

---

# Secure Split Learning against Property Inference and Data Reconstruction Attacks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Split learning of deep neural networks (SplitNN) has provided a promising solu-  
2 tion to learning jointly for the mutual interest of a guest and a host, which may  
3 come from different backgrounds, holding features partitioned vertically. How-  
4 ever, SplitNN creates a new attack surface for the adversarial participant, holding  
5 back its practical use in the real world. By investigating the adversarial effects  
6 of two highly threatening attacks, i.e., property inference and data reconstruction,  
7 adapted from security studies of federated learning, we identify the underlying  
8 vulnerability of SplitNN. To prevent potential threats and ensure learning guaran-  
9 tees of SplitNN, we design a privacy-preserving tunnel for information exchange  
10 between the guest and the host. The intuition behind our design is to perturb the  
11 propagation of knowledge in each direction with a controllable unified solution.  
12 To this end, we propose a new activation function named  $R^3eLU$ , transferring pri-  
13 vate smashed data and partial loss into randomized responses in forward and back-  
14 ward propagations, respectively. Moreover, we give the first attempt to achieve  
15 a fine-grained privacy budget allocation scheme for SplitNN. The analysis of pri-  
16 vacy loss proves that our privacy-preserving SplitNN solution provides a tight  
17 privacy budget, while the experimental result shows that our solution outperforms  
18 existing solutions in attack defense and model usability.

## 19 1 Introduction

20 Collaborative learning enables participants from different backgrounds to learn jointly for mutual  
21 interests. A well-known collaborative learning paradigm is federated learning (FL) [25], focusing  
22 on the coordination of distributed participants. Meanwhile, another paradigm, split neural network  
23 learning (SplitNN for short) [19, 11, 4], is designed for vertically partitioned features. The emer-  
24 gence of SplitNN provides a promising solution to building cooperative models such as two-towers  
25 recommendation systems [37, 39]. By combining different features, SplitNN is supposed to be more  
26 expressive for tasks like user profiling and recommendation.

27 However, collaborative learning paradigms are faced with severe security issues. Roughly speak-  
28 ing, there are three kinds of known threats, inference attack [9, 30], reconstruction attack [17, 31],  
29 and poisoning attack [33, 18]. The inference attack discloses properties or membership information  
30 of data samples, while the reconstruction attack seeks to recover participants' private data samples.  
31 Unlike these two kinds of threats, the poisoning attack [33, 18] aims to put harmful data into collab-  
32 orative learning for malicious purposes rather than stealing private information. As a result, some  
33 defensive solutions have been proposed for federated learning. According to their techniques, these  
34 solutions can be classified into three categories: differential privacy solutions [36, 32], homomorphic  
35 encryption solutions [38, 21], and secure multiparty computation solutions [12, 27].

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

36 Unfortunately, defense solutions for SplitNN are barely discussed, but threatening attacks are con-  
37 tinuously emerging. We note that the workflow of SplitNN has a unique asymmetric design, which  
38 is quite different from FL. Therefore, most solutions for secure FL are not suitable for SplitNN.  
39 Secure multiparty computation solutions and homomorphic encryption solutions can achieve ideal  
40 data confidentiality in SplitNN. But the overhead introduced is still far away from practical uses.  
41 Thus, for the first attempt at a secure SplitNN solution, we will concentrate on privacy leakage is-  
42 sues caused by property inference and data reconstruction attacks in this paper because they share  
43 similar adversarial goals. Instead, poisoning attacks need to be studied separately [7, 24].

44 We note an inherent contradiction [6] between privacy preservation and model utility in SplitNN.  
45 To investigate the privacy leakage risk of SplitNN, we evaluate the adversarial effect of property  
46 inference and data reconstruction attacks against recommendation and image classification models  
47 built using SplitNN. The results in semi-honest and malicious settings show that both attacks have  
48 sufficiently high success rates when disclosing the privacy of the guest or the host. To secure SplitNN  
49 against these attacks [17, 30, 29], we give the first privacy-preserving SplitNN solution, stemming  
50 privacy leakage from either direction with dynamic privacy budget allocation.

51 **Contribution.** Through a thorough investigation of the privacy leakage issue in SplitNN, we con-  
52 firm that the best option of defensive solution for SplitNN is to construct a privacy-preserving tunnel  
53 between the split surfaces of the host and guest sides. To this end, we propose a novel activation  
54 function named R<sup>3</sup>eLU, responding to forward and backward propagations in a randomized manner.  
55 Furthermore, we propose a fine-grained privacy budget allocation scheme for SplitNN to achieve  
56 more efficient perturbations, dynamically allocating the privacy budget in spatial dimension (feature  
57 level) and temporal dimension (epoch level). The analysis of privacy loss shows that our solution  
58 provides a differential privacy guarantee regarding activations. The evaluation results regarding rec-  
59 ommendation and classification models show that our solution can outperform the existing solutions  
60 in privacy preservation and model usability.

## 61 2 Problem Statement

### 62 2.1 Split Learning

63 Given a training dataset  $\mathbf{X}$  and model parameters  $\theta$ , a learning task is to find approximately optimal  
64  $\theta$  by minimizing a pre-defined loss function  $\mathcal{L}$ . We assume that the optimizer used is a mini-batch  
65 stochastic gradient descent (SGD) algorithm, which updates  $\theta$  with a batch input of  $\mathbf{X}$  iteratively.  
66 Assuming the batch size is  $N$ , then the total loss of  $\theta$  for a batch input  $\mathbf{x} = \{x_i | x_i \in \mathbf{X}, i \in$   
67  $[1, N]\}$  should be  $\sum_{x \in \mathbf{x}} \mathcal{L}(\theta, x)$  in the  $t$ -th training iteration. The gradients of  $\theta$  for model updating  
68 should be estimated by  $\frac{1}{N} \sum_{x \in \mathbf{x}} \nabla_{\theta} \mathcal{L}(\theta, x)$  approximately. Hence, parameters  $\theta$  can be updated  
69 as  $\theta^{t+1} = \theta^t - \frac{1}{N} \sum_{x \in \mathbf{x}} \nabla_{\theta} \mathcal{L}(\theta, x)$ . This mini-batch SGD based optimizing procedure should be  
70 repeated until the model usability meets the requirement or the maximal count of iterations reaches.

71 Generally, there are two roles in a two-party SplitNN. We denote by the guest who holds features  
72 only and the host who holds both features and labels. According to related studies [34, 14], there  
73 exist several split configurations of neural networks. We will focus on the SplitNN designed for  
74 vertically partitioned features with labels held by the host only [4, 11]. We denote by  $\theta^g$ ,  $\theta^h$  the  
75 partial models of the guest and the host after split and  $\theta^k$  the rest part of the original model. Then  
76 the forwarding result  $\text{Forward}(\mathbf{x}^g, \theta^g)$  of the guest should be evaluated locally and passed to the  
77 host. The host should merge  $\text{Forward}(\mathbf{x}^h, \theta^h)$  with  $\text{Forward}(\mathbf{x}^g, \theta^g)$  using a predefined strategy  
78 such as concatenating and averaging. Then the host finishes the rest of the forward propagation  
79  $\mathcal{L}(\text{Merge}(\text{Forward}(\mathbf{x}^g, \theta^g), \text{Forward}(\mathbf{x}^h, \theta^h)), \theta^k)$  and initiates backward propagation, sending  
80 the partial loss regarding  $\mathbf{x}^g$  and  $\theta^g$  back to the guest.

81 We take a SplitNN based recommendation system as an example and give a benchmark in Table 1.  
82 Different merging strategies and two public datasets, MovieLens [15] and BookCrossing [40], have  
83 been evaluated. In case of misaligned features, zero padding will be used to retain the shape of fea-

84 ture vectors. We notice that concatenating and element-wise averaging with padding have relatively  
 85 stable and desirable performance, which will be used as the default setting in the rest of this work.

Table 1: Top-10 hit ratio of SplitNN based recommendation using different merging strategies.

		concatenate	element-wise				no split
			max	sum	average	min	
MovieLens	padding	56.62%	56.26%	56.35%	56.89%	57.19%	57.21%
	non-padding	55.38%	54.75%	54.95%	55.72%	55.08%	
Book Crossing	padding	61.70%	60.84%	60.21%	61.16%	60.98%	61.92%
	non-padding	58.80%	59.34%	59.44%	59.10%	59.02%	

85

## 86 2.2 Threat Model

87 In SplitNN, interactions between the guest and the host pose threats to each other. Thus, we will  
 88 investigate private data leakage threats from either direction. We first assume that the host and  
 89 the guest are honest but curious about the private data of each other. Moreover, we will consider  
 90 a more powerful threat model [29], where the guest or the host could be malicious, hijacking the  
 91 feature space during split learning. In both cases, we take into account property inference and data  
 92 reconstruction attacks, which are highly threatening attacks identified in collaborative learning.

93 **Property inference attack.** Since either role of SplitNN has access to the output of the other’s local  
 94 model, the adversary can mount a property inference attack [9, 22], inferring properties of private  
 95 data through observing query input and the corresponding output. By constructing shadow models  
 96 elaborately, the adversary can steal substantial information from the target. In this way, the adversary  
 97 acquires the capability of inferring some properties (such as gender and age) of the data samples used  
 98 for training. Denoted by  $F$ ,  $T$  and  $\mathcal{L}_F$  the inference model, target model, and the loss function used  
 99 for  $F$ , the adversarial goal of property inference attack is

$$\mathcal{A}_{PIA} = \arg \min_F \sum_{x_i \in \mathbf{x}} \mathcal{L}_F(F(T(x_i)), l_i), l_i \in \{0, 1\}. \quad (1)$$

100 **Data reconstruction attack.** By taking advantage of generative adversarial networks (GANs) [13], a  
 101 data reconstruction attack [17, 29] becomes possible in collaborative learning. To mount the attack,  
 102 the adversary augments the training data per iteration by inserting fake samples  $\mathbf{z}$  generated by a  
 103 generator  $G$ . The target model will serve as a discriminator  $D$ . The adversary affects the target  
 104 model by deceiving the target with fake training samples. For correcting the adversary, the target  
 105 participant is supposed to put more private information into the learning. In this game-style training,  
 106 the adversary can obtain a generator to reconstruct data samples similar to target private data. The  
 107 adversarial goal of the data reconstruction attack can be given as

$$\mathcal{A}_{DRA} = \min_G \max_D \frac{1}{|\mathbf{x}|} \sum_{x \in \mathbf{x}} \log D(x) + \frac{1}{|\mathbf{z}|} \sum_{z \in \mathbf{z}} \log(1 - D(G(z))). \quad (2)$$

108 **Feature space hijacking.** The property inference and data reconstruction attacks adapted from FL  
 109 can be mounted by a semi-honest participant, who follows the split learning protocol normally. How-  
 110 ever, a recent attack study dedicated to SplitNN has revealed that a malicious host can achieve more  
 111 impressive adversarial effects on property inference or data reconstruction by explicitly distorting  
 112 the objective of split learning. Thus, we also consider the defensive effect of our solution against  
 113 property inference and data reconstruction attacks using this feature space hijacking approach.

## 114 3 Privacy-Preserving Split Learning

115 Our goal is to design an unified defense solution to preserving privacy from both the host and guest’s  
 116 perspectives. Ideally, the guest wants to collaborate with the host under the condition that the host  
 117 could disclose no private information and vice versa. According to recent studies [3, 23, 16], artifi-  
 118 cial perturbations of data samples or parameters could prevent privacy leakage effectively. However,

119 different from conventional model publishing scenarios, the host and guest in SplitNN will keep  
 120 exchanging intermediate results during training. These continuous queries significantly increase the  
 121 risk of privacy leakage for both sides. Moreover, the attack surface of splitNN is in the middle of  
 122 neural network propagations, which makes things tricky. Thus, our primary idea to tackle this prob-  
 123 lem is to construct a bidirectional privacy-preserving tunnel for interactions. Recent studies such as  
 124 [10] have proved that activation functions are more adaptive for perturbed operands than other neu-  
 125 ral network components. Moreover, activation functions have various forms [2], which are flexible  
 126 for configuration. As a result, we propose a new variant of ReLU as a privacy-preserving interface.

### 127 3.1 R<sup>3</sup>eLU: Randomized-Response ReLU

128 Inspired by a randomized response approach [35], we propose a new activation function named  
 129 R<sup>3</sup>eLU. The original randomized response method is good at statistical analysis of item sets. But  
 130 the result of an activation function is commonly a continuous variable. Hence, they cannot be easily  
 131 combined together. It should also be noted that it is risky to perturb activation functions directly  
 132 because non-activated results may be flipped unexpectedly. Recall that the original formula of ReLU  
 133 is  $f(x) = \max(0, x)$ ,  $x \in \mathcal{R}$ . Our randomized-response variant will yield a proper substitute for  
 134 replacing the real activation with a probability of  $p$ . If we yield 0 as the substitute, then we can  
 135 randomly inactivate a part of ReLU results. But nothing has been changed for  $x \leq 0$ . Because  
 136  $f(x) = 0$  when  $x \leq 0$ . Hence, the variant is not completely privacy preserved since  $f(x) = 0$  also  
 137 reveals useful information to the adversary. For the completeness of the variant, we generate noisy  
 138 activations  $x' \leftarrow \text{Laplace}(0, \sigma)$  for  $x \leq 0$ . Now, we can give the definition of R<sup>3</sup>eLU as

$$R^3eLU(x) = \begin{cases} \max(0, x + x'), & \text{with probability } p, \\ 0, & \text{with probability } (1 - p). \end{cases} \quad (3)$$

139 We remark that the way R<sup>3</sup>eLU handles non-activated results is dangerous, although learning ac-  
 140 curacy is often traded off for privacy. But we will mitigate the side effect through privacy budget  
 141 allocation. Then the risk of applying R<sup>3</sup>eLU will not be a problem.

### 142 3.2 Forward Propagation with R<sup>3</sup>eLU

143 Now we show how to apply R<sup>3</sup>eLU in the forward pass of SplitNN regarding the guest’s privacy.  
 144 Generally, when the  $t$ -th training iteration begins, the guest randomly samples a mini-batch  $\mathbf{x}$  from  
 145 private training dataset  $\mathcal{X}^{(g)}$ . Assuming that an embedding procedure  $\text{Embedding}() : \mathcal{R}^M \leftarrow$   
 146  $\mathcal{X}$  is publicly available, raw data samples in the mini-batch can be encoded into feature vectors  
 147  $V = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\}$ ,  $\mathbf{v}_i = \{v_{i1}, v_{i2}, \dots, v_{iM}\}$ ,  $i \in [1, N]$ , and  $M$  is the dimension of feature  
 148 representing space. From the functional perspective, we assume that  $\mathcal{L}(\mathbf{v}^h, \mathbf{v}^g, \boldsymbol{\theta})$  is equivalent to  
 149  $\mathcal{L}(\text{Forward}(\mathbf{v}^h, \boldsymbol{\theta}^h), \text{Forward}(\mathbf{v}^g, \boldsymbol{\theta}^g), \boldsymbol{\theta}^k)$ , where  $\text{Forward}() : \mathcal{R}^{N_s} \leftarrow \mathcal{R}^M$  is a function to yield  
 150 the feedforward result of hidden units for a given neural network, and  $N_s$  indicates the output shape.

151 We now make the minimal modification of the forward propagation of SplitNN while leaving the  
 152 rest part unchanged. Denoted by  $\text{Forward-R}^3\text{eLU}()$  the feedforward result of neural network by  
 153 replacing the activation functions of the guest’s output layer with R<sup>3</sup>eLU. Then the feedforward  
 154 result transmitted to the host should be  $\mathbf{a}^g = \text{Forward-R}^3\text{eLU}(\mathbf{v}^g, \boldsymbol{\theta}^g)$ . Next, the host executes a pre-  
 155 defined aggregation procedure taking as input  $\mathbf{a}^g$  and  $\mathbf{a}^h$ . Finally, the loss function will be evaluated  
 156 by  $\mathcal{L}(\text{Forward-R}^3\text{eLU}(\mathbf{v}^g, \boldsymbol{\theta}^g), \text{Forward}(\mathbf{v}^h, \boldsymbol{\theta}^h), \boldsymbol{\theta}^k)$ . Algorithm 1 in the appendix integrates the  
 157 above forward propagation using R<sup>3</sup>eLU into SplitNN.

### 158 3.3 Privacy-Preserving Backward Propagation

159 According to recent studies of privacy leakage in backward propagation [26, 31, 8], model updating  
 160 information may cause severe leakage of private training data. Since the partial loss will be propa-  
 161 gated to the guest in SplitNN, it is crucial to protect the host’s privacy from being disclosed. Fortu-  
 162 nately, our R<sup>3</sup>eLU is adaptable to noisy partial losses and we design R<sup>3</sup>eLU-Diff in a randomized-

163 response manner for the derivative of R<sup>3</sup>eLU as

$$\text{R}^3\text{eLU-Diff}(\delta^g, \mathbf{a}^g, \mathbf{v}^g) = \begin{cases} \delta^g \times \text{ReLU-Diff}(\mathbf{a}^g, \mathbf{v}^g) + x', & \text{with probability } (1 - p), \\ 0, & \text{with probability } p. \end{cases}$$

164 Noting that the derivative value of ReLU for any input is either one or zero, randomly flipping the  
 165 derivative value may still disclose  $\delta_t^g$  when value ones are not flipped. Therefore, we integrate a  
 166 Laplace mechanism into the R<sup>3</sup>eLU-Diff. Please also note that the constructions of R<sup>3</sup>eLU-Diff and  
 167 R<sup>3</sup>eLU are in a similar way. This is helpful to obtain uniform analysis results of two parties, which  
 168 will be shown in the privacy analysis part.

169 We sketch the backward propagation using R<sup>3</sup>eLU-Diff in Algorithm 2 in the appendix. Briefly,  
 170 when the  $t$ -th backward propagation begins, the host calculates the partial loss  $\delta^k = \nabla_{\mathbf{a}^k} \mathcal{L}(l, \bar{\mathbf{a}}, \theta_t^k)$   
 171 for  $\theta^k$  regarding the total loss  $l$  obtained in forward propagation, where  $\bar{\mathbf{a}}$  is the averaged activation  
 172 result. Then a `Backward()`:  $\mathcal{R}^{N_r} \leftarrow \mathcal{R}^{N_s}$  procedure calculates the gradients of model parameters,  
 173 where  $N_r$  is the shape of parameters. To update the guest model, the host sends partial loss  $\delta^g$  for  
 174  $\theta^g$  to the guest. To preserve the host’s privacy within  $\delta^g$ , we disturb the partial loss  $\delta^g$  propagating  
 175 to the guest and keep the partial loss on the host side unchanged.

### 176 3.4 Dynamic Privacy Budget Allocation

177 It has been proved that the importance of a parameter can be quantified by the error introduced when  
 178 it is removed from the model [28]. Thus, we define the importance  $I_j$  of a given SplitNN parameter  
 179  $\theta_j \in \theta$  as the squared difference of prediction errors caused by removing  $\theta_j$ ,

$$I_j = (\mathcal{L}(\mathbf{x}, \theta) - \mathcal{L}(\mathbf{x}, \theta \setminus \{\theta_j\}))^2. \quad (4)$$

180 Due to the consideration of efficiency, it is suggested in [28] to estimating the importance by a  
 181 first-order Taylor expansion approximately. Then the importance of  $\theta_j$  is estimated as

$$\hat{I}_j = (g_j \cdot \theta_j)^2, \quad (5)$$

182 where  $g_j$  is the gradient of the parameter  $\theta_j$  regarding a specific sample when  $\theta$  is well-trained.  
 183 Given parameter importance, the importance of features can be derived further. Specifically, the  
 184 importance of a neuron (or a feature)  $U_j$ ,  $j \in [1, N_u]$ , where  $N_u$  is the total number of neurons  
 185 in the model, can be calculated as a joint importance of relevant parameters by summing up the  
 186 importance of all relevant parameters. Hence,  $U_j = \sum_{\theta_k \in \hat{\theta}_j} \hat{I}_k$ , where  $\hat{\theta}_j$  denotes the set of all  
 187 parameters directly connected to the  $j$ -th neuron.

188 However, the above importance estimation method is designed for a well-trained model and cannot  
 189 be directly applied to intermediate models during training. To tackle this problem, we give a dynamic  
 190 estimation method by deriving the original method into a cumulative form. The importance of a  
 191 feature will be accumulated as the training epoch increases. Specifically, the importance of the  $j$ -th  
 192 feature in the  $q$ -th training epoch is

$$U_j^q = \frac{\sum_{\theta_k \in \hat{\theta}_j} \hat{I}_k + U_j^{q-1} \times (q \times \lfloor T/n_t \rfloor + (t \bmod n_t) - 1)}{q \times \lfloor T/n_t \rfloor + (t \bmod n_t)}, \quad (6)$$

193 where  $n_t$  indicates the iteration number within a training epoch,  $T$  is the maximum training iteration  
 194 number, and  $t$  is the current training iteration globally. Generally, we assume that  $T/n_t = n_q$ ,  
 195  $n_q \in \mathcal{N}$ , which also means that  $q \in [1, n_q]$ . We give the importance estimation results of different  
 196 neurons in Figure 3 in the appendix, showing the correctness of our dynamic importance estimation  
 197 method and the existence of unbalanced feature importance.

198 Given the parameter importance estimated dynamically, we are capable of allocating privacy budgets  
 199 regarding different features. The intuition is to give larger budgets to more important features while  
 200 smaller budgets to less important ones. Before the  $q$ -th training epoch begins, we estimate the  
 201 feature importance vector  $\mathbf{U} = \{U_1^q, U_2^q, \dots, U_{N_u}^q\}$ . Based on  $\mathbf{U}$ , the corresponding privacy budget  
 202 to be allocated should be  $\epsilon_j \times U_j^q$ ,  $j \in [1, N_u]$ . And the total privacy budget for all features is

203  $\epsilon_F = \sum_{j \in [1, N_u]} \epsilon_j$ . On the other hand, we can also dynamically allocate privacy budgets for  
 204 different iterations to optimize the total privacy budget further. Given the total privacy budget  $\epsilon_T$   
 205 for all iterations, we assign the privacy budget  $\epsilon_i = \frac{\epsilon_T}{2^i}$  to the  $i$ -th iteration. Since  $\sum_{i=1}^{\infty} \frac{\epsilon_T}{2^i} = \epsilon_T$ ,  
 206 according to the sequential composition theory of differential privacy, we can still ensure that the  
 207 whole training process achieves  $\epsilon_T$ -differential privacy.

### 208 3.5 Privacy Analysis

209 We will give the privacy analysis for the host and the guest respectively. In forward propagation,  
 210 we recall that  $v^g$ ,  $\theta^g$ , and  $\mathbf{a}^g$  are input features, model parameters, and activation results of the  
 211 guest. According to the definition of  $R^3eLU$ ,  $\mathbf{a}^g$  should be randomly flipped with probability  $p$ . For  
 212 brevity, we denote by  $f(\cdot)$  the evaluation of neural network before activation. Then the probability  
 213 of observing any activation  $a_o$  for a given input  $v$  should be

$$P(a_o = 0|v) = p + (1-p) \int_{-\infty}^{-f(v)_o} \frac{1}{2b} \exp\left(-\frac{|v|}{b}\right) dv = p + \frac{1-p}{2} \exp\left(\frac{-f(v)_o}{b}\right),$$

$$P(a_o > 0|v) = \frac{1-p}{2b} \exp\left(-\frac{|a_o - f(v)_o|}{b}\right).$$

214 Then we can give the following conclusion regarding the guest model. The proof is in the appendix.

215 **Corollary 1.** *When forward propagation of the guest model is activated by  $R^3eLU$  in split learning,*  
 216 *the activation result is  $\epsilon$ -DP, given Laplace noise scale  $\sigma_g$ .*

217 For a fine-grained privacy budget allocation, we introduce dynamic budget allocation in our solution.  
 218 We now give the analysis of feature-specific privacy budget. Given the estimated feature importance  
 219 vector  $U = \{U_1^q, U_2^q, \dots, U_{N_u}^q\}$ , we allocate privacy budget  $\epsilon_i = \epsilon U_i^q$  to each feature. Thus, if the  
 220 divergence of features can be bounded by their privacy budgets, then the total privacy budget will  
 221 be bounded. Denoted by  $b_i$  and  $c_i$  the noise parameter and bound of the  $i$ -th feature. Then we can  
 222 account the divergence of features as

$$\frac{P(a_i > 0|v)}{P(a_i > 0|v')} = \frac{\frac{1-p}{2b} \exp\left(-\frac{|a_i - f(v)_i|}{b_i}\right)}{\frac{1-p}{2b} \exp\left(-\frac{|a_i - f(v')_i|}{b_i}\right)} \leq \exp\left(\frac{|a_i - f(v')_i| - |a_i - f(v)_i|}{b_i}\right) \leq \exp\left(\frac{2c_i}{b_i}\right)$$

$$\frac{P(a_i = 0|v)}{P(a_i = 0|v')} = \frac{p + \frac{1-p}{2} \exp\left(\frac{-f(v)_i}{b_i}\right)}{p + \frac{1-p}{2} \exp\left(\frac{-f(v')_i}{b_i}\right)} \leq \frac{p + \frac{1-p}{2} \exp\left(\frac{c_i}{b_i}\right)}{p + \frac{1-p}{2} \exp\left(\frac{-c_i}{b_i}\right)} \leq \exp\left(\frac{2c_i}{b_i}\right)$$

223 Basing on this result, we can conclude the following corollary.

224 **Corollary 2.** *In forward propagation with  $R^3eLU$ , the privacy budget of the  $i$ -th feature can be*  
 225 *bounded by  $\epsilon_i$  if we choose  $b_i$  to satisfy  $\exp\left(\frac{2c_i}{b_i}\right) \leq \exp(\epsilon_i)$ ,  $\forall p \in [0, 1]$ .*

226 The last analysis result of the guest is privacy budget allocation during the whole training stage.  
 227 Since we have allocated  $\epsilon = \sum_{i \in [1, N_u]} \epsilon_i$  for all features, we can directly conclude that each training  
 228 step is  $\gamma\epsilon$ -DP by following the privacy amplification theory, where  $\gamma = \frac{B}{N}$  is the sampling ratio of a  
 229 batch regarding the whole training dataset. Now we can give the total privacy budget of the whole  
 230 training stage using the strong composition theorem.

231 **Corollary 3.** *The total privacy budget of the whole training process using  $R^3eLU$  is  $(\epsilon'_g, \delta'_g)$ -DP,*  
 232 *where  $\epsilon'_g = \gamma\epsilon \sqrt{2T \ln\left(\frac{1}{\delta'_g}\right)} + \gamma\epsilon T(e^{\gamma\epsilon} - 1)$ .*

233 So far, we have given the privacy analysis from the guest perspective. Since we construct  $R^3eLU$  and  
 234  $R^3eLU$ -Diff using the same way and same technique, these two procedures have the same analysis  
 235 result if we choose the noise scale of host  $\sigma_h = \sigma_g$ . In this way, we can conclude the following  
 236 result for the host.

237 **Corollary 4.** *In backward propagation with  $R^3eLU$ -Diff, the privacy budget of the host can be*  
 238 *bounded by  $\epsilon_h$  if we choose  $\sigma_h = \sigma_g$ . The total privacy budget of the host in the whole training*  
 239 *process is  $(\epsilon'_h, \delta'_h)$ -DP, where  $\epsilon'_h = \gamma\epsilon_h \sqrt{2T \ln\left(\frac{1}{\delta'_h}\right)} + \gamma\epsilon_h T(e^{\gamma\epsilon_h} - 1)$ .*

240 **4 Evaluation**

241 We evaluate our privacy-preserving SplitNN solution from two aspects model usability and privacy  
 242 loss. To be comprehensive, we will compare our solution with the baseline (without any protection)  
 243 and the most relevant defense solutions, i.e., a primitive Laplace mechanism [5] and DPSGD [1],  
 244 the most well-known privacy-preserving deep learning solution, in the same setting. We will use  
 245 the same fixed total privacy budget and the same split way (shown in the appendix) for all solutions.  
 246 We adapt solutions into recommendation models using two real-world datasets, MovieLens [15],  
 247 BookCrossing [40], and an image classification model using MNIST [20] dataset. The MovieLens  
 248 1-M dataset contains 1 million ratings of 4,000 movies collected from 6,000 users and users’ de-  
 249 mographic information such as gender and age. The BookCrossing dataset includes 278,858 users’  
 250 demographic information and 1,149,780 ratings of 271,379 books. The MNIST database has 70,000  
 251 handwriting image examples. We will use a 32 batch size, a 0.01 learning rate and an Adam opti-  
 252 mizer as default. Since different datasets and defense solutions may require various epochs for split  
 253 learning, we calculate the metrics when the learning converges, or the privacy budget is drained. All  
 experimental results are averaged across multiple runs.

Table 2: Model usability results while preserving the privacy of the guest.

$\epsilon$	MovieLens			BookCrossing			MNIST		
	Laplace	DPSGD	Ours	Laplace	DPSGD	Ours	Laplace	DPSGD	Ours
0.1	30.84%	32.29%	<b>34.03%</b>	57.02%	55.89%	<b>58.18%</b>	17.43%	30.21%	<b>32.41%</b>
0.5	41.25%	43.69%	<b>43.87%</b>	57.67%	56.14%	<b>58.54%</b>	27.33%	58.43%	<b>60.38%</b>
1.0	48.16%	49.09%	<b>50.56%</b>	58.02%	56.56%	<b>58.42%</b>	31.05%	75.58%	<b>76.60%</b>
2.0	49.32%	50.38%	<b>50.49%</b>	58.74%	56.91%	<b>59.24%</b>	38.92%	92.90%	<b>93.53%</b>
4.0	49.26%	<b>50.86%</b>	50.73%	59.01%	57.16%	<b>59.26%</b>	95.37%	<b>95.87%</b>	94.12%

Table 3: Model usability results while preserving the privacy of the host.

$\epsilon$	MovieLens			BookCrossing			MNIST		
	Laplace	DPSGD	Ours	Laplace	DPSGD	Ours	Laplace	DPSGD	Ours
0.1	31.47%	30.68%	<b>33.98%</b>	57.37%	57.46%	<b>58.26%</b>	27.64%	<b>33.45%</b>	32.36%
0.5	41.75%	42.31%	<b>42.67%</b>	58.62%	58.24%	<b>58.59%</b>	55.38%	65.28%	<b>67.83%</b>
1.0	47.43%	48.29%	<b>50.39%</b>	59.49%	58.44%	<b>59.77%</b>	71.95%	<b>89.74%</b>	88.14%
2.0	49.86%	50.43%	<b>51.47%</b>	59.34%	59.97%	<b>60.27%</b>	89.15%	<b>92.66%</b>	92.52%
4.0	49.57%	50.09%	<b>51.62%</b>	59.55%	<b>60.75%</b>	60.66%	94.61%	<b>95.37%</b>	95.01%

254

255 **4.1 Model Usability**

256 Since artificial perturbation may affect the learning procedure, we will evaluate how SplitNN is  
 257 affected by privacy-preserving solutions. Moreover, two asymmetric parties of SplitNN may have  
 258 different privacy concerns and affect learning differently. Thus, we will evaluate model usability  
 259 concerning privacy from the perspective of the guest or the host. We use an averaged test accuracy  
 260 across all test samples for the evaluation. Precisely, the test accuracy of a recommendation model  
 261 is calculated using a top-10 hit ratio, while the test accuracy of an image classifier is the prediction  
 262 accuracy. In Table 2 and Table 3, we show the model usability results regarding various privacy  
 263 budget values of two parties. We note that model accuracy baselines of MovieLens, BookCrossing,  
 264 and MNIST are 56.62%, 61.70%, and 98.00%, respectively.

265 For the MovieLens model, our solution achieves the best model usability in most cases, especially  
 266 with less privacy budget. DPSGD has a better result when  $\epsilon = 4$  for the guest. But DPSGD will  
 267 cause a significant privacy leakage in this case. For the BookCrossing model, the model usability of  
 268 our solution is relatively high in cases of protecting the guest and the host. Similarly, DPSGD has a  
 269 better result when  $\epsilon = 4$  for the host, sacrificing the privacy guarantee. Things are a bit different for  
 270 the MNIST model. DPSGD gets some better results when protecting the host. The reason is that split  
 271 learning for an image classification model segments image samples roughly, making our dynamic  
 272 budget allocation approach malfunction. Meanwhile, DPSGD is not designed for protecting partial  
 273 loss in SplitNN, leading to an optimistic estimation of the threat against the host. Apart from these  
 274 exceptional cases, our solution outperforms other solutions on model usability.

275 **4.2 Defense Against Inference and Reconstruction Attacks**

276 We will evaluate the performance of privacy preservation by comparing attack results against  
 277 SplitNN with and without the defense. We will mount property inference and data reconstruction  
 278 attacks against the guest and the host, respectively. Prediction accuracy of the adversary’s inference  
 279 model will be used to measure the performance of the property inference attack. As for the data re-  
 280 construction attack, the adversary tries to generate data samples as similar as possible to the target’s  
 281 private data. In this case, we can use a mean squared error (MSE) between a generated sample and  
 a target data sample to measure the adversarial effect.

Table 4: Results of defending the guest against property inference attack.

$\epsilon$	MovieLens			BookCrossing			MNIST		
	Laplace	DPSGD	Ours	Laplace	DPSGD	Ours	Laplace	DPSGD	Ours
0.1	66.99%	77.71%	<b>60.99%</b>	<b>54.76%</b>	73.29%	55.78%	<b>43.27%</b>	53.95%	44.33%
0.5	66.16%	74.23%	<b>64.16%</b>	<b>54.97%</b>	74.52%	56.33%	46.92%	54.23%	<b>45.59%</b>
1.0	<b>67.19%</b>	78.65%	68.18%	<b>55.03%</b>	74.96%	58.65%	47.58%	54.26%	<b>47.51%</b>
2.0	68.65%	73.06%	<b>68.56%</b>	<b>54.85%</b>	74.26%	58.14%	<b>48.06%</b>	54.65%	52.87%
4.0	<b>69.14%</b>	76.18%	71.91%	<b>54.92%</b>	74.33%	60.76%	<b>48.47%</b>	54.57%	55.73%

Table 5: Results of defending the host against property inference attack.

$\epsilon$	MovieLens			BookCrossing			MNIST		
	Laplace	DPSGD	Ours	Laplace	DPSGD	Ours	Laplace	DPSGD	Ours
0.1	53.46%	78.59%	<b>51.86%</b>	<b>54.55%</b>	74.35%	59.42%	60.34%	80.29%	<b>48.74%</b>
0.5	53.46%	75.64%	<b>51.89%</b>	<b>54.62%</b>	74.36%	59.42%	59.82%	81.92%	<b>49.71%</b>
1.0	53.46%	73.54%	<b>52.75%</b>	<b>54.95%</b>	74.39%	59.52%	59.74%	82.80%	<b>50.48%</b>
2.0	<b>53.47%</b>	75.05%	59.77%	<b>54.40%</b>	74.39%	58.13%	60.38%	88.88%	<b>50.57%</b>
4.0	<b>53.48%</b>	79.28%	56.52%	<b>54.95%</b>	74.39%	62.04%	60.62%	89.73%	<b>50.47%</b>

282

283 **Defense against property inference attack.** We give evaluation results of the defensive effect of the  
 284 guest and the host in Table 4 and Table 5, respectively, inferring the property of users in MovieLens  
 285 and BookCrossing and an image patch in MNIST. The attack accuracy against baselines of Movie-  
 286 Lens, BookCrossing, and MNIST models can achieve above 80%, 79%, and 94% by an adversarial  
 287 host, 80%, 78%, and 57% by an adversarial guest, respectively. However, our SplitNN solution  
 288 can effectively mitigate the adversarial effect during training and decrease the attack accuracy sig-  
 289 nificantly. It should be noted that the primitive Laplace mechanism frustrates the inference attack  
 290 badly because the artificial noise added by the mechanism is indiscriminate, leading to conspicuous  
 291 damages to the model usability. Even so, our solution has significant advantages on MovieLens and  
 292 MNIST datasets. In contrast, the primitive Laplace mechanism cannot cover image classification  
 cases, while DPSGD cannot defeat the attack.

Table 6: Results of defending the guest against data reconstruction attack.

$\epsilon$	MovieLens			BookCrossing			MNIST		
	Laplace	DPSGD	Ours	Laplace	DPSGD	Ours	Laplace	DPSGD	Ours
0.1	0.2459	0.2455	<b>0.3223</b>	0.3216	0.2907	<b>0.3329</b>	1.8849	1.8885	<b>2.0181</b>
0.5	0.2453	0.2451	<b>0.3222</b>	0.3202	0.2902	<b>0.3329</b>	1.8024	1.8137	<b>1.9875</b>
1.0	0.2453	0.2451	<b>0.3222</b>	0.3202	0.2902	<b>0.3221</b>	1.7857	1.7509	<b>1.9533</b>
2.0	0.2452	0.2451	<b>0.3222</b>	0.3202	0.2902	<b>0.3221</b>	1.7336	1.7469	<b>1.9391</b>
4.0	0.2452	0.2451	<b>0.3222</b>	0.3202	0.2902	<b>0.3221</b>	1.7014	1.7440	<b>1.9206</b>

293

294 **Defense against data reconstruction attack.** We show the defense results against an adversarial host  
 295 and an adversarial guest in Table 6 and Table 7, respectively. We note that the MSE is measured  
 296 after the attack model has been trained sufficiently in all cases. The MSEs measured for the attack  
 297 against baselines of MovieLens, BookCrossing, and MNIST models are 0.2412, 0.2629, and 0.9612  
 298 by an adversarial host, 0.2369, 0.2402, and 1.6998 by an adversarial guest, respectively. Please note  
 299 that these attack results against the baselines are frustrating because the reconstruction attack is hard  
 300 to succeed in the semi-honest setting. Meanwhile, data samples in two recommendation datasets are  
 301 similar and embedded with the same feature vectors. This leads to similar reconstruction results and  
 302 similar MSEs. But we can still conclude from the results that our solution has dominant performance  
 303 in the defense against reconstruction attacks mounted by either side.



Table 7: Results of defending the host against data reconstruction attack.

$\epsilon$	MovieLens			BookCrossing			MNIST		
	Laplace	DPSGD	Ours	Laplace	DPSGD	Ours	Laplace	DPSGD	Ours
0.1	0.4032	0.2417	<b>0.5486</b>	0.4237	0.2758	<b>0.5066</b>	1.2887	1.0875	<b>1.8257</b>
0.5	0.4024	0.2419	<b>0.5357</b>	0.4222	0.2756	<b>0.5149</b>	1.2778	1.0685	<b>1.7758</b>
1.0	0.4008	0.2422	<b>0.5285</b>	0.4217	0.2743	<b>0.5235</b>	1.2602	1.0422	<b>1.7528</b>
2.0	0.3982	0.2421	<b>0.5083</b>	0.4214	0.2697	<b>0.5150</b>	1.2613	1.0333	<b>1.7334</b>
4.0	0.3960	0.2422	<b>0.4819</b>	0.4194	0.2683	<b>0.5046</b>	1.2549	0.9996	<b>1.7262</b>

304 *Defense against feature space hijacking attacks (FSHA)*. Please note that property inference and  
 305 data reconstruction attacks implemented in [29] hijack the learning objective, offering the adversary  
 306 an advantage over the previous attacks we straightforwardly adapted from FL. In this setting, the  
 307 malicious attacker trains a generator using the split neural network as a discriminator during the  
 308 split learning process and uses a gradient-scaling trick to train the generator. Since the generating  
 309 part is critical to FSHA, we will focus on the defense against the reconstruction. If the generating  
 310 part fails, the inference attack will be impossible. Here we give the defense results of the MNIST  
 311 model because FSHA [29] is mainly evaluated using this dataset. We also evaluate our solution for  
 312 other datasets against FSHA, and the results are given in the appendix. In Figure 1 and Figure 2,  
 313 we give the reconstruction results of FSHA mounted by an adversarial host and guest against target  
 314 samples used in [29]. The second row of two figures shows the results of FSHA against baselines.  
 315 The following rows show that our solution can effectively preserve private data for both the guest and  
 316 the host, even the privacy budget is relaxed to 4. More practical privacy budget values for defending  
 317 against FSHA are presented in the appendix.



Figure 1: Reconstruction results of FSHA against the guest’s data in the first row. The following rows are attack results against the original SplitNN and our solution ( $\epsilon = 0.1, 1.0, 4.0$ ), respectively.

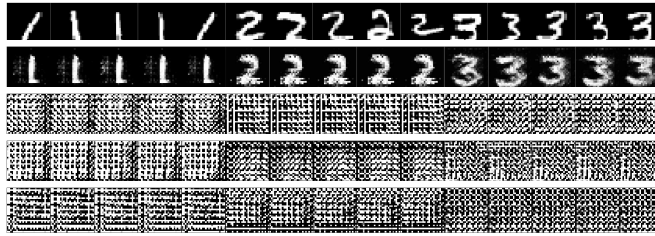


Figure 2: Reconstruction results of FSHA against the host’s data in the first row. The following rows are attack results against the original SplitNN and our solution ( $\epsilon = 0.1, 1.0, 4.0$ ), respectively.

## 318 5 Conclusion and Limitation

319 Our privacy-preserving SplitNN solution, built upon a new activation function  $R^3eLU$  and its deriva-  
 320 tive  $R^3eLU-Diff$  in a randomized-response manner, significantly reduces privacy leakage risk for  
 321 both the guest and the host. We show that our solution can provide a tight privacy budget for split  
 322 learning through the privacy analysis. The model usability and privacy loss can be further balanced  
 323 by our dynamic privacy budget allocation. The experimental evaluation using different datasets  
 324 shows that our solution outperforms the existing privacy-preserving SplitNN solutions in model us-  
 325 ability and privacy protection. We note that our solution deals with property inference and data  
 326 reconstruction attacks in a feature level, but a clustering-based label inference attack [8] is out of  
 327 our reach, which is an interesting topic to be studied in future work.

## References

- [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *ACM SIGSAC conference on computer and communications security*, pages 308–318, 2016.
- [2] Andrea Apicella, Francesco Donnarumma, Francesco Isgrò, and Roberto Prevete. A survey on modern trainable activation functions. *Neural Networks*, 2021.
- [3] Eugene Bagdasaryan, Omid Poursaeed, and Vitaly Shmatikov. Differential privacy has disparate impact on model accuracy. *Advances in Neural Information Processing Systems*, 32: 15479–15488, 2019.
- [4] Iker Ceballos, Vivek Sharma, Eduardo Mugica, Abhishek Singh, Alberto Roman, Praneeth Vepakomma, and Ramesh Raskar. Splitnn-driven vertical partitioning. *arXiv preprint arXiv:2008.04137*, 2020.
- [5] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, pages 211–407, 2014.
- [6] Ege Erdogan, Alptekin Kupcu, and A Ercument Cicek. Unsplit: Data-oblivious model inversion, model stealing, and label inference attacks against split learning. *arXiv preprint arXiv:2108.09033*, 2021.
- [7] Minghong Fang, Neil Zhenqiang Gong, and Jia Liu. Influence function based data poisoning attacks to top-n recommender systems. In *Proceedings of The Web Conference*, pages 3019–3025, 2020.
- [8] Chong Fu, Xuhong Zhang, Shouling Ji, Jinyin Chen, Jingzheng Wu, Shanqing Guo, Jun Zhou, Alex X Liu, and Ting Wang. Label inference attacks against vertical federated learning. In *USENIX Security Symposium*, 2022.
- [9] Karan Ganju, Qi Wang, Wei Yang, Carl A Gunter, and Nikita Borisov. Property inference attacks on fully connected neural networks using permutation invariant representations. In *ACM SIGSAC conference on computer and communications security*, pages 619–633, 2018.
- [10] Hongyang Gao, Lei Cai, and Shuiwang Ji. Adaptive convolutional relus. In *AAAI Conference on Artificial Intelligence*, volume 34, pages 3914–3921, 2020.
- [11] Yansong Gao, Minki Kim, Sharif Abuadbba, Yeonjae Kim, Chandra Thapa, Kyuyeon Kim, Seyit A Camtepe, Hyounghick Kim, and Surya Nepal. End-to-end evaluation of federated learning and split learning for internet of things. In *International Symposium on Reliable Distributed Systems*, pages 91–100, 2020.
- [12] Ran Gilad-Bachrach, Nathan Dowlin, Kim Laine, Kristin Lauter, Michael Naehrig, and John Wernsing. Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. In *International conference on machine learning*, pages 201–210, 2016.
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [14] Otkrist Gupta and Ramesh Raskar. Distributed learning of deep neural network over multiple agents. *Journal of Network and Computer Applications*, 116:1–8, 2018.
- [15] Attribution F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *ACM Transactions on Interactive Intelligent Systems*, 5(4), 2015.
- [16] Xinlei He and Yang Zhang. Quantifying and mitigating privacy risks of contrastive learning. In *ACM Conference on Computer and Communications Security*, 2021.

- 372 [17] Briland Hitaj, Giuseppe Ateniese, and Fernando Perez-Cruz. Deep models under the gan: in-  
373 formation leakage from collaborative deep learning. In *ACM SIGSAC Conference on Computer*  
374 *and Communications Security*, pages 603–618, 2017.
- 375 [18] Hai Huang, Jiaming Mu, Neil Zhenqiang Gong, Qi Li, Bin Liu, and Mingwei Xu. Data poison-  
376 ing attacks to deep learning based recommender systems. In *Annual Network and Distributed*  
377 *System Security Symposium*, 2021.
- 378 [19] Juyong Kim, Yookoon Park, Gunhee Kim, and Sung Ju Hwang. Splitnet: Learning to seman-  
379 tically split deep networks for parameter reduction and model parallelization. In *International*  
380 *Conference on Machine Learning*, pages 1866–1874, 2017.
- 381 [20] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning  
382 applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 383 [21] Yang Liu, Zhuo Ma, Yilong Yang, Ximeng Liu, Jianfeng Ma, and Kui Ren. Revfrf: Enabling  
384 cross-domain random forest training with revocable federated learning. *IEEE Transactions on*  
385 *Dependable and Secure Computing*, 2021.
- 386 [22] Xinjian Luo, Yuncheng Wu, Xiaokui Xiao, and Beng Chin Ooi. Feature inference attack on  
387 model predictions in vertical federated learning. In *IEEE International Conference on Data*  
388 *Engineering*, pages 181–192, 2021.
- 389 [23] Yunlong Mao, Wenbo Hong, Boyu Zhu, Zhifei Zhu, Yuan Zhang, and Sheng Zhong. Secure  
390 deep neural network models publishing against membership inference attacks via training task  
391 parallelism. *IEEE Transactions on Parallel and Distributed Systems*, 2021.
- 392 [24] Yunlong Mao, Xinyu Yuan, Xinyang Zhao, and Sheng Zhong. Romoa: Robust model ag-  
393 gregation for the resistance of federated learning to model poisoning attacks. In *European*  
394 *Symposium on Research in Computer Security*, pages 476–496, 2021.
- 395 [25] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Ar-  
396 cas. Communication-efficient learning of deep networks from decentralized data. In *Artificial*  
397 *intelligence and statistics*, pages 1273–1282, 2017.
- 398 [26] Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov. Exploiting  
399 unintended feature leakage in collaborative learning. In *IEEE Symposium on Security and*  
400 *Privacy*, pages 691–706. IEEE, 2019.
- 401 [27] Payman Mohassel and Yupeng Zhang. Secureml: A system for scalable privacy-preserving  
402 machine learning. In *IEEE symposium on security and privacy*, pages 19–38, 2017.
- 403 [28] Pavlo Molchanov, Arun Mallya, Stephen Tyree, Iuri Frosio, and Jan Kautz. Importance esti-  
404 mation for neural network pruning. In *IEEE/CVF Conference on Computer Vision and Pattern*  
405 *Recognition*, pages 11264–11272, 2019.
- 406 [29] Dario Pasquini, Giuseppe Ateniese, and Massimo Bernaschi. Unleashing the tiger: Inference  
407 attacks on split learning. *ACM SIGSAC Conference on Computer and Communications Secu-*  
408 *rity*, 2021.
- 409 [30] Ahmed Salem, Yang Zhang, Mathias Humbert, Mario Fritz, and Michael Backes. MI-leaks:  
410 Model and data independent membership inference attacks and defenses on machine learning  
411 models. In *Annual Network and Distributed Systems Security Symposium*, 2019.
- 412 [31] Ahmed Salem, Apratim Bhattacharya, Michael Backes, Mario Fritz, and Yang Zhang. Updates-  
413 leak: Data set inference and reconstruction attacks in online learning. In *USENIX Security*  
414 *Symposium*, pages 1291–1308, 2020.

- 415 [32] Lichao Sun, Jianwei Qian, and Xun Chen. LDP-FL: practical private aggregation in federated learning with local differential privacy. In *International Joint Conference on Artificial Intelligence*, pages 1571–1578, 2021.
- 416
- 417
- 418 [33] Vale Tolpegin, Stacey Truex, Mehmet Emre Gursoy, and Ling Liu. Data poisoning attacks against federated learning systems. In *European Symposium on Research in Computer Security*, pages 480–501, 2020.
- 419
- 420
- 421 [34] Praneeth Vepakomma, Otkrist Gupta, Tristan Swedish, and Ramesh Raskar. Split learning for health: Distributed deep learning without sharing raw patient data. *arXiv preprint arXiv:1812.00564*, 2018.
- 422
- 423
- 424 [35] Stanley L Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60(309):63–69, 1965.
- 425
- 426 [36] Kang Wei, Jun Li, Ming Ding, Chuan Ma, Howard H Yang, Farhad Farokhi, Shi Jin, Tony QS Quek, and H Vincent Poor. Federated learning with differential privacy: Algorithms and performance analysis. *IEEE Transactions on Information Forensics and Security*, 15:3454–3469, 2020.
- 427
- 428
- 429
- 430 [37] Ji Yang, Xinyang Yi, Derek Zhiyuan Cheng, Lichan Hong, Yang Li, Simon Xiaoming Wang, Taibai Xu, and Ed H Chi. Mixed negative sampling for learning two-tower neural networks in recommendations. In *Proceedings of the Web Conference*, pages 441–447, 2020.
- 431
- 432
- 433 [38] Chengliang Zhang, Suyi Li, Junzhe Xia, Wei Wang, Feng Yan, and Yang Liu. Batchcrypt: Efficient homomorphic encryption for cross-silo federated learning. In *USENIX Annual Technical Conference*, 2020.
- 434
- 435
- 436 [39] Yin Zhang, Derek Zhiyuan Cheng, Tiansheng Yao, Xinyang Yi, Lichan Hong, and Ed H Chi. A model of two tales: Dual transfer learning framework for improved long-tail item recommendation. In *Proceedings of the Web Conference*, pages 2220–2231, 2021.
- 437
- 438
- 439 [40] Cai-Nicolas Ziegler, Sean M McNeen, Joseph A Konstan, and Georg Lausen. Improving recommendation lists through topic diversification. In *International Conference on World Wide Web*, pages 22–32, 2005.
- 440
- 441

## 442 Checklist

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

- 447
- Did you include the license to the code and datasets? **[Yes]** See Section ??.
  - 448 • Did you include the license to the code and datasets? **[No]** The code and the data are  
 449 proprietary.
  - 450 • Did you include the license to the code and datasets? **[N/A]**

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

454 1. For all authors...

- 455 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
 456 contributions and scope? **[Yes]**
- 457 (b) Did you describe the limitations of your work? **[Yes]**

- 458 (c) Did you discuss any potential negative societal impacts of your work? [N/A]  
459 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
460 them? [Yes]
- 461 2. If you are including theoretical results...
- 462 (a) Did you state the full set of assumptions of all theoretical results? [Yes]  
463 (b) Did you include complete proofs of all theoretical results? [Yes]
- 464 3. If you ran experiments...
- 465 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
466 mental results (either in the supplemental material or as a URL)? [Yes]  
467 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
468 were chosen)? [Yes]  
469 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
470 ments multiple times)? [Yes]  
471 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
472 of GPUs, internal cluster, or cloud provider)? [Yes]
- 473 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 474 (a) If your work uses existing assets, did you cite the creators? [Yes]  
475 (b) Did you mention the license of the assets? [Yes]  
476 (c) Did you include any new assets either in the supplemental material or as a URL? [No]  
477  
478 (d) Did you discuss whether and how consent was obtained from people whose data  
479 you're using/curating? [N/A]  
480 (e) Did you discuss whether the data you are using/curating contains personally identifi-  
481 able information or offensive content? [N/A]
- 482 5. If you used crowdsourcing or conducted research with human subjects...
- 483 (a) Did you include the full text of instructions given to participants and screenshots, if  
484 applicable? [N/A]  
485 (b) Did you describe any potential participant risks, with links to Institutional Review  
486 Board (IRB) approvals, if applicable? [N/A]  
487 (c) Did you include the estimated hourly wage paid to participants and the total amount  
488 spent on participant compensation? [N/A]

489 **A Appendix**

490 **A.1 Split Learning with R<sup>3</sup>eLU**

491 We give the forward propagation procedure of SplitNN using our R<sup>3</sup>eLU in Algorithm 1, and the  
 492 backward propagation procedure of SplitNN using our R<sup>3</sup>eLU-Diff in Algorithm 2. The two algo-  
 rithms are supplementary materials for Section 3.2 and 3.3.

---

**Algorithm 1:** forward propagation with R<sup>3</sup>eLU

---

**Input:** batch size  $N$ , data batch  $\mathbf{x}^g, \mathbf{x}^h$ , encoded label  $\mathbf{y}$ , training indicator  $I$ , noise scale  $\sigma$ .

**Output:** loss or prediction.

```

1  $\theta^h \leftarrow \mathcal{N}(0, 1), \theta^k \leftarrow \mathcal{N}(0, 1), \theta^g \leftarrow \mathcal{N}(0, 1)$  // initial
Guest:
2 for  $i \leftarrow 1$  to  $N$  do
3 |  $v_i^g \leftarrow \text{Embedding}(x_i^g);$ 
4 end
5  $\mathbf{v}^g \leftarrow \{v_1^g, v_2^g, \dots, v_N^g\}$ 
6  $\mathbf{a}^g \leftarrow \text{Forward-R}^3\text{eLU}(\mathbf{v}^g, \theta^g)$  // send  $\mathbf{a}^g$  to the host
Host:
7 for  $i \leftarrow 1$  to  $N$  do
8 |  $v_i^h \leftarrow \text{Embedding}(x_i^h);$ 
9 end
10  $\mathbf{v}^h \leftarrow \{v_1^h, v_2^h, \dots, v_N^h\}$ 
11  $\mathbf{a}^h \leftarrow \text{Forward}(\mathbf{v}^h, \theta^h)$  // wait for  $\mathbf{a}^g$  from the guest
12  $\bar{\mathbf{a}} \leftarrow \text{Average}(\mathbf{a}^g, \mathbf{a}^h)$ 
13  $\sigma^k \leftarrow \text{Forward}(\bar{\mathbf{a}}, \theta^k)$ 
14 if  $I == \text{train}$  then
15 |  $\text{loss} \leftarrow \mathcal{L}(\sigma^k, \mathbf{y})$  // training
16 else
17 |  $\text{pred} \leftarrow \text{Softmax}(\sigma^k)$  // predicting
18 end

```

---

493

---

**Algorithm 2:** backward propagation with R<sup>3</sup>eLU-Diff

---

**Input:** batch size  $N$ , feature vectors  $\mathbf{v}^g, \mathbf{v}^h$ , forward results  $\mathbf{a}^g, \mathbf{a}^h, \bar{\mathbf{a}}$ , total loss  $l$ , encoded  
 label  $\mathbf{y}$ , learning rate  $\eta^g, \eta^h, \eta^k$ , clipping bound  $C$ , noise scale  $\sigma$ .

**Output:** updated parameters  $\theta_{t+1}^g, \theta_{t+1}^h, \theta_{t+1}^k$ .

```

Host:
1  $\delta^k \leftarrow \nabla_{\mathbf{a}^k} \mathcal{L}(l, \bar{\mathbf{a}}, \theta_t^k)$ 
2  $\mathbf{g}_t^k \leftarrow \text{Backward}(\delta^k, \bar{\mathbf{a}}, \theta_t^k)$ 
3  $\delta^h \leftarrow \nabla_{\mathbf{a}^h} \mathcal{L}(\delta^k, \mathbf{a}^h, \theta_t^h), \delta^g \leftarrow \nabla_{\mathbf{a}^g} \mathcal{L}(\delta^k, \mathbf{a}^g, \theta_t^k)$ 
4  $\mathbf{g}_t^h \leftarrow \text{Backward}(\delta^h, \mathbf{v}^h, \theta_t^h)$ 
5  $\bar{\mathbf{g}}_t^k \leftarrow \frac{1}{N} \sum_{i \in [1, N]} g_{i,t}^k, \bar{\mathbf{g}}_t^h \leftarrow \frac{1}{N} \sum_{i \in [1, N]} g_{i,t}^h$ 
6  $\theta_{t+1}^k \leftarrow \theta_t^k - \eta^k \bar{\mathbf{g}}_t^k, \theta_{t+1}^h \leftarrow \theta_t^h - \eta^h \bar{\mathbf{g}}_t^h$  // host updates
7  $\hat{\delta}^g \leftarrow \delta^g / \max(1, \frac{\|\delta^g\|_1}{C})$ 
8  $\tilde{\delta}^g \leftarrow \text{R}^3\text{eLU-Diff}(\hat{\delta}^g, \mathbf{a}^g, \mathbf{v}^g)$  // send  $\tilde{\delta}^g$  to the guest
Guest:
9  $\mathbf{g}_t^g \leftarrow \text{Backward}(\tilde{\delta}^g, \mathbf{a}^g, \theta_t^g)$ 
10  $\bar{\mathbf{g}}_t^g \leftarrow \frac{1}{N} \sum_{i \in [1, N]} g_{i,t}^g$ 
11  $\theta_{t+1}^g \leftarrow \theta_t^g - \eta^g \bar{\mathbf{g}}_t^g$  // guest updates

```

---

494 **A.2 Privacy Analysis**

495 Here we give the proof of our Corollary 1 in detail for the privacy analysis in Section 3.5.

496 *Proof.* Without loss of generality, we assume an activation result  $\mathbf{a} = \{a_1 = 0, a_2 = 0, \dots, a_o =$   
 497  $0, a_{o+1} > 0, \dots, a_{N_s} > 0\}$ . Then the difference of activations for two neighboring feature vectors  
 498  $v, v'$  can be bounded by

$$\begin{aligned}
 & \frac{P(\mathbf{a}|v)}{P(\mathbf{a}|v')} \\
 &= \prod_{i=1}^o \frac{p + \frac{1-p}{2} \exp(-\frac{|f(x)_i|}{b})}{p + \frac{1-p}{2} \exp(-\frac{|f(v')_i|}{b})} \times \prod_{i=o+1}^{N_s} \frac{\frac{1-p}{2b} \exp(-\frac{|a_i - f(v)_i|}{b})}{\frac{1-p}{2b} \exp(-\frac{|a_i - f(v')_i|}{b})} \\
 &= \prod_{i=1}^o \frac{p \exp(\frac{|f(v')_i|}{b}) + \frac{1-p}{2} \exp(\frac{|f(v')_i| - |f(v)_i|}{b})}{p \exp(\frac{|f(v')_i|}{b}) + \frac{1-p}{2}} \times \\
 & \quad \prod_{i=o+1}^{N_s} \exp(\frac{|a_i - f(v')_i| - |a_i - f(v)_i|}{b}) \\
 &\leq \prod_{i=1}^o \frac{p \exp(\frac{|f(v')_i|}{b}) + \frac{1-p}{2} \exp(\frac{|f(v')_i - f(v)_i|}{b})}{p \exp(\frac{|f(v')_i|}{b}) + \frac{1-p}{2}} \times \prod_{i=o+1}^{N_s} \exp(\frac{|f(v')_i - f(v)_i|}{b}) \\
 &\leq \prod_{i=1}^o \frac{p \exp(\frac{|f(v')_i|}{b}) \exp(-\frac{|f(v)_i|}{b}) + \frac{1-p}{2} \exp(\frac{|f(v')_i - f(v)_i|}{b})}{p \exp(\frac{|f(v')_i|}{b}) \exp(-\frac{|f(v)_i|}{b}) + \frac{1-p}{2}} \times \\
 & \quad \prod_{i=o+1}^{N_s} \exp(\frac{|f(v')_i - f(v)_i|}{b}) \\
 &\leq \prod_{i=1}^o \frac{\frac{1+p}{2}}{p \exp(\frac{|f(v')_i| - |f(v)_i|}{b}) + \frac{1-p}{2}} \times \prod_{i=o+1}^{N_s} \exp(\frac{|f(v')_i - f(v)_i|}{b})
 \end{aligned}$$

499 We assume  $|f(v)_i| \leq c_g$  and  $b \leq \frac{c_g}{\ln 2}$ , then the RHS

$$\begin{aligned}
 &\leq \prod_{i=1}^o \frac{\frac{1+p}{2}}{p \exp(\frac{-c_g}{b}) + \frac{1-p}{2}} \times \prod_{i=o+1}^{N_s} \exp(\frac{|f(v')_i - f(v)_i|}{b}) \\
 &\leq \prod_{i=1}^o \frac{\frac{1+p}{2}}{\exp(\frac{-c_g}{b})} \times \prod_{i=o+1}^{N_s} \exp(\frac{|f(v')_i - f(v)_i|}{b}) \\
 &\leq (\frac{1+p}{2 \exp(\frac{-c_g}{b})})^o \times \exp(\frac{\Delta f}{b}) \\
 &= \exp(\epsilon)
 \end{aligned}$$

500 Thus, we have  $N_S(\frac{c_g}{\sigma_g} + \ln(\frac{1+p}{2})) + \frac{\Delta f}{\sigma_g} = \epsilon$ . □

501 **A.3 Evaluation of Dynamic Parameter Importance Estimation**

502 Additionally, we evaluate the effectiveness of our feature importance estimation method and give  
 503 the result here. We compare the estimating results between parameter importance estimated in the  
 504 model's finally stable state and parameter importance estimated in our dynamic manner. The result  
 505 is shown in Figure 3. We can conclude from the result that our dynamic parameter importance  
 506 estimation approach can achieve desirable effectiveness.

507 We aim to estimate the importance of parameters precisely during privacy-preserving split learning.  
 508 But it is difficult to get the same result with the estimation in a non-perturbed case. By carefully  
 509 constructing the dynamic estimation method, we can obtain dynamic estimation results of noisy pa-  
 510 rameters quite close to a baseline estimation result in the stable state without any privacy protection.  
 511 In Figure 4, we show the result of our dynamic importance estimation approach. We can find that

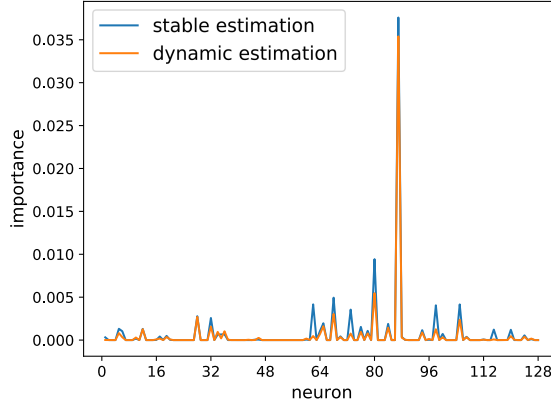


Figure 3: Comparison of parameter importance estimation results.

512 artificial noise introduced by our privacy-preserving SplitNN solution will affect the estimation of  
 513 parameter importance. But the growth tendency keep the same as the baseline. This is good enough  
 514 for us since we use the proportionality factor of each parameter to calculate the budget allocation,  
 which will not be influenced by the magnitude.

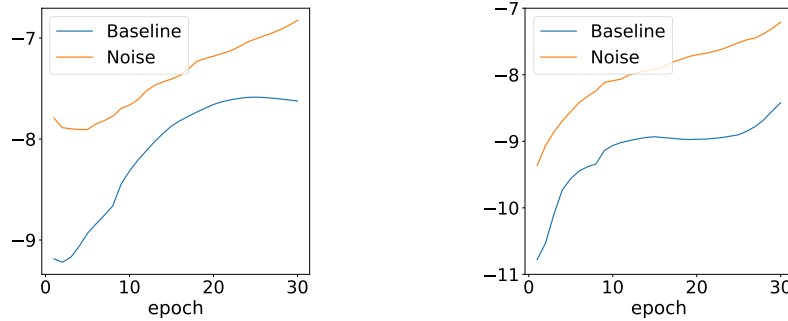


Figure 4: Importance of noisy parameters estimated dynamically in logarithm.

515

#### 516 A.4 Supplementary Results for Evaluation

517 The neural networks we used for the MovieLens, BookCrossing, and MNIST datasets after split are  
 518 shown in Tab 8, Table 9, and Table 10. These networks are commonly used in related studies. We  
 519 split them according to the interpretation of SplitNN in [4, 11, 29].

Guest Layer	Output Shape	Param #
Linear(160,128)+ReLU	(None,128)	20608
Host Layer	Output Shape	Param #
Linear(160,128)+ReLU	(None,128)	20608
Merge Guest		
Linear(128,128)+ReLU	(None,128)	16512
Linear(128,64)+ReLU	(None,64)	8256
Linear(64,3952)+Softmax	(None,3952)	256880

Table 8: The MovieLens model.

520 To further investigate how our solution affects the learning process of SplitNN, we report learning  
 521 results of a MovieLens recommendation model protecting the privacy of the guest and the host  
 522 in Figure 5 and 6, respectively. In each plot, we show trends of training loss, training accuracy,  
 523 and testing accuracy as the training epoch increases. We can conclude from these figures that our  
 524 solution achieve satisfying model usability with small privacy budget for either side of SplitNN.



<b>Guest Layer</b>	<b>Output Shape</b>	<b>Param #</b>
Linear(160,128)+ReLU	(None,128)	20608

<b>Host Layer</b>	<b>Output Shape</b>	<b>Param #</b>
Linear(160,128)+ReLU	(None,128)	20608
Merge Guest		
Linear(128,256)+ReLU	(None,256)	33024
Linear(256,128)+ReLU	(None,128)	32896
Linear(128,17384)+Softmax	(None,10)	2242536

Table 9: The BookCrossing model.

<b>Guest Layer</b>	<b>Output Shape</b>	<b>Param #</b>
Linear(28*14,128)	(None,128)	50304
BatchNorm+ReLU	(None,128)	
Linear(128,64)	(None,64)	8256

<b>Host Layer</b>	<b>Output Shape</b>	<b>Param #</b>
Linear(28*14,128)	(None,128)	50304
BatchNorm+ReLU	(None,128)	
Linear(128,64)	(None,64)	8256
BatchNorm+ReLU	(None,64)	
Merge Guest		
Linear(64,64)	(None,64)	4160
BatchNorm+ReLU	(None,64)	
Linear(64,10)+Softmax	(None,10)	650

Table 10: The MNIST model.

525 We notice that FSHA is sensitive to our privacy-preserving SplitNN solution. In Figure 1 and Fig-  
526 ure 2, we find that neither an adversarial host nor an adversarial guest can reconstruct meaningful  
527 samples, even if the target’s privacy budget is  $\epsilon = 4$ . Therefore, we are curious about a practical  
528 choice of the privacy budget when dealing with FSHA. To this end, we give more defense results  
529 of our solution against FSHA using various privacy budget values in Table 11. We can conclude  
530 from the MSE results that two recommendation models prefer relatively low privacy budget, such  
531 as  $\epsilon = 1.0$ . However, it is interesting to see that FSHA attack against an image classification model  
532 using SplitNN can be frustrated by our solution using relatively high privacy budget, which also  
means a high model usability.

Table 11: Defense results against FSHA mounted by an adversarial host.

$\epsilon$	MovieLens		BookCrossing		MNIST	
	baseline	<b>Ours</b>	baseline	<b>Ours</b>	baseline	<b>Ours</b>
	1200 Epochs	5000 Epochs	1200 Epochs	5000 Epochs	9000 Epochs	9000 Epochs
0.1		<b>455.5676</b>		<b>500.1487</b>		<b>1.98257</b>
0.25		<b>73.7824</b>		<b>83.7606</b>		<b>1.9788</b>
0.5		<b>21.1706</b>		<b>30.2905</b>		<b>1.9703</b>
0.75		<b>9.5312</b>		<b>8.5984</b>		<b>1.9534</b>
1.0	$0.2652 \times 10^{-3}$	<b>4.4903</b>	$0.2365 \times 10^{-3}$	<b>6.1267</b>	0.0206	<b>1.9442</b>
2.0		<b>1.1559</b>		<b>1.7047</b>		<b>1.9283</b>
4.0		<b>0.2719</b>		<b>0.2478</b>		<b>1.9209</b>
6.0		<b>0.1091</b>		<b>0.1638</b>		<b>1.9135</b>
8.0		<b>0.0975</b>		<b>0.0846</b>		<b>1.9116</b>

533

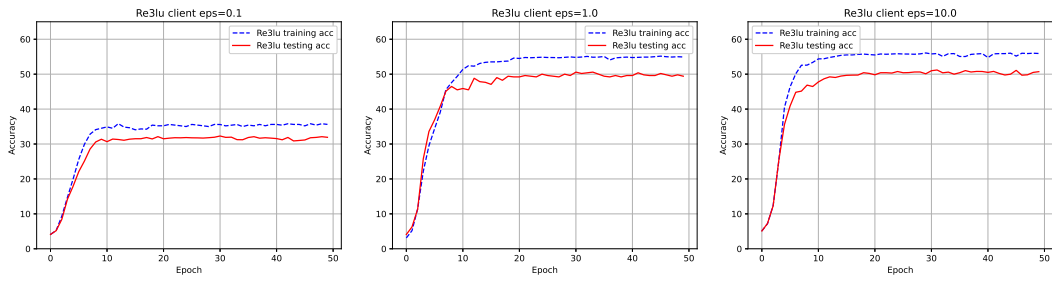


Figure 5: SplitNN learning curve with the guest's privacy protected by our solution.

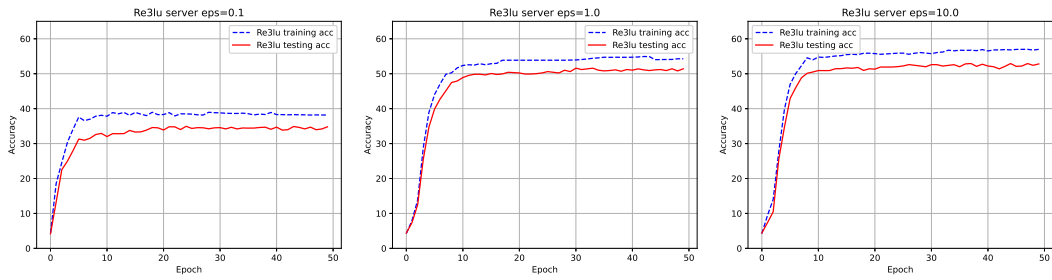


Figure 6: SplitNN learning curve with the host's privacy protected by our solution.