

• Supplementary File •

# Efficient Privacy-Preserving Federated Learning under Dishonest-Majority Setting

Yinbin Miao<sup>1</sup>, Da Kuang<sup>1</sup>, Xinghua Li<sup>1</sup>, Tao Leng<sup>2,3,4\*</sup>, Ximeng Liu<sup>5</sup> & Jianfeng Ma<sup>1</sup>

<sup>1</sup>*School of Cyber Engineering, Xidian University, Xi'an 710071, China;*

<sup>2</sup>*Intelligent Policing Key Laboratory of Sichuan Province, Sichuan Police College, Luzhou 646000, China;*

<sup>3</sup>*Institute of Information Engineering, Chinese Academy of Sciences, Beijing 100085, China;*

<sup>4</sup>*School of Cyber Security, University of Chinese Academy of Sciences, Beijing 100049, China;*

<sup>5</sup>*Key Laboratory of Information Security of Network Systems, Fuzhou University, Fuzhou 350108, China*

## Appendix A Technical

### Appendix A.1 Preprocessing Operation

Sample-based Secret Sharing. The method is designed according to the characteristics of computing Hamming distance, which can efficiently compute Hamming distance under ciphertext and also support weighted average without the process of Bit2A. And the process is shown in Algorithm A1. Given  $n$  participants  $\{P_i | i \in [0, n-1]\}$  and a secret distributor  $D$  holding a binary vector  $M = (m_0, \dots, m_{k-1})_2$  of length  $k$ ,  $D$  first initializes  $n$  vectors  $\llbracket m \rrbracket_i$  of length  $k$  with 0, then initializes a random seed locally to generate a random sequence  $Ind = (Ind_0, \dots, Ind_{k-1})$  with a length of  $k$ , where  $Ind_{k-1} \in [0, k-1]$  and each  $Ind_{k-1}$  is not repeated.<sup>1)</sup> Then  $D$  divides  $Ind$  into  $n$  parts as  $P_i$ 's position set  $ind_i = (ind_{i,0}, \dots, ind_{i,j})$  with a length of  $j$ , where  $j = k/n$ .  $D$  then samples  $M$  according to  $ind_i$ , that is, for the secret share  $\llbracket m \rrbracket_i$ , assign  $\llbracket M \rrbracket [ind_{i,j}]$  to the corresponding  $\llbracket m \rrbracket_i [ind_{i,j}]$ , and the other values in  $\llbracket m \rrbracket_i$  remain unchanged and are still equal to 0 (Lines 3-9 in Algorithm A1). Finally  $D$  distributes  $\llbracket m \rrbracket_i$  to each participant  $P_i$ .

Next, we take two parties as an example to explain the process of generating secret shares. Given a secret  $M = [100101]_2$ ,  $D$  will generate two secret shares  $\llbracket m \rrbracket_0$  and  $\llbracket m \rrbracket_1$ . First  $D$  initializes  $\llbracket m \rrbracket_0 = [000000]$ ,  $\llbracket m \rrbracket_1 = [000000]$ , then scrambles the sequence  $[0, 1, 2, 3, 4, 5]$  and divides it into 2 position sets,  $ind_0 = [0, 1, 5]$ ,  $ind_1 = [2, 3, 4]$ . Then  $D$  samples  $M$ , then we have  $\llbracket m \rrbracket_0 [ind_{0,0}] = M[ind_{0,0}] \rightarrow \llbracket m \rrbracket_0 [0] = M[0] = 1$  and  $\llbracket m \rrbracket_0 [1] = 0$ ,  $\llbracket m \rrbracket_0 [5] = 1$  can be deduced in the same way. And the value of  $\llbracket m \rrbracket_0 [2]$ ,  $\llbracket m \rrbracket_0 [3]$  and  $\llbracket m \rrbracket_0 [4]$  are still equal to 0, and finally we have  $\llbracket m \rrbracket_0 = [100001]$ ,  $\llbracket m \rrbracket_1 = [000100]$ .

---

#### Algorithm A1 Gradient Sharing GS( $\cdot$ )

---

**Input:** Secret  $M = (m_0, \dots, m_{k-1})_2$

**Output:** Secret sharing  $\llbracket m \rrbracket_0, \dots, \llbracket m \rrbracket_{n-1}$

- |   |  |
|---|--|
| <p>1: The secret distributor <math>D</math> initializes <math>n</math> secret shares <math>\llbracket m \rrbracket_i = (0_0, \dots, 0_{k-1})_2</math> of <math>M</math>;</p> <p>2: <math>D</math> randomly reorders the sequence with values from 0 to <math>k-1</math> to get <math>Ind</math>, and then divides <math>Ind</math> into <math>n</math> parts of position set <math>ind_i = (ind_{i,0}, \dots, ind_{i,j})</math> for each <math>P_i</math> on average.</p> | <p>3: <b>for</b> <math>i = 0, 1, \dots, n-1</math> <b>do</b></p> <p>4:   <b>for</b> <math>t = 0, 1, \dots, k-1</math> <b>do</b></p> <p>5:     <b>if</b> <math>t \in ind_i</math> <b>then</b></p> <p>6:       <math>\llbracket m \rrbracket_i [t] \leftarrow \llbracket M \rrbracket [t]</math>;</p> <p>7:     <b>end if</b></p> <p>8:   <b>end for</b></p> <p>9: <b>end for</b></p> <p>10: <b>return</b> <math>\llbracket m \rrbracket_0, \dots, \llbracket m \rrbracket_{n-1}</math>.</p> |
|---|--|
- 

---

#### Algorithm A2 Preprocessing Operation PO( $\cdot$ )

---

**Input:** Current global model  $w^\tau$

**Output:** Preprocessed gradient  $\llbracket \nabla w_i^e \rrbracket_0, \llbracket \nabla w_i^e \rrbracket_1$

- |   |  |
|---|--|
| <p>1: <math>w_i^\tau \leftarrow \text{LocalUpdate}(w^\tau)</math></p> <p>2: <math>\nabla w_i^\tau = w^\tau - w_i^\tau</math></p> <p>3: <b>if</b> <math>\nabla w_i^\tau &gt; 0</math> <b>then</b></p> <p>4:   <math>\nabla w_i^q \leftarrow 1</math></p> <p>5: <b>else</b></p> | <p>6:   <math>\nabla w_i^q \leftarrow -1</math></p> <p>7: <b>end if</b></p> <p>8: <math>\nabla w_i^e \leftarrow [(1 - \nabla w_i^q)/2]</math></p> <p>9: <math>\llbracket \nabla w_i^e \rrbracket_0, \llbracket \nabla w_i^e \rrbracket_1 \leftarrow F.GS(\nabla w_i^e)</math></p> <p>10: Send <math>\llbracket \nabla w_i^e \rrbracket_0, \llbracket \nabla w_i^e \rrbracket_1</math> to <math>\mathcal{S}_0</math> and <math>\mathcal{S}_1</math> respectively</p> <p>11: <b>return</b> <math>\llbracket \nabla w_i^e \rrbracket_0, \llbracket \nabla w_i^e \rrbracket_1</math></p> |
|---|--|
- 

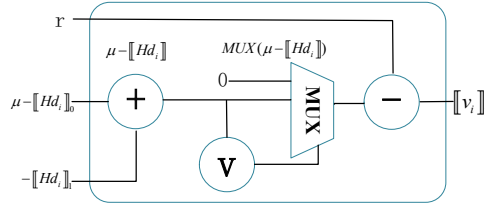
\* Corresponding author (email: lengtao@iie.ac.cn)

1) Here we use the Fisher-Yates shuffling algorithm to generate a random sequence.

## Appendix A.2 GMW

We use GMW protocol [1] to construct Garbled Circuit (GC) [2] for two-party communication. The GMW protocol consists of three stages and its objective function is composed of XOR gates, AND gates and NOT gates. Assume there are two parties  $P_0, P_1$  holding two bit strings  $x = (x_0, \dots, x_{k-1})_2, y = (y_0, \dots, y_{k-1})_2$  of length  $k$  respectively. The implementation of GMW is described below.

1. *Circuit generation.* The generator needs to convert functions into circuits composed of XOR gates, AND gates and NOT gates.
2. *Secret sharing.*  $P_0$  and  $P_1$  generate secret shares  $\llbracket x \rrbracket = (\llbracket x_0 \rrbracket, \dots, \llbracket x_{k-1} \rrbracket), \llbracket y \rrbracket = (\llbracket y_0 \rrbracket, \dots, \llbracket y_{k-1} \rrbracket)$  of  $x$  and  $y$  respectively and then send one of the shared shares to each other.
3. *Circuit execution.*  $P_0$  and  $P_1$  input secret shares into the circuit respectively and then execute each gate in the circuit in sequence. We describe their execution processes for different gates, where  $\llbracket z \rrbracket = (\llbracket z_0 \rrbracket, \dots, \llbracket z_{k-1} \rrbracket)$  is the value to be calculated.
  - $P_1$  executes the XOR gate by calculating  $\llbracket z_i \rrbracket_j = \llbracket x_i \rrbracket_j \oplus \llbracket y_i \rrbracket_j$  locally, where  $0 \leq i \leq k-1, j = 0, 1$ .
  - $P_0$  executes NOT gate by calculating  $\llbracket z_i \rrbracket_j = \llbracket x_i \rrbracket_j \oplus j$  locally.
  - To execute the AND gate,  $P_0$  first lists the truth table and then  $P_0$  and  $P_1$  perform the calculation together through OT protocol.



**Figure A1** Garbled circuit for calculating client weight.  $\llbracket v_i \rrbracket_0 = r, \llbracket v_i \rrbracket_1 = MUX(\mu - \llbracket Hd_i \rrbracket) - r$ , where  $r$  is a random number generated by  $S_0$  and  $MUX()$  is a two-to-one data selector. When the input parameter is greater than 0, garbled circuit return  $\mu - \llbracket Hd_i \rrbracket$ , otherwise the output is 0.

## Appendix A.3 Generate Multiplication Triples

we use Correlated Oblivious Product Evaluation (COPE) [3] to securely generate  $\llbracket a \rrbracket, \llbracket b \rrbracket, \llbracket c \rrbracket$ .

COPE is a protocol based on OT. Its purpose is to allow multiple participants to obtain the secret shares of the product of all input values through COPE after their respective inputs, note that the data held by each participant can only be known by themselves in Algorithm 4. Next, we will take two participants as an example to illustrate the whole process in detail. Suppose there are two parties  $P_0, P_1$  holding  $a, b \in F_p$  respectively and they aim to calculate  $\llbracket c \rrbracket = \llbracket a \rrbracket * \llbracket b \rrbracket$ , where  $p$  is a large prime number.

1.  $P_0$  takes  $a$  as input and  $P_1$  first converts  $b$  to binary stream  $(b_0, b_1, \dots, b_{k-1}) \in (0, 1)^k, b = \sum_i^k b_i * 2^{i-1}$  and then takes it as input.
2.  $P_0$  and  $P_1$  execute the OT protocol  $k$  times. Note that  $P_0$  is the sender,  $P_1$  is the receiver,  $r_i$  is a random value selected by  $P_0$  from  $F_p$  in each round of OT protocol (Line 4 in Algorithm A3).
3.  $P_0$  and  $P_1$  compute the secret shares  $\llbracket c \rrbracket_0, \llbracket c \rrbracket_1$  of  $\llbracket c \rrbracket$  respectively (Lines 6-7 in Algorithm A3).

---

### Algorithm A3 COPE( $\cdot$ )

---

**Input:**  $P_0$  has  $a, P_1$  has  $b$

**Output:**  $\llbracket c \rrbracket = a * b$

- |   |   |
|---|---|
| <ol style="list-style-type: none"> <li>1: <math>P_1</math> converts <math>b</math> to binary <math>(b_0, b_1, \dots, b_{k-1})</math>;</li> <li>2: <b>for</b> <math>i = 0, 1, \dots, k-1</math> <b>do</b></li> <li>3:     <math>P_0</math> generates a random number <math>r_i</math>;</li> <li>4:     <math>P_0</math> and <math>P_1</math> jointly execute <math>q_i \leftarrow OT(r_i, r_i + a, b_i)</math>;</li> </ol> | <ol style="list-style-type: none"> <li>5: <b>end for</b></li> <li>6: <math>P_0</math> computes <math>\llbracket c \rrbracket_0 \leftarrow - \sum_i^k t_i * 2^{i-1}</math> as a secret share of <math>\llbracket c \rrbracket</math>;</li> <li>7: <math>P_1</math> computes <math>\llbracket c \rrbracket_1 \leftarrow \sum_i^k q_i * 2^{i-1}</math> as a secret share of <math>\llbracket c \rrbracket</math>;</li> <li>8: <b>return</b> <math>\llbracket c \rrbracket</math>.</li> </ol> |
|---|---|
- 

Finally, in the offline phase, we use COPE to generate  $\llbracket a \rrbracket, \llbracket b \rrbracket, \llbracket c \rrbracket$  as shown in Algorithm A4 and the steps are as follows:

1. Each participant  $P_i$  randomly generates a random number  $\llbracket a \rrbracket_i, \llbracket b \rrbracket_i \in F_p$ , where  $\llbracket a \rrbracket_i, \llbracket b \rrbracket_i$  are the secret shares of  $\llbracket a \rrbracket, \llbracket b \rrbracket$  respectively.
2. The COPE protocol is implemented between two different parties. Participants  $P_i, P_j$  respectively take  $\llbracket a \rrbracket_i, \llbracket b \rrbracket_j$  as the inputs of COPE,  $P_i, P_j$  obtain  $\llbracket t \rrbracket_{i,j}, \llbracket q \rrbracket_{j,i}$  respectively (Line 5 in Algorithm A4).
3. Each  $P_i$  calculates the secret share of  $\llbracket c \rrbracket$  by Eq.A1 (Line 9 in Algorithm A4).

$$\llbracket c \rrbracket_i = \llbracket a \rrbracket_i * \llbracket b \rrbracket_i + \sum_j^n (\llbracket t \rrbracket_{i,j} + \llbracket q \rrbracket_{j,i}). \quad (\text{A1})$$

---

**Algorithm A4** Generate Multiplication Triples GMT( $\cdot$ )
 

---

**Input:**  $\{P_i | 0 \leq i \leq n-1\}$

**Output:** Multiplication triples  $\llbracket a \rrbracket, \llbracket b \rrbracket, \llbracket c \rrbracket$

1: Each $P_i$ generates a pair of random numbers $(\llbracket a \rrbracket_i, \llbracket b \rrbracket_i)$ as secret share of $\llbracket a \rrbracket, \llbracket b \rrbracket$ ; 2: <b>for</b> $i = 0, \dots, n-1$ <b>do</b> 3: <b>for</b> $j = 0, \dots, n-1$ <b>do</b> 4: <b>if</b> $i \neq j$ <b>then</b>	5: $\llbracket t \rrbracket_{i,j}, \llbracket q \rrbracket_{j,i} \leftarrow \text{COPE}(\llbracket a \rrbracket_i, \llbracket b \rrbracket_j)$ ; 6: <b>end if</b> 7: <b>end for</b> 8: <b>end for</b> 9: Each $P_i$ gets $\llbracket c \rrbracket_i \leftarrow \llbracket a \rrbracket_i * \llbracket b \rrbracket_i + \sum_j \left( \llbracket t \rrbracket_{i,j} + \llbracket q \rrbracket_{j,i} \right)$ ; 10: <b>return</b> $\llbracket a \rrbracket, \llbracket b \rrbracket, \llbracket c \rrbracket$ .
--	--

---

## Appendix A.4 secMul

Secret sharing [4] is a technique for sharing secrets among  $n$  participants  $\{P_i | i \in [0, n-1]\}$ . Specifically, given an additive secret sharing  $\llbracket x \rrbracket = \sum_{i=0}^{n-1} \llbracket x \rrbracket_i$ , where  $\llbracket x \rrbracket_i$  is the share held by participant  $P_i$ , an additive secret sharing usually has the following four security operations.

- Given secret shares  $\llbracket x \rrbracket = \sum_{i=0}^{n-1} \llbracket x \rrbracket_i, \llbracket y \rrbracket = \sum_{i=0}^{n-1} \llbracket y \rrbracket_i$ , we have  $\llbracket z \rrbracket = \llbracket x \rrbracket + \llbracket y \rrbracket \rightarrow \sum_{i=0}^{n-1} (\llbracket x \rrbracket_i + \llbracket y \rrbracket_i)$ .
- Given two secret shares  $\llbracket x \rrbracket, \llbracket y \rrbracket$ , we can get  $\llbracket z \rrbracket = \llbracket x \rrbracket * \llbracket y \rrbracket$  by multiplying triples [34] as shown in Algorithm A5, where  $\llbracket a \rrbracket, \llbracket b \rrbracket, \llbracket c \rrbracket$  are multiplicative triples.
- Given a secret share  $\llbracket x \rrbracket = \sum_{i=0}^{n-1} \llbracket x \rrbracket_i$  and a constant  $c$ , we have  $\llbracket z \rrbracket = \llbracket x \rrbracket + c \rightarrow \sum_{i=0}^{n-1} \llbracket x \rrbracket_i + c$ .
- Given a secret share  $\llbracket x \rrbracket$  and a constant  $c$ , we have  $\llbracket z \rrbracket = \llbracket x \rrbracket * c \rightarrow \sum_{i=0}^{n-1} (\llbracket x \rrbracket_i * c)$ .

---

**Algorithm A5** Secure Multiplication SecMul( $\cdot$ )
 

---

**Input:** Secret shares  $\llbracket x \rrbracket, \llbracket y \rrbracket, \llbracket a \rrbracket, \llbracket b \rrbracket, \llbracket c \rrbracket$

**Output:** Secret share  $\llbracket z \rrbracket = \llbracket x \rrbracket * \llbracket y \rrbracket$

1: <b>for</b> $i \in [0, n-1]$ <b>do</b> 2: $P_i$ computes $\llbracket \varepsilon \rrbracket_i = \llbracket x \rrbracket_i - \llbracket a \rrbracket_i$ and $\llbracket \rho \rrbracket_i = \llbracket y \rrbracket_i - \llbracket b \rrbracket_i$ ; 3:     Specify $P_0$ as the main participant; 4: $P_i$ receives the pair $(\llbracket \varepsilon \rrbracket_j, \llbracket \rho \rrbracket_j)$ from $P_j, j \neq i$ ;	5: <b>if</b> $i == 0$ <b>then</b> 6: $\llbracket z \rrbracket_i \leftarrow \llbracket c \rrbracket_i + \llbracket \rho \rrbracket * \llbracket a \rrbracket_i + \llbracket \varepsilon \rrbracket * \llbracket b \rrbracket_i + \llbracket \varepsilon \rrbracket * \llbracket \rho \rrbracket$ ; 7: <b>else</b> 8: $\llbracket z \rrbracket_i \leftarrow \llbracket c \rrbracket_i + \llbracket \rho \rrbracket * \llbracket a \rrbracket_i + \llbracket \varepsilon \rrbracket * \llbracket b \rrbracket_i$ ; 9: <b>end if</b> 10: <b>end for</b> 11: <b>return</b> $\llbracket z \rrbracket_i$ .
--	---

---

## Appendix B Theoretical Analysis

### Appendix B.1 Robustness Analysis

Similar to [5], we give a robustness analysis of Hamming distance. Cosine similarity is one of the best metrics for measuring the similarity between two vectors and here we show that using Hamming distance to detect malicious models can achieve equivalent results as cosine similarity by Eq.B1, where  $cs_i$  is the cosine similarity of the gradient  $g_i^q$  and the root gradient  $g_s^q$ .

$$cs_i = \frac{\langle \nabla w_i^q, \nabla w_s^q \rangle}{\|\nabla w_i^q\| \cdot \|\nabla w_s^q\|}. \quad (\text{B1})$$

According to Eq.??, we can get Eq.B2, where  $d$  is the length of  $\nabla w_i$ .

$$\begin{aligned}
 cs_i &= \frac{\sum_{j=1}^d \nabla w_{i,j}^q \cdot \nabla w_{s,j}^q}{\sqrt{d} \cdot \sqrt{d}} = \frac{1}{d} \cdot \left( \sum_{j=1}^d (1 - 2 * \nabla w_{i,j}^e) \right) \cdot \left( 1 - 2 * \nabla w_{s,j}^e \right) \\
 &= 1 - \frac{2}{d} \cdot \left( \sum_{j=1}^d \nabla w_{i,j}^e \oplus \nabla w_{s,j}^e \right) = 1 - 2 \cdot \frac{Hd_i}{d}.
 \end{aligned} \quad (\text{B2})$$

Therefore, we can deduce that when the threshold is  $\tau = d/2$ , the client weight  $v_i$  and  $cs_i$  are consistent, implying  $v_i > 0 \Leftrightarrow cs_i > 0$ .

$$cs_i > 0 \Leftrightarrow 1 - 2 \cdot \frac{Hd_i}{d} > 0 \Leftrightarrow Hd_i < \frac{d}{2} \quad (\text{B3})$$

Based on the above, this scheme uses the Hamming distance to detect malicious models to ensure that the robustness of the scheme is feasible. In addition, using Hamming distance will be more flexible than using cosine similarity, and we can vary the value of the threshold  $\tau$  for different task metrics to obtain the best Byzantine robustness.

### Appendix B.2 Security Analysis

Now given two secret distributors  $D_0, D_1$  hold secrets  $M_0 = (M_{0,0}, \dots, M_{0,k})_2, M_1 = (M_{1,0}, \dots, M_{1,k})_2$  respectively, and two participants  $P_0, P_1$ . First,  $D_0$  and  $D_1$  jointly negotiate a random seed through secret key exchange technology and digital signature algorithm, and use the shuffling algorithm to randomly sort the sequence  $[0, \dots, k-1]$  locally, and obtain the same sequence  $Ind$ . After that,  $D_0$  and  $D_1$  divide  $Ind$  into two parts, namely  $ind_0 = (Ind_0, \dots, Ind_s)$  and  $ind_1 = (Ind_{s+1}, \dots, Ind_k)$ , where  $s = \lfloor \frac{k-1}{2} \rfloor$ . Finally,  $D_0$  samples  $M_0$  according to  $ind_0$  and  $ind_1$  to obtain  $\llbracket M_0 \rrbracket_0 = (M_{0,0}, \dots, M_{0,Ind_s}, \dots, 0)_2$  and  $\llbracket M_0 \rrbracket_1 = (0, \dots, M_{0,Ind_{s+1}}, \dots, M_{0,k})_2$ . Similarly,  $D_1$  generates  $\llbracket M_1 \rrbracket_0 = (M_{1,0}, \dots, M_{1,Ind_s}, \dots, 0)_2$  and  $\llbracket M_1 \rrbracket_1 = (0, \dots, M_{1,Ind_{s+1}}, \dots, M_{1,k})_2$ . Among them,  $P_0$  holds  $\llbracket M_0 \rrbracket_0, \llbracket M_1 \rrbracket_0$ , and  $P_1$  holds  $\llbracket M_1 \rrbracket_0, \llbracket M_1 \rrbracket_1$ . We analyze the correctness of the sample-based secret sharing method, as Theorem 1 and Theorem 2.

**Table C1** Model descriptions of MNIST(Fashion-MNIST) and CIFAR-10.

MNIST	Conv <sub>1</sub>	in_channels=1,out_channels=10,kernel_size=5
	Conv <sub>2</sub>	in_channels=10,out_channels=20,kernel_size=5
	FC <sub>1</sub>	in_features=320, out_features=50
	FC <sub>2</sub>	in_features=50, out_features=10
CIFAR-10	Conv <sub>1</sub>	in_channels=3,out_channels=6,kernel_size=5,padding=1
	Pool	kernel_size=2, stride=2
	Conv <sub>2</sub>	in_channels=6,out_channels=16,kernel_size=5,padding=1
	FC <sub>1</sub>	in_features=16*5*5,out_features=120
	FC <sub>2</sub>	in_features=120,out_features=84
	FC <sub>3</sub>	in_features=84,out_features=10

**Theorem 1** (Correctness of additive properties). The sample-based secret sharing methods have all the properties of additive secret sharing, as described in sec. 3.2.

*Proof.* Given  $\llbracket M_0 \rrbracket_0$  and  $\llbracket M_0 \rrbracket_1$ , we have

$$\llbracket M_0 \rrbracket_0 + \llbracket M_0 \rrbracket_1 = (M_{0,0}, \dots, M_{0,Ind_s}, \dots, M_{0,Ind_k})_2 = M_0. \quad (\text{B4})$$

According to Eq.B4, given a constant  $c$ , we have  $(\llbracket M_0 \rrbracket_0 + c) + \llbracket M_0 \rrbracket_1 = (\llbracket M_0 \rrbracket_0 + \llbracket M_0 \rrbracket_1) + c = M_0 + c$  and  $\llbracket M_0 \rrbracket_0 \cdot c + \llbracket M_0 \rrbracket_1 \cdot c = (\llbracket M_0 \rrbracket_0 + \llbracket M_0 \rrbracket_1) \cdot c = M_0 \cdot c$ . From the above, the sampled-secret sharing can be used in Algorithm A5, such as calculating  $\llbracket M_0 \cdot M_1 \rrbracket = \llbracket M_0 \rrbracket \cdot \llbracket M_1 \rrbracket$  through  $\llbracket a \rrbracket, \llbracket b \rrbracket, \llbracket c \rrbracket$ .

**Theorem 2** (Correctness of Hamming distance). The sample-based secret sharing can correctly calculate  $Hd = \llbracket Hd \rrbracket_0 + \llbracket Hd \rrbracket_1$ , where the Hamming distance  $Hd$  is between  $M_0$  and  $M_1$ .

*Proof.* Then we prove Theorem 2. as follows:

$$\begin{aligned} \llbracket Hd \rrbracket_0 + \llbracket Hd \rrbracket_1 &= H(\llbracket M_0 \rrbracket_0 \oplus \llbracket M_1 \rrbracket_0) + H(\llbracket M_0 \rrbracket_1 \oplus \llbracket M_1 \rrbracket_1) \\ &= H(M_{0,0} \oplus M_{1,0}, \dots, M_{0,Ind_s} \oplus M_{1,Ind_s}, \dots, M_{0,Ind_k} \oplus M_{1,Ind_k}) = Hd. \end{aligned} \quad (\text{B5})$$

The security of sample-based secret sharing mainly depends on the randomness and privacy of  $Ind$ . We use secure key exchange technology and digital signatures to negotiate random seeds to ensure that the random seeds will not be leaked, thereby ensuring that  $Ind$  cannot be obtained by the adversary. Secondly, we use a secure shuffling algorithm to ensure the randomness of  $Ind$  so that  $Ind$  can only be cracked violently by adversary with negligible probability.<sup>2)</sup>

In addition, we analyze the privacy of FESD against semi-honest adversaries (servers), In Theorem 3.

**Theorem 3** (Privacy of FESD). FESD guarantees that semi-honest adversaries can learn nothing except what they can infer from the aggregated result  $(\llbracket \nabla w_i^r \rrbracket_0, \llbracket \nabla w_i^r \rrbracket_1)$  with probability  $1 - \varepsilon$ .

*Proof.* Since the security of any combination of different protocols is guaranteed in the general composability framework [6], we only need to prove the security of a single protocol.

*Preprocessing.* In this stage,  $\mathcal{C}$  and  $\mathcal{S}_0$  perform operations locally without interaction. In the real world,  $\mathcal{C}$ 's view is  $\{w^t, \nabla w_s^q, \llbracket \nabla w_i^e \rrbracket, \llbracket \nabla w_i^e \rrbracket_0, \llbracket \nabla w_i^e \rrbracket_1\}$  and  $\mathcal{S}_0$ 's view is  $\{w^t, \nabla w_s^q, \llbracket \nabla w_s^e \rrbracket, \llbracket \nabla w_s^e \rrbracket_0, \llbracket \nabla w_s^e \rrbracket_1\}$ . In an ideal world, since  $\mathcal{C}$  and  $\mathcal{S}_0$  can perform preprocessing locally and  $\text{Sim}$  does not need to simulate this process, their views are consistent with those in the real world, thereby satisfying privacy security. Similarly, in calculating Hamming distance, decoding and model aggregation, each participant performs operations locally without interaction. These stages ensure privacy security by local operations.

*Calculate client weight.* In this stage, each participant obtains their weights by using garbled circuits which are zero-knowledge and secure. The view composed of tags in the real world by each participant is indistinguishable from that simulated by  $\text{Sim}$ . This stage ensures privacy security by garbled circuits.

*Multiplication.* In this stage,  $\mathcal{S}_0$  and  $\mathcal{S}_1$  generate multiplication triples securely through COPE protocol and local operations. In addition, the COPE protocol has proved to be secure in MASCOT [3], so the multiplication process is secure by using multiplication triples.

Therefore, we guarantee that the adversaries learns nothing beyond what can be inferred from the aggregated results with an overwhelming probability.

## Appendix C Performance Analysis

### Appendix C.1 Experiment Setup

This solution is implemented in Python 3.8. In the training model stage, we rely on TorchVision and Torch two libraries to implement. In addition, unless otherwise specified, we test the CNN model structure used in the MNIST (Fashion-MNIST) and CIFAR-10 datasets as described in Table C1.

#### Appendix C.1.1 Dataset

We test our scheme with multiple datasets under Byzantine attacks, including 3 datasets for image classification. It should be noted here that all our experiments are under the IID setting. We distribute datasets in the same way that datasets are distributed among clients in FITrust. Assume that there are  $M$  classes in the dataset, we first randomly divide clients into  $M$  groups and then set the probability of assigning the dataset with the  $i$ -th label to the  $i$ -th group as  $p$ , note that the probability of assigning it to other groups as  $(1 - p)/(M - 1)$ . In the same group, the dataset will be evenly distributed to clients.

The data sets we used in this experiment mainly include MNIST [7], Fashion-MNIST [8] and CIFAR-10 [9]. Next we will briefly introduce root dataset.

**Root dataset:** For the server, a root dataset is required to obtain a trusted gradient. By default, we assume that the root dataset and the client's local training data have the same distribution, specifically the samples of the root dataset are random and

2) Here we use Diffie–Hellman and an RSA-based signature algorithm to jointly negotiate a secret seed.

uniform from the local clean training data of all clients. For example, we randomly and uniformly sample the root dataset from its 60,000 training samples when using MNIST.

### Appendix C.1.2 Experiment Parameters

We assume that a total of  $N$  clients participate in federated learning and the number of communication rounds is  $T$ , where  $\alpha$  is the global model learning rate. And the proportion of clients participating in each round of communication is set to  $\sigma$ , which includes malicious clients and non-malicious clients. Among them, corrupt clients account for a proportion of  $\delta$  of all clients in this round. For local model training, the local model is updated using momentum SGD. In addition, the local batch size is  $B$ , the number of local updates is  $E$ , the local learning rate is  $\beta$  and the momentum parameter is  $\gamma$ . The test parameter settings we use by default are shown in Table C2.

**Table C2** System parameter settings.

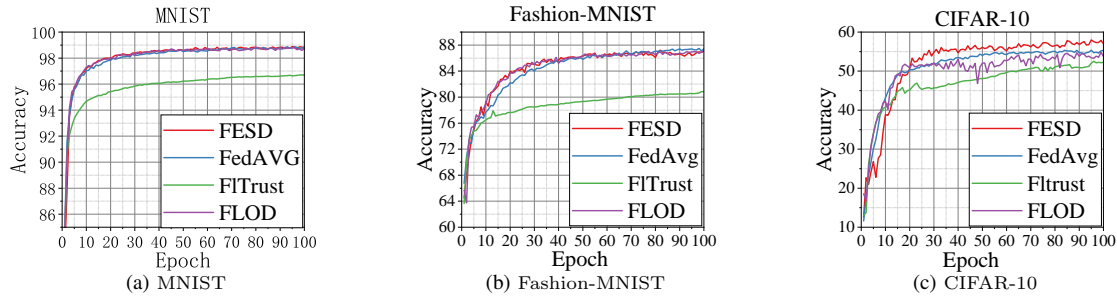
Notation	Parameter	Notation	Parameter
$N$	100	$\sigma$	0.1
$E$	5	$T$	100
$B$	10	$\alpha$	0.25
$\beta$	0.01	$\delta$	0.1

## Appendix C.2 Experimental Results

In this comparison experiment, we will use four baseline comparison solutions, namely FedAvg, FLTrust, FLOD and FLDetector. Among them, FedAvg is the basic algorithm of federated learning. The latter three solutions detect malicious clients by calculating cosine similarity, Hamming distance, and Euclidean distance respectively, and support dishonest-majority settings. To test the effect of this scheme on these three design goals (Section 4.4), we carry out the following experiments.

### Appendix C.2.1 Fidelity

A Byzantine robust aggregation method should also be suitable for attack-free situations in practical applications. When  $\delta = 0$ , it should have almost no loss in model performance. Therefore, we first conduct model training on the FedAvg scheme without malicious model attacks. When FedAvg achieves the best performance, we compare FESD with it to show how good or bad this scheme is in terms of fidelity. In addition, we also measure the accuracies of FITrust and FLOD. In the test, we successively use three different datasets, as shown in Figure C1(a) and Figure C1(b), when we use the MNIST and Fashion-MNIST datasets for training, the accuracy of FESD is almost the same as the FedAvg and FLOD schemes and higher than FITrust. And as shown in Figure C1(c), the accuracy of FESD is even higher than those of other schemes. This is because the root model update in FESD also participates in model aggregation, which is more conducive to model training to a certain extent and the effect is more prominent on the CIFAR-10 dataset, which is closer to the real dataset and difficult to identify. Finally, we present the highest accuracy and loss measured by each scheme in Table C3. Based on the above, the effect of this scheme in terms of fidelity is feasible.



**Figure C1** Model performance testing without attacks

### Appendix C.2.2 Robustness

In order to fully evaluate the Byzantine robust performance of this scheme, we continuously change the value of  $\delta$  of corrupt clients from 0.1 to 0.8 and compare the model accuracy measured in different datasets when compared with the FITrust and FLDetector. It should be noted that the value of the global learning rate  $\alpha$  is also different in the case of different proportions of corrupt clients.

We show the model accuracy for different proportions of corrupt clients in Figure C2. First of all, we find that FESD shows good stability under the MNIST, Fashion-MNIST and Cifar-10 datasets. When  $\delta = 0.1$ , the model accuracy on the MNIST dataset is 98.75%, when  $\delta = 0.8$ , the accuracy of the global model remains above 98.37% and only decreases by 0.38%. Similarly, in the Fashion-MNITS dataset, the model accuracy does not vary by more than 2%. In the Cifar-10 dataset, the model accuracy does not vary by more than 5%. In addition, under different proportions of corrupt clients, the model accuracy of FESD is generally equal to or higher than that of FITrust, FLOD and FLDetector, which shows that it is effective for us to participate in the global model aggregation process of root model update, which can improve the model accuracy to a certain extent. In addition, when  $\delta > 0.5$ ,

**Table C3** Accuracy and loss without malicious attack.

$\delta = 0$	dataset	FESD	FedAvg	FITrust	FLOD
Accuracy	MNIST	98.889	98.879	96.71	98.779
	Fashion-MNIST	87.220	87.489	80.86	87.040
	CIFAR-10	57.990	55.290	52.951	55.819
Loss	MNIST	0.033	0.037	0.103	0.038
	Fashion-MNIST	0.376	0.345	0.507	0.385
	CIFAR-10	1.793	1.828	1.942	1.839

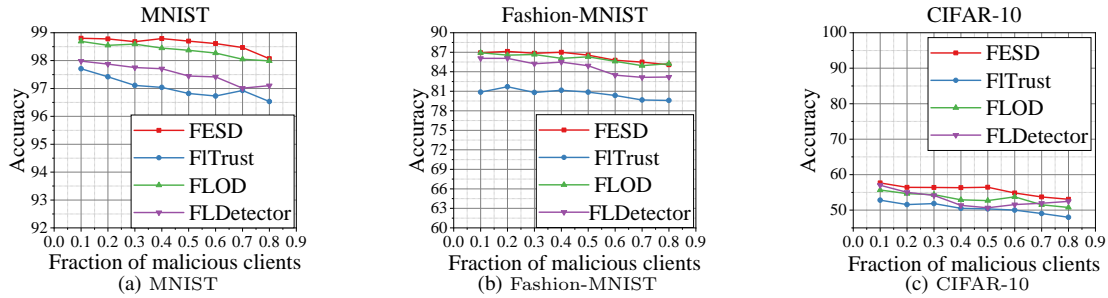


Figure C2 Model performance test with different fraction of corrupt clients

Table C4 Client overhead with communication rounds

Dataset		FESD			FdeAvg		
		MNIST	Fashion	CIFAR-10	MNIST	Fashion	CIFAR-10
Rounds	10	169.942s	166.817s	146.022s	163.118s	165.895s	144.157s
	50	849.716s	835.687s	731.673s	816.732s	829.783s	721.931s
	100	1699.672s	1672.458s	1463.796s	1635.986s	1661.342s	1445.426s

the trend of model accuracy decline is slower than other schemes. This shows that it is effective for us to use the proportion of detected corrupt clients to adaptively scale the global learning rate and to some extent make up for the problem of missing effective information caused by malicious clients. In addition, when a malicious client is detected by FLDetector, the malicious client will be directly removed from the model aggregation. Although it effectively prevents the influence of malicious data on the model, it also causes the data information of some clients to be missing, which leads to the convergence speed of FLDetector being much lower than that of FESD, FITrust and FLOD. In Figure C2, FESD, FITrust and FLOD basically converge in 100 rounds, while FLDetector needs more than 300 rounds to complete convergence. Based on the above, the robustness of the scheme is feasible.

Table C5 The overhead of several schemes on PCBit2A

Operation		Bit2A					
Model		FNet			ResNET-18		
Phase	Rounds	FESD	FLOD	FLGUARD	FESD	FLOD	FLGUARD
Online	10	-	1.43s	2.90s	-	7.44s	15.28s
	50	-	7.44	14.77s	-	33.86s	91.74s
	100	-	14.07s	30.60s	-	67.93s	153.13s
Offline	10	-	58.16s	151.22s	-	237.17s	616.64s
	50	-	290.06s	754.16s	-	1185.06s	3081.16s
	100	-	580.16s	1508.42s	-	2370.45s	6163.17s

### Appendix C.2.3 Efficiency

Since the client in federated learning is a resource-constrained device, it is necessary to minimize its load. To measure user efficiency, we use client runtime as a metric. We give the total time for all clients to run locally under different rounds in Table C4, from which we can calculate the average time of each client. After calculation, it can be concluded that the average duration of each client in this solution or FedAvg does not exceed 0.03s. On the premise of ensuring data privacy and security, it is desirable not to increase the user load too much. In addition to highlighting the advantages of this scheme, we also measure the extra overhead used by the FLOD scheme at PCBit2A. It can be seen from Table C5 that this scheme adopts the sampling secret sharing method, which can greatly reduce the computational cost and improve efficiency.

### References

- O. Goldreich, S. Micali, and A. Wigderson, "How to play any mental game," in *Proc. ACM Symposium on the Theory of Computing (STOC'87)*, 1987, pp. 218–229.
- A. C.-C. Yao, "How to generate and exchange secrets," in *Proc. IEEE Annual Symposium on Foundations of Computer Science (FOCS'86)*, 1986, pp. 162–167.
- M. Keller, E. Orsini, and P. Scholl, "Mascot: faster malicious arithmetic secure computation with oblivious transfer," in *Proc. ACM Conference on Computer and Communications Security (CCS'16)*, 2016, pp. 830–842.
- J. Cha, S. K. Singh, T. W. Kim, and J. H. Park, "Blockchain-empowered cloud architecture based on secret sharing for smart city," *Journal of Information Security and Applications*, 2021, doi:https://doi.org/10.1016/j.jisa.2020.102686.
- Y. Dong, X. Chen, K. Li, D. Wang, and S. Zeng, "Flod: Oblivious defender for private byzantine-robust federated learning with dishonest-majority," in *European Symposium on Research in Computer Security (ESORICS'21)*, 2021, pp. 497–518.
- R. Canetti, "Universally composable security: A new paradigm for cryptographic protocols," in *Proc. IEEE Annual Symposium on Foundations of Computer Science (FOCS'01)*, 2001, pp. 136–145.
- L. Deng, "The mnist database of handwritten digit images for machine learning research [best of the web]," *IEEE signal processing magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint*, 2017, arXiv:1708.07747.
- A. Krizhevsky, G. Hinton et al., "Learning multiple layers of features from tiny images," 2009.