potential negative societal impacts of your work? [Yes] As
445 discussed in Section 1, this paper provides attacking methods to gain proprietary model
446 information which breaches Intellectual Property. However, we intend to point out
447 the vulnerability of SFL scheme, which makes entities aware of the danger and hence
448 prevents them from being attacked. We also provide potential defensive methods to
449 mitigate our proposed attacks.
- 450 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
451 them? [Yes]
- 452 2. If you are including theoretical results...
- 453 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
454 (b) Did you include complete proofs of all theoretical results? [N/A]
- 455 3. If you ran experiments...
- 456 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
457 mental results (either in the supplemental material or as a URL)? [Yes] We will include
458 them in the supplemental material.
- 459 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
460 were chosen)? [Yes] We will open-source our code which includes all the required
461 information.
- 462 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
463 ments multiple times)? [No]
- 464 (d) Did you include the total amount of compute and the type of resources used (e.g., type
465 of GPUs, internal cluster, or cloud provider)? [Yes]
- 466 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 467 (a) If your work uses existing assets, did you cite the creators? [Yes]
468 (b) Did you mention the license of the assets? [No]
469 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
470 (d) Did you discuss whether and how consent was obtained from people whose data you're
471 using/curating? [No]
472 (e) Did you discuss whether the data you are using/curating contains personally identifiable
473 information or offensive content? [No]
- 474 5. If you used crowdsourcing or conducted research with human subjects...=
- 475 (a) Did you include the full text of instructions given to participants and screenshots, if
476 applicable? [N/A]
477 (b) Did you describe any potential participant risks, with links to Institutional Review
478 Board (IRB) approvals, if applicable? [N/A]
479 (c) Did you include the estimated hourly wage paid to participants and the total amount
480 spent on participant compensation? [N/A]