• Supplementary File •

FUSE: A \underline{F} ederated Learning and \underline{U} -shape \underline{S} plitLearning-based \underline{E} lectricity Theft DetectionFramework

Xuan LI¹, Naiyu WANG¹, Liehuang ZHU², Shuai YUAN³ & Zhitao GUAN^{1*}

¹School of Control and Computer Engineering, North China Electric Power University, Beijing 102206, China;
 ²School af Cyberspace Science and Technology, Beijing Institute of Technology, Beijing 100081, China;
 ³Department of Finance, Operations, and Information Systems (FOIS), Brock University, St.Catharines L2S 3A1, Canada

Appendix A Preliminary

Appendix A.1 Federated Learning

Federated learning is a method of distributed machine learning where a global model is trained collaboratively across multiple clients, coordinated by a central server [1]. A vanilla federated learning typically consists of four steps: Initialization, Local Training, Model Aggregation, and Model Update.

• Initialization: During this phase, participating clients (e.g., user devices) reach a consensus on training a high-performing model using data distributed across the clients. An initial global model is created on the server, which serves as the beginning of the federated learning cycle. Each client then downloads the initial global model as their local model for the first round.

• Local Training: Each client trains their local model using their own data following the predefined hyperparameters and training procedures. Then after a set number of iterations, either the local model parameters or the model updates are sent to the server by the clients.

• Model Aggregation: Upon receiving all local model parameters or updates, the server aggregates them to generate an updated global model for the next round of global communication.

• Model Update: All clients download the newly aggregated global model to update their local models. The cycle of Local Training, Model Aggregation, and Model Update repeats iteratively until the global model achieves convergence or meets the predefined number of communication rounds.

Based upon the aggregation strategy on the server, federated learning can be categorized into synchronous, asynchronous, and semi-asynchronous federated learning [2]. In synchronous federated learning, the server waits to aggregate the global model until it has received all local models from the clients, as discussed above [1]. To improve efficiency, asynchronous federated learning enables the clients to submit their local models asynchronously, with the server updating the global model upon each receipt [3]. However, the high communication overhead in the asynchronous aggregation strategy and the extended waiting periods in the synchronous aggregation strategy have led to the development of semi-asynchronous federated learning [4–6]. This approach stores early-arriving local models and aggregates them within a specific timeframe. Furthermore, local models collected within this timeframe may come from different communication rounds and are aggregated with variable weights. Thus it effectively combines the advantages of both synchronous aggregation strategies.

Appendix A.2 Split Learning

Split learning is designed to minimize the substantial computational overhead while preserving the privacy of the raw data [7–9]. In this context, two distinct roles are involved: the client, typically with limited computational capability, and the server, introduced to alleviate the computational burden of the client. Specifically, in a basic split learning configuration, a neural network with a large number of parameters is partitioned into two sections. The lower layers, closer to input, are assigned to the client, while the remaining layers are managed by the server. This division occurs at a designated "cut layer". The client trains the lower layers using its data and submits the output from the cut layer to the server. The server then completes the forward propagation process without needing to access the raw data of the client. Backpropagation begins at the server, with the gradients of the cut layer transmitted back to the client for subsequent steps. This process repeats until the model achieves the desired performance level. Through split learning, the bulk of the model training computation is performed by the server with sufficient computational resources. Meanwhile, the client retains the raw data locally, ensuring its privacy throughout the entire training process.

Appendix B Problem Formulation & Threat Model

Appendix B.1 Problem Formulation

Consider M transformer districts $\text{TD} = \{\text{TD}_i | i \in [1, M]\}$ in different energy service providers (ESP) are involved in the federated learning process to train a theft detection model jointly. For each TD, there are N smart meters SM_i^j , where j is from 1 to N. SM_i^j maintains the daily consumption data x_i^j of a given consumer. y_i^j represents the consumer label of x_i^j , denoting whether or not the theft behavior exists on SM_i^j . Thus, it is assigned and managed by TD_i. The local dataset of TD_i, therefore, can be represented

^{*} Corresponding author (email: guan@ncepu.edu.cn)

as $D_i = \left\{ \begin{pmatrix} x_i^i, y_i^j \end{pmatrix} | i \in [1, M], j \in [1, N] \right\}$ with the size of $|D_i|$. The theft detection model is then defined as $\hat{y} = f(w, x)$, where \hat{y} denotes the prediction of the model, w denotes the global model parameter, and x denotes the model input, i.e. electricity consumption data. Accordingly, the loss function of TD_i is represented as $\mathcal{L}_i = \frac{1}{|D_i|} \sum_{j \in [1,N]} \ell\left(\hat{y}_i^j, y_i^j\right)$, where $\ell\left(\hat{y}_i^j, y_i^j\right)$ is the loss function for data point $\begin{pmatrix} x_i^j, y_i^j \end{pmatrix}$ and $\ell(\cdot) \in (0, +\infty)$. As a result, the global loss function is defined as:

$$F(w) = \sum_{i=1}^{M} \frac{|D_i|}{S} \mathcal{L}_i$$
(B1)

where $S = \sum |D_i|$.

Therefore, the final objective of the framework is to solve the optimization problem as follows:

$$w^* = \arg\min F\left(w\right). \tag{B2}$$

where w^* denotes the optima of the global parameter.

Meanwhile, the straggler issue arises due to variations in computation and communication resources among TDs from different ESPs. Consequently, TDs send their local models asynchronously, resulting in the global model being updated based on the aggregation of local models trained in different rounds. Note that w_i^k denotes the local model trained by TD_i based on the global model in round k, w^r denotes the global model in round r, and \mathcal{A} denotes an aggregation algorithm. Thus the update of the global model in round r can be expressed as $w^r = w^{r-1} + \mathcal{A}(w_i^k | i \in [1, M], k \in [0, r-1])$. We also define the staleness of w_i^k as $\tau = r - k$.

Appendix B.2 Threat Model and Objectives

The threat model considered in this paper addresses both internal and external issues as follows:

Internal Threats: We assume that the assisted cloud server and federated learning aggregator are honest-but-curious, i.e., they will perform their tasks honestly but may attempt to infer sensitive information from the raw data through model updates.

External Threats: The adversary may attempt to obtain sensitive data, such as the consumer consumption data or information regarding electricity theft, either directly or indirectly by eavesdropping on the communication link between the entities in our framework.

Therefore, our framework is designed with the following objectives in mind:

(1) Lightweight model training on TDs and SMs.

(2) Accommodation of asynchronous model submissions from heterogeneous clients without compromising model performance.

(3) Preservation of training data, including the consumption data as input in SMs and the consumer labels in TDs, at the

location where it is generated to the greatest extent possible.

(4) Prevention of the assisted cloud, aggregator, and eavesdropper from inferring sensitive information throughout the entire process.

Appendix C The Proposed Framework

Appendix C.1 Entities in FUSE

Four entities participate in the FUSE framework:

• Smart Meter (SM): The SM collects and stores the electricity consumption data, i.e., the input of the model. Due to the constrained computational capacity and communication bandwidth, only the Feature Extractor is allocated to a given SM and the forward propagation of Feature Extractor is executed by SM. In each local iteration, the output of Feature Extractor is sent to the associated TD.

• Transformer District (TD): The TD stores the label on whether or not the theft behavior exists on a given SM in the region. Thus, the Feature Classifier is allocated to each TD in the training process. Meanwhile, to reduce the computational requirements of SM, Feature Extractor is also assigned to each TD such that backpropagation and parameter updates can also be performed.

• Cloud Server (CS): The CS is assumed to have sufficient computational capacity and communication bandwidth. Note that there is a 1:1 mapping of each TD to a CS logically, where the main component of the detection framework, i.e., Feature Learner, is allocated to perform local model training. As a result, a local theft detection model in federated learning consists of the Feature Extractor in the TD (or SM), the Feature Classifier in the TD, and the Feature Learner in the TD's corresponding CS.

• Aggregator: To address additional privacy concerns, the ciphertext of the serialized model parameters is submitted such that the global model is securely aggregated by the Aggregator. This aggregation process can be performed by the parent company of ESPs.

Notation	Definition			
TD_i	The <i>i</i> th Transformer District			
SM_i^j	The j th Smart Meter in the i th Transformer District			
CS_i	The Cloud Server which completes the local model training for TD_i			
x_i^j	The electric consumption data collected by SM_i^j			
y_i^j	The consumer label of x_i^j which is assigned and managed by $ ext{TD}_i$			
D_i	The dataset on TD_i			
$w_{\rm FE},\!w_{\rm FL},\!w_{\rm FC}$	The parameter of the Feature Extractor, the Feature Learner, and the Feature Classifier			
$v^k_{\mathrm{TD}_i}$, $v^k_{\mathrm{CS}_i}$	The serialized local model parameter of TD_i and CS_i in round k			

Appendix C.2 Table C1: Key Notations

Table C1 Key Notations

$ \begin{array}{ll} pk_{\mathrm{TD}} \ , sk_{\mathrm{TD}} & \mbox{The public key and the private key of all TDs} \\ pk_{\mathrm{CS}} \ , sk_{\mathrm{CS}} & \mbox{The public key and the private key of all CSs} \\ \hline pk_{\mathrm{TD}} & \mbox{The buffer of serialized Feature Extractors and Feature Classifiers} \\ \hline puff_{\mathrm{TD}} & \mbox{The buffer of serialized Feature Learners} \\ \hline T & \mbox{The maximum waiting time for a communication round} \\ \hline E & \mbox{The number of local iterations} \end{array} $	v^k_{TD} , v^k_{CS}	The serialized global model parameter in round k
$ \begin{array}{ll} pk_{\rm CS} \ , sk_{\rm CS} & \mbox{The public key and the private key of all CSs} \\ Buff_{\rm TD} & \mbox{The buffer of serialized Feature Extractors and Feature Classifiers} \\ Buff_{\rm CS} & \mbox{The buffer of serialized Feature Learners} \\ T & \mbox{The maximum waiting time for a communication round} \\ \hline E & \mbox{The number of local iterations} \end{array} $	pk_{TD} , sk_{TD}	The public key and the private key of all TDs
$Buff_{TD}$ The buffer of serialized Feature Extractors and Feature Classifiers $Buff_{CS}$ The buffer of serialized Feature Learners T The maximum waiting time for a communication round E The number of local iterations	${pk_{\rm CS}}$, ${sk_{\rm CS}}$	The public key and the private key of all CSs
$Buff_{CS}$ The buffer of serialized Feature Learners T The maximum waiting time for a communication round E The number of local iterations	$Buff_{\rm TD}$	The buffer of serialized Feature Extractors and Feature Classifiers
T The maximum waiting time for a communication round E The number of local iterations	$Buff_{\rm CS}$	The buffer of serialized Feature Learners
<i>E</i> The number of local iterations	T	The maximum waiting time for a communication round
	E	The number of local iterations

Appendix C.3 Algorithm C1: FUSE

Algorithm C1 FUSE
Input: $TD_i \in TD$, $CS_i \in CS$, global communication round R , max waiting time T
Output: $\llbracket v_{\text{TD}} \rrbracket_{pk_{\text{TD}}}, \llbracket v_{\text{CS}} \rrbracket_{pk_{\text{CS}}}$
▷ initialize the global model
1: Select TD_i from $TD = \{TD_1, TD_2,\}$ at random.
2: TD_i generates $(pk_{\text{TD}}, sk_{\text{TD}})$ and broadcasts to $\text{TD}_i \in \text{TD}$.
3: TD _i initializes w_{FE}^0 , w_{FC}^0 and sends $[v_{\text{TD}}^0]_{pk_{\text{TD}}} = encrypt(Serialized(w_{\text{FE}}, w_{\text{FC}}), pk_{\text{TD}})$ to the Aggregator.
4: Select CS_i from $CS = \{CS_1, CS_2,\}$ at random.
5: CS_i generates (pk_{CS}, sk_{CS}) and broadcasts to $CS_i \in CS$.
6: CS_i initializes w_{FL}^0 and sends $[v_{CS}^0]_{pk_{CS}} = encrypt(Serialized(w_{FL}), pk_{CS})$ to the Aggregator.
⊳ global training process
7: for round $r = 1, 2, 3,, R$ do
8: The Aggregator broadcasts $[\![v_{\text{TD}}^{r-1}]\!]_{pk_{\text{TD}}}$ to $\text{TD}_i \in \text{TD}, [\![v_{\text{CS}}^{r-1}]\!]_{pk_{\text{CS}}}$ to $\text{CS}_i \in \text{CS}$.
▷ local training process
9: for each (TD_i, CS_i) in parallel do
10: $LocalSplitTraining([v_{TD}^{r-1}]_{pk_{TD}}, [v_{CS}^{r-1}]_{pk_{CS}})$
11: end for
\triangleright initialize the local model buffer
12: $Buff_{TD} = \emptyset, Buff_{CS} = \emptyset, j = 0$
13: while not reaching the max waiting time T do
14: The Aggregator retrieves $[\![v_{\text{TD}_i}^k]\!]_{pk_{\text{TD}}}$ and $[\![v_{\text{CS}_i}^{k'}]\!]_{pk_{\text{CS}}}$
15: $\cos_{\mathrm{TD}}^{j} = \frac{\left[v_{\mathrm{TD}_{i}}^{k}\right]_{pk_{\mathrm{TD}}} \cdot \left[v_{\mathrm{TD}_{i}}^{r-1}\right]_{pk_{\mathrm{TD}}}}{\left\ \left[v_{\mathrm{TD}_{i}}^{k}\right]_{pk_{\mathrm{TD}}}\right\ \left\ \left[v_{\mathrm{TD}_{i}}^{r-1}\right]_{pk_{\mathrm{TD}}}\right\ } \right } > \text{ calculate the cosine similarity of the model from TD}$
16: $\cos_{\mathrm{CS}}^{j} = \frac{\left[v_{\mathrm{CS}_{i}}^{k'} \right]_{pk_{\mathrm{CS}}} \cdots \left[v_{\mathrm{CS}}^{r-1} \right]_{pk_{\mathrm{CS}}}}{\left\ v_{\mathrm{CS}_{i}}^{k'} \right\ _{pk_{\mathrm{CS}}} \left\ v_{\mathrm{CS}_{i}}^{r-1} \right\ _{pk_{\mathrm{CS}}}} > \text{calculate the cosine similarity of the model from CS}$
17: $Buff_{TD} = Buff_{TD} \cup (\llbracket v_{TD_i}^k \rrbracket_{pk_{TD}}, cos_{TD}^j, i, k, j) \triangleright \text{ add the model from TD to the buffer}$
18: $Buff_{CS} = Buff_{CS} \cup ([v_{CS_i}^{k'}]]_{pk_{CS}}, cos_{CS}^j, i, k', j) \triangleright$ add the model from CS to the buffer
19: $j = j + 1$
20: end while
\triangleright perform aggregation process to obtain the new global model
21: $ [\![v_{\text{CS}}^r]\!]_{pk}_{\text{CS}} = \textit{Two-StageAggregation}(Buff_{\text{CS}}, [\![v_{\text{CS}}^{r-1}]\!]_{pk}_{\text{CS}}, r) $
22: $[\![v_{\text{TD}}^r]\!]_{pk_{\text{TD}}} = \textit{Two-StageAggregation}(Buff_{\text{TD}}, [\![v_{\text{TD}}^{r-1}]\!]_{pk_{\text{TD}}}, r)$
23: end for
24: $\llbracket \boldsymbol{v}_{\mathrm{TD}} \rrbracket_{pk_{\mathrm{TD}}} = \llbracket \boldsymbol{v}_{\mathrm{TD}}^r \rrbracket_{pk_{\mathrm{TD}}}, \llbracket \boldsymbol{v}_{\mathrm{CS}} \rrbracket_{pk_{\mathrm{CS}}} = \llbracket \boldsymbol{v}_{\mathrm{CS}}^r \rrbracket_{pk_{\mathrm{CS}}}$
25: return $\llbracket v_{\mathrm{TD}} \rrbracket_{pk_{\mathrm{TD}}}, \llbracket v_{\mathrm{CS}} \rrbracket_{pk_{\mathrm{CS}}}$

Appendix C.4 Initialization

At the commencement of the training, an arbitrary $\text{TD}_i \in \text{TD}$ is selected to initialize the global parameter of the Feature Extractor w_{FE}^0 , the global parameter of the Feature Classifier w_{FC}^0 , and the public-private key pair $(pk_{\text{TD}}, sk_{\text{TD}})$ for homomorphic encryption. Then, $(pk_{\text{TD}}, sk_{\text{TD}})$ is broadcasted to all $\text{TD}_i \in \text{TD}$ through secure channels. Meanwhile, an arbitrary $\text{CS}_i \in \text{CS}$ is also selected to initialize the global parameter of the Feature Learner w_{FL}^0 and the public-private key pair $(pk_{\text{CS}}, sk_{\text{CS}})$. Then $(pk_{\text{CS}}, sk_{\text{CS}})$ is broadcasted to all $\mathrm{CS}_i \in \mathrm{CS}$ through secure channels. To ensure that the model parameters and structure are not disclosed to the Aggregator, each local update is submitted after the encryption by the public key of TD or CS and serialized. The serialization operation refers to the transformation from the origin parameter $w \in \mathbb{R}^{M \times N}$ to $v \in \mathbb{R}^{1 \times MN}$, such that the structure of the model remains unknown. Furthermore, at the beginning of each communication round r, the serialized and encrypted model parameters $[v_{\mathrm{TD}}^{r-1}]_{pk_{\mathrm{TD}}}$ and $[v_{\mathrm{CS}}^{r-1}]_{pk_{\mathrm{CS}}}$ are broadcasted to all TD_i and CS_i respectively.

Appendix C.5 Three-tier U-shape Split Learning-based Local Training

Each local model is spilt and trained by different entities in our framework following Algorithm C2. In this section, we elaborate on the three-tier U-shape split learning-based local model training from the perspective of a given TD_i and its corresponding CS_i as follows.

Algorithm C2 LocalSplitTraining

Input: $[\![v_{\text{TD}}]\!]_{pk_{\text{TD}}}, [\![v_{\text{CS}}]\!]_{pk_{\text{CS}}}$ **Output:** $[v_{TD_i}]_{pk_{TD}}, [v_{CS_i}]_{pk_{CS}}$ \triangleright initialize the local model 1: TD_i : $v_{\text{TD}} = decrypt(\llbracket v_{\text{TD}} \rrbracket_{pk_{\text{TD}}}, sk_{\text{TD}})$, $w_{\text{FE}}, w_{\text{FC}} = Deservative(v_{\text{TD}})$ $2: \ \mathrm{CS}_i: \ v_{\mathrm{CS}} = decrypt([\![v_{\mathrm{CS}}]\!]_{pk}_{\mathrm{CS}}, sk_{\mathrm{CS}}), \ w_{\mathrm{FL}} = Deserialize(v_{\mathrm{CS}})$ 3: $w_{\text{FE}}^0 \leftarrow w_{\text{FE}}, w_{\text{FL}}^0 \leftarrow w_{\text{FL}} \ w_{\text{FC}}^0 \leftarrow w_{\text{FC}}$ ▷ three-tier U-shape split learning process 4: for $e \in [1, 2, ..., E]$ do TD_i broadcasts w_{FE}^{e-1} to $\mathrm{SM}_i^j \in \mathrm{SM}_i$ 5:▷ forward propagation process for $\mathrm{SM}_i^j \in \mathrm{SM}_i$ in parallel do 6: $output_{\text{FE}}^{e-1} \leftarrow ForwardPropagation(w_{\text{FE}}^{e-1}, x_i^j)$ 7: send $output_{\text{FE}_i}^{j}$ to TD_i 8: 9: end for
$$\begin{split} & \text{TD}_i \colon \textit{output}_{\text{FE}}^{e-1} = \textit{Concatenate}(\textit{output}_{\text{FE}j}^{e-1}), j \in [1, \dots, n], \, \text{send } \textit{output}_{\text{FE}}^{e-1} \text{ to } \text{CS}_i. \\ & \text{CS}_i \colon \textit{output}_{\text{FL}}^{e-1} \leftarrow \textit{ForwardPropagation}(w_{\text{FL}}^{e-1}, \textit{output}_{\text{FE}}^{e-1}) \end{split}$$
10: 11: CS_i : send $output_{FL}^{e-1}$ to TD_i 12: $\text{TD}_i: \hat{y} \leftarrow \textit{ForwardPropagation}(w_{\text{FC}}^{e-1}, output_{\text{FL}}^{e-1}), L = \ell(\hat{y}, y) \triangleright \text{ calculate the loss of the data}$ 13:▷ backpropagation and local model update process $\text{TD}_i: g_{\text{FC}} \leftarrow BackPropagation(\bigtriangledown L, w_{\text{FC}}^{e-1}) \triangleright$ calculate the gradient of the FC's parameters 14: $\text{TD}_i: \text{ update Feature Classifier by } w_{\text{FC}}^e \leftarrow w_{\text{FC}}^{e-1} - \eta g_{\text{FC}} \text{ , send } g_{\text{FC}}^*, \text{ the gradient of the first layer of } w_{\text{FC}}^e, \text{ to } \text{CS}_i \text{ and } y_{\text{FC}}^* \text{ and } y_{\text{F$ 15: $\mathrm{CS}_i \colon g_{\mathrm{FL}} \leftarrow BackPropagation(g_{\mathrm{FC}}^*, w_{\mathrm{FL}}^{e-1}) \triangleright \text{ calculate the gradient of the FL's parameters}$ 16: CS_i : update Feature Learner by $w_{\text{FL}}^e \leftarrow w_{\text{FL}}^{e-1} - \eta g_{\text{FL}}$, send g_{FL}^* , the gradient of the first layer of w_{FL}^e , to TD_i 17: $\mathrm{TD}_i : \ g_{\mathrm{FE}} \leftarrow BackPropagation(g^*_{\mathrm{FL}}, w^{e-1}_{\mathrm{FE}}) \triangleright \text{ calculate the gradient of the FE's parameters}$ 18:TD_i: update Feature Extractor by $w_{\text{FE}}^{e} \leftarrow w_{\text{FE}}^{e-1} - \eta g_{\text{FE}}$ 19:20: end for 21: TD_i : send $[\![v_{\text{TD}_i}]\!]_{pk_{\text{TD}}} = encrypt(Serialize(w_{\text{FE}}, w_{\text{FC}}), pk_{\text{TD}})$ to the Aggregator 22: CS_i: send $[v_{CS_i}]_{pk_{CS}} = encrypt(Serialize(w_{FL}), pk_{CS})$ to the Aggregator

$$w^e \leftarrow w^{e-1} - \eta g \tag{C1}$$

where w^e denotes the updated parameter in local round e, w^{e-1} denotes the parameter in last round, g denotes the gradient of weights, and η denotes the learning rate. It is worth noting that the backward propagation and the update of Feature Extractor are performed by TD_i, rather than SMⁱ_j, such that the risk of incurring high computational overhead on SMs is mitigated.

After a pre-determined number of local iteration rounds, the serialized split model in each entity is encrypted and sent to the Aggregator. Since parameters are updated ahead, the model can be submitted by CS_i without the necessity of waiting for TD_i such that the asynchronous problem between TD_i and its corresponding CS_i is addressed.

Once $[v_{\text{TD}}]_{pk_{\text{TD}}}$ is received, TD_i decrypts it by the private key sk_{TD} . And $v_{\text{TD}} \in \mathbb{R}^{1 \times MN}$ is describing to obtain parameters w_{FE} and w_{FC} . Similarly, $[v_{\text{CS}}]_{pk_{\text{CS}}}$ is decrypted by sk_{CS} in CS_i and w_{FL} is obtained through describing describing the vector.

In a given local iteration, TD_i first distributes w_{FE} to all SM_i^j in its region. Then for each SM_i^j , the consumption data is fed into the Feature Extractor where the forward propagation is carried out. Accordingly, the outputs of all Feature Extractors from the last layer are collected by TD_i and concatenated as the input to the Feature Learner in CS_i such that the forward propagation can be performed. The output of the Feature Learner is then fed into the Feature Classifier in TD_i where the predicted label is derived. In addition, the backpropagation is performed by feeding the gradient of the first layer of the partitioned model into the prior one. For instance, the gradient of the first layer of the Feature Classifier is sent from TD_i to CS_i . Hence, the parameters are updated once the gradients are calculated as follows:



Appendix C.6 Two-stage Semi-asynchronous Aggregation



The updated model from different TDs and CSs may be received by the Aggregator asynchronously, as demonstrated in Figure C1. In classical federated learning models, either longer training time occurs in waiting for the stragglers, or the model performance is reduced because of the removals of such stragglers. On the other hand, while computation time can be reduced by asynchronous aggregation during the training process, extensively higher communication overhead is incurred on TD since each submission of a new local model necessitates the update of the global model. Therefore, a two-stage semi-asynchronous aggregation mechanism is proposed in our work to address these issues. Meanwhile, to guarantee that both the structure of the model and the plaintext of the parameters are not disclosed to the Aggregator, homomorphic encryption is implemented. Furthermore, considering the computation involved in the subsequent aggregation process, fully homomorphic encryption is employed. Hence the serialized and encrypted model parameters $[v^{r-1}]_{pk}$ will be aggregated in each communication round. The two-stage semi-asynchronous aggregation is shown in Algorithm C3.

Algorithm C3 Two-StageAggregation

Input: Buff, $[v^{r-1}]_{pk}$, r

Output: $\llbracket v^r \rrbracket_{pk}$

- 1: Sort Buff in descending order of cos^j
- $2: \textit{Buff} \gets \textit{Buff}[0:m] \triangleright \text{ only top-m models are retained in } \textit{Buff}$

 \triangleright pre-aggregate the models in the buffer

3:
$$\llbracket v^{new} \rrbracket_{pk} = \sum_{i,k,j \in Buff} \frac{\cos^j}{\sum_{j \in Buff} \cos^j} \llbracket v_i^k \rrbracket_{pk}$$

4:
$$|\overline{D}| = \sum_{i,j \in Buff} \frac{\cos^j}{\sum_{j \in Buff} \cos^j} |D_i|$$

5:
$$\overline{\tau} = \sum_{k,j \in Buff} \frac{\cos j}{\sum_{j \in Buff} \cos j} (r-k)$$

 \triangleright aggregate the local models to obtain the new global

6: $[\![v^r]\!]_{pk} \leftarrow (1 - \theta_r) [\![v^{r-1}]\!]_{pk} + \theta_r [\![v^{new}]\!]_{pk}$, where $\theta_r = \theta g(\overline{\tau}, |\overline{D}|), g(\cdot)$ is a function of staleness and dataset size, respectively. $g(\overline{\tau}, |\overline{D}|) = \frac{1}{(1 + \overline{\tau} + \sum |D_i|) \frac{1}{c}}$

model

7: return
$$\llbracket v^r \rrbracket_{pk}$$

Once a local update is received, the model parameters are stored in the corresponding buffer, i.e., the models from TDs are stored in $Buff_{\rm TD}$ while those from CSs are stored in $Buff_{\rm CS}$. Taking the varying communication capabilities of different devices into account, models in the buffers with lower quality or higher variance from the global optima have a direct impact on the model performance. The cosine similarity between current local model and previous global model can effectively indicate the discrepancy in the convergence direction [10]. This similarity measure serves as a valuable tool to identify and filter out malicious models. Building upon this conclusion, we employ cosine similarity to evaluate the quality of the model, which provides the basis for obtaining the aggregation weights. Thus, before the aggregation with the global model from prior round, we pre-aggregate the

models in each buffer in terms of cosine similarity as follows:

$$cos = \frac{\llbracket v_k^i \rrbracket_{pk} \cdot \llbracket v^{r-1} \rrbracket_{pk}}{\lVert \llbracket v_k^i \rrbracket_{pk} \rVert \, \lVert \llbracket v^{r-1} \rrbracket_{pk} \rVert} \tag{C2}$$

such that the quality can be measured, where k represents the round from which the basic global model derived, r represents the current round, and the staleness of the model can be derived as $\tau = r - k$. To further improve the efficiency, only top-m models with the highest cos(the model with the cos < 0 will be dropped directly) will be pre-aggregated, thereby filtering out the models with low quality. This pre-aggregation is comprised of deriving the weighted average parameters, dataset size and the staleness regarding, as follows:

$$\llbracket v^{new} \rrbracket_{pk} = \sum_{i,k,j \in Buff} \frac{\cos^j}{\sum_{j \in Buff} \cos^j} \llbracket v_i^k \rrbracket_{pk}$$
(C3)

$$|\overline{D}| = \sum_{i,j \in Buff} \frac{\cos^j}{\sum_{j \in Buff} \cos^j} |D_i|$$
(C4)

$$\overline{\tau} = \sum_{k,j \in Buff} \frac{\cos^j}{\sum_{j \in Buff} \cos^j} (r - k) \tag{C5}$$

For the update of the global model, we define a hyper-parameter $\theta \in [0, 1]$. Therefore the final weight of the updated model θ_r is determined by the product of θ and $g(\cdot)$, which is a function of staleness and dataset size, i.e., $\theta_r = \theta g(\overline{\tau}, |\overline{D}|)$. $g(\cdot)$ is defined as follows where a smaller weight will be assigned to the model with larger staleness and smaller dataset size:

$$g(\overline{\tau}, |\overline{D}|) = \frac{1}{(1 + \overline{\tau} + \frac{\sum |D_i|}{M|\overline{D}|})^{\frac{1}{e}}}$$
(C6)

Hence, the global model is updated as:

$$\llbracket v^r \rrbracket_{pk} \leftarrow (1 - \theta_r) \llbracket v^{r-1} \rrbracket_{pk} + \theta_r \llbracket v^{new} \rrbracket_{pk}$$
(C7)

Note that all the algorithms including cosine similarity calculations and model aggregations, are performed in ciphertext because of the fully homomorphic encryption applied. In addition, since the comparisons in ciphertext can be approximated by addition and multiplication according to [11], the detailed operations on the ciphertext will not be disclosed in our framework.

Appendix D Experimental Evaluation

Appendix D.1 Experimental Setups

Implementation. All the experiments are carried out on Windows 10 operating system equipped with AMD Ryzen 3600 CPU, 16 GB RAM, and NVIDIA GeForce GTX 1660 SUPER. The federated learning algorithm is conducted on top of Pytorch. And the implementation of the homomorphic encryption for secure aggregation is based upon the library of TenSEAL.

Dataset. To evaluate the performance of FUSE on electricity theft detection tasks, we conduct experiments on a real-world dataset published by State Grid Cooperation of China (SGCC). The dataset consists of daily electricity consumption data collected by the SMs deployed in 42,372 electricity consumers' homes. For each consumer, it contains the data for 1,035 days from 2014 to 2016 where a regular consumer is labeled as 0 while the abnormal one is labeled as 1. There are 38,757 regular consumers and 3,615 abnormal consumers. Accordingly, the entire dataset is partitioned into a training set and a test set with a ratio of 9:1.

Theft Detection Model. In our work, to simulate the scenario wherein large-scale models are trained on those resourceconstrained devices, we choose an attention-based theft detection model with over 50,000,000 parameters [12]. As mentioned in Appendix C, the theft detection model is split into the Feature Extractor, the Feature Learner, and the Feature Classifier, with 106,785 parameters, 50,808,706 parameters, and 2,050 parameters respectively. The local training process follows Algorithm C2.

Federated Setting. Three TDs and corresponding CSs are considered in our experiments. In addition, the number of SMs differs over these TDs, thereby resulting in different local models with varying sample sizes. Therefore, the data follows the quantity-based skew situation of non-IID [13].

Evaluation Metrics. Due to the sample imbalance phenomenon in the dataset, the loss and accuracy (ACC) as the only evaluation criteria are not sufficient. Thus, Area under the ROC Curve (AUC), recall rate, F1 score, and the Matthews correlation coefficient (MCC) are also incorporated in our analysis.

Appendix D.2 Hyperparamater Learning

In this subsection, we explore the hyperparameters from local training to global aggregation, such that the optimized hyperparameters can be applied in following experiments on the performance evaluation of FUSE. For the purpose of simulating the asynchronous communications between TD, CS, and the Aggregator, the maximum staleness of model submission for TD is set to 6, and for CS is set to 4.

Appendix D.2.1 Impact of Hyperparameter on Local Training

1) Impact of local optimizer: In classical federated learning, Stochastic Gradient Descent (SGD) is the most commonly used local optimization algorithm [1]. However, RAdam is deployed in [12] instead of SGD in a centralized learning setting. We thereby first conduct an experiment to explore whether RAdam performs better than SGD in a federated learning setting. We set batchsize = 100, m = 3, and $\theta = 0.6$, where m denotes the maximum number of models in buffers that will be aggregated and θ denotes the basic aggregation weight of the new models in each round. The results are presented in both Figure D1 and Table D1. We verify that RAdam with a learning rate $\eta = 0.001$ outperforms SGD with a learning rate $\eta = 0.1$. Furthermore, RAdam has a faster convergence rate and better performance in terms of the criteria used in our experiments. Therefore, RAdam is chosen for subsequent experiments in the federated learning context.



 ${\bf Figure \ D1} \quad {\rm Model \ comparisons \ with \ different \ local \ optimizers}$



 ${\bf Figure \ D2} \quad {\rm Model \ comparisons \ with \ different \ batch \ sizes}$



Figure D3 Model comparisons by varying m



Figure D4 Model comparisons by varying θ

Sci China Inf Sci 8

Item	Setting	ACC	Recall	$\mathbf{F1}$	AUC	MCC
Optimizer	SGD	0.930	0.255	0.379	0.881	0.408
	RAdam	0.951	0.582	0.670	0.929	0.653
	50	0.949	0.557	0.645	0.930	0.627
Batchsize	100	0.951	0.582	0.670	0.929	0.653
	200	0.950	0.548	0.653	0.930	0.642
	1	0.946	0.521	0.617	0.921	0.600
n	3	0.951	0.582	0.670	0.929	0.653
	5	0.949	0.562	0.651	0.931	0.633
	0.5	0.949	0.548	0.657	0.932	0.646
θ	0.6	0.951	0.582	0.670	0.929	0.653
	0.8	0.948	0.557	0.640	0.932	0.620

Table D1 Comparisons on the Performance between Models with Different Settings

2) Impact of batch size: Batch size is another significant factor affecting the performance of the model. Thus, an appropriate batch size should be determined. We set m = 3 and $\theta = 0.6$, and batch size is set three levels: 50, 100, and 200 respectively. The result are shown in Figure D2 and Table D1. As Figure D2 illustrates, the batch size has a negligible impact on the convergence rate, since the number of rounds to reach the convergence remains almost unchanged under each setting. However, best performance regarding the recall rate, F1 score, and MCC can be observed at the model trained with *batchsize* = 100. Therefore we set *batchsize* = 100 for the rest of the experiments.

Appendix D.2.2 Impact of Hyperparameters during Model Aggregation

1) Impact of m in top-m: Since m determines the number of the models in the buffers to be aggregated in a new round, we set $\theta = 0.6$ and set m = 1, 3, and 5 respectively in this experiment. Figure D3 and Table D1 compare the performance across the models where the worst convergence occurs when m = 1. Though a slightly faster convergence can be achieved when m = 5, it does not work better than m = 3 in terms of accuracy, recall rate, F1 score, and MCC. This is because the update of the global model is negligible during each aggregation when n is relatively small. On the other hand, the accumulated variance from the global optimum becomes larger as m increases, thereby leading to negative impact on the performance. Consequently, we set m = 3 for the following experiments.

2) Impact of θ : θ has a direct impact on the basic weight of local models for aggregation in each communication round. To explore such influence, we compare the model convergence and the other metrics over $\theta = 0.5, 0.6, 0.8$, respectively. The results are presented in Figure D4 and Table D1. As illustrated, the models with different θ are on par in convergence. Meanwhile, $\theta = 0.6$ performs better on accuracy, recall rate, F1 score, and MCC. Hence we set θ to 0.6 in the later experiments.

Appendix D.3 Result and Performance Analysis

In this section, we carry out experimental studies to compare how FUSE performs vis-à-vis a set of benchmark works. All the hyperparameters are obtained from the discussions in Apprendix D.2.

Appendix D.3.1 Comparisons on Performance

We first compare our framework with a synchronous aggregation algorithm FedAvg [1], an asynchronous aggregation algorithm FedAsync [3], and a semi-asynchronous aggregation algorithm FedSA [14]. The local training process for all models follows Algorithm C2 and the maximum staleness of model submission is set to 6 for TD, 4 for CS. Furthermore, we assume FedAvg removes the model when a straggler issue occurs [15]. As in FedAsync, we consider two strategies to address the asynchronous problem between TD and its corresponding CS during aggregation for the purpose of exploring the effects of aggregating different components of the full model asynchronously: FedAsync (wait) is to wait until all components of the model have arrived, whereas FedAsync is to perform aggregation asynchronously based upon the arrivals of components. In addition, FedSA filters the models according to a predefined threshold on staleness such that the models are aggregated based on the number of samples.

Figure D5 and Table D2 compare the model performance between our framework and the benchmarks. As expected, a faster convergence is achieved by FedAsync than FedAsync(wait). We also find that better performance and faster convergence are fulfilled through FUSE in comparison with FedAvg and FedAsync. This is because FedAvg removes models with high level of latency, thereby resulting in loss of features from the stragglers. Whereas only one model is aggregated in each communication round in FedAsync, leading to a lower rate of convergence. Meanwhile, better performance in terms of accuracy, recall rate, F1 score, and MCC is obtained by FUSE in contrast to FedSA, since the models with latencies are pre-aggregated in FUSE based upon cosine similarity, instead of imposing strong assumptions on threshold. On the other hand, not only the number of samples but the degree of model latency is taken into account in our framework, thus incentivizing a significantly better performance.

Appendix D.3.2 Overhead Analysis

Homomorphic Overhead: In this section, we evaluate the additional overhead introduced by the homomorphic encryption for secure aggregation in our framework. More specifically, Cheon-Kim-Kim-Song (CKKS) scheme is used in our experiments to encrypt the parameters [16]. The overhead brought on by the homomorphic encryption includes parameter encryption and decryption on TDs and CSs, and ciphertext computing of aggregation on the Aggregator.

Both encryption and decryption overheads are presented in Table D3, where the full model denotes the solution without U-shape split learning for local training. We find that our framework is more cost-effective, while the overhead on CS is significantly greater



Figure D5 Comparison on model performance

Table D2 Comparison on Model Performance between FUSE and the Benchmarks

Methods	Accuracy	Recall	$\mathbf{F1}$	AUC	MCC
FedAvg	0.938	0.332	0.476	0.911	0.505
FedAsync(wait)) 0.935	0.288	0.432	0.888	0.476
FedAsync	0.937	0.357	0.490	0.912	0.500
FedSA	0.950	0.524	0.636	0.931	0.627
FUSE	0.951	0.582	0.670	0.929	0.653

 Table D3
 Computation Time for Encryption & Decryption

	Encryption	Decryption
$v_{\scriptscriptstyle \mathrm{TD}}$	$129.12 \mathrm{ms}$	$31.03 \mathrm{ms}$
$v_{\rm \scriptscriptstyle CS}$	$60740.20\mathrm{ms}$	$13759.50\mathrm{ms}$
Full model	$64819.94\mathrm{ms}$	$14637.31\mathrm{ms}$

than TD. This is not surprising because the majority of the overhead is transferred from TD to CS such that large scale of models are processed on CS with higher capacities. To further explore the overhead of ciphertext aggregation in depth, we decouple the operations into atomic ones based on addition and multiplication before the testing. The results are shown in Table D4, where a denotes a scalar, $v_{\rm TD}$ and $v_{\rm CS}$ denote the serialized model parameters of TD and CS, respectively. Considering the enhancement on privacy-preserving and the scenario wherein such operations are performed by the Aggregator with sufficient capacities instead of TD, it would be worthwhile to compromise on the overhead.

 Table D4
 Computation Time for Operations in Ciphertext

Operation	Time Cost	Operation	Time Cost
$\llbracket a \rrbracket \cdot \llbracket a \rrbracket$	$3.00\mathrm{ms}$	$[\![v_{\scriptscriptstyle \mathrm{TD}}]\!] \cdot [\![v_{\scriptscriptstyle \mathrm{TD}}]\!]$	2313.10ms
$[\![v_{\scriptscriptstyle \mathrm{CS}}]\!] \cdot [\![v_{\scriptscriptstyle \mathrm{CS}}]\!]$	$1060607.01\mathrm{ms}$	$a \cdot \llbracket v_{\scriptscriptstyle TD} \rrbracket$	$22.02 \mathrm{ms}$
$a \cdot [\![v_{\scriptscriptstyle \mathrm{CS}}]\!]$	$9393.54 \mathrm{ms}$	$[\![a]\!] + [\![a]\!]$	$\approx 0.00 \mathrm{ms}$
$[\![v_{\scriptscriptstyle \mathrm{TD}}]\!]+[\![v_{\scriptscriptstyle \mathrm{TD}}]\!]$	4.00ms	$[\![v_{\scriptscriptstyle \mathrm{CS}}]\!]+[\![v_{\scriptscriptstyle \mathrm{CS}}]\!]$	$933.85 \mathrm{ms}$
$[\![a]\!]\cdot[\![v_{\scriptscriptstyle \mathrm{TD}}]\!]$	$588.54\mathrm{ms}$	$[\![a]\!]\cdot[\![v_{\scriptscriptstyle \mathrm{CS}}]\!]$	$263427.03\mathrm{ms}$

Communication and Computing Overhead: We evaluate the communication and computational overhead from the perspective of the entities in FUSE. Since both of them are positively related to the number of model parameters, such overhead can be closely approximated by the model size. Assuming the number of local iterations is E we present the results in Table D5. To demonstrate that FUSE enables reducing the computational overheads on resource-constrained devices, we compare our framework with a vanilla federated learning-based framework, i.e., the whole local training process executed by SM or TD. Given a traditional approach without U-shape split learning, the computational overhead can be represented as $E \cdot [FP(|w|) + BP(|w|)]$, where $|w| = |w_{\rm FE}| + |w_{\rm FL}| + |w_{\rm FC}|$. Since $|w_{\rm FL}| \gg |w_{\rm FE}| + |w_{\rm FC}|$, the computational overheads of TD or SM are significantly lower than $E \cdot [FP(|w|) + BP(|w|)]$. As a result, the analysis highlights that the computational overheads on SM and TD are migrated to assisted cloud servers, in contrast to the solution without U-shape split learning. However, this trade-off of reducing computational overhead on resource-constrained devices comes at the cost of increased overall communication overhead. The communication overhead of TD can be represented as $|w| = |w_{\rm FE}| + |w_{\rm FC}|$ while the whole local training process is executed completely by TD. In our scheme, the communication overhead of TD can be further represented as $(E \cdot |\rm{SM}| + 1) \cdot |w_{\rm FE}| + E \cdot (|output_{\rm FE}| \cdot |\rm{SM}| + |output_{\rm FL}|) + |w_{\rm FC}|$ according to Table D5, since $|y_{\rm FC}^*| = |output_{\rm FL}|$, where $|y_{\rm FL}^*|$

denotes the gradient of the first layer of the Feature Learner. Therefore, the communication overhead of TD is reduced if and only if $|w_{\rm FL}| > E \cdot |\rm SM| \cdot |w_{\rm FE}| + E \cdot (|output_{\rm FE}| \cdot |\rm SM| + |output_{\rm FL}|)$. This inequality holds true in most cases, as the parameters of the Feature Learner constitute the majority of the overall model parameters. Therefore, our approach can significantly reduce communication overhead for TD in most scenarios.

 Table D5
 Communication & Computing Overhead In One Communication Round

	Computing	Communication
SM	$E \cdot FP(w_{{}_{\mathrm{FE}}})$	$E \cdot output_{\scriptscriptstyle \mathrm{FE}} $
		to TD
TD	$\frac{E \cdot [FP(w_{_{\rm FC}}) + BP(w_{_{\rm FC}}) + BP(w_{_{\rm FC}}) + BP(w_{_{\rm FE}})]}{BP(w_{_{\rm FE}})]}$	$\underbrace{E \cdot w_{\scriptscriptstyle \rm FE} \cdot {\rm SM} }_{E \cdot} + \underbrace{E \cdot \mathit{output}_{\scriptscriptstyle \rm FE} \cdot {\rm SM} + E \cdot g_{\scriptscriptstyle \rm FC}^* }_{E \cdot} + \underbrace{ w_{\scriptscriptstyle \rm FE} + w_{\scriptscriptstyle \rm FC} }_{E \cdot}$
		to SM to CS to Aggregator
\mathbf{CS}	$E \cdot [\mathit{FP}(w_{\scriptscriptstyle \mathrm{FL}}) + \mathit{BP}(w_{\scriptscriptstyle \mathrm{FL}})]$	$\underbrace{E \cdot output_{_{\rm FL}} + E \cdot g_{_{\rm FL}}^* }_{\text{FL}} + \underbrace{ w_{_{\rm FL}} }_{\text{FL}}$
		to TD to Aggregator

* $FP(\cdot)$ and $BP(\cdot)$ denotes the function of computing overhead of forward propagation and back-propagation respectively.

Appendix E Related Work

Appendix E.1 Distributed Theft Detection Method

Typically, the increasing demands on privacy-preserving may not be fulfilled by the centralized methods satisfactorily, and even though studies are conducted on enhancing the security through cryptography techniques, they are not suitable for the scenario of Advanced Metering Infrastructure(AMI) wherein computational resources are limited due to the extensive overhead required for decryption and encryption. Thus, federated learning and split learning have emerged as major offerings for theft detection models by providing commitment on protecting sensitive data. FedDetect is the first federated learning-based schema [17] where data privacy is guaranteed by homomorphic encryption and local differential privacy. However, the work has not considered the substantial computational overhead for the clients and asynchronous model updates. Ashraf et al. [18] subsequently introduce FedDP, which presents a novel federated voting classifier specifically for theft detection. In contrast, a decentralized federated learning-based theft detector is proposed in [19]. This approach designs a secure and communication-efficient functional encryption-based aggregation scheme, eliminating the need for a centralized key distribution center to generate security parameters. However, all these schemes involve collecting the electricity consumption data from SMs to detection stations, which exposes them to the risk of raw data leakage. Alromih et al. [20] employ split learning in the smart grid to ensure better privacy protection. However, they model the detection as an unsupervised task, which is not appropriate for the scenarios where labeled data is available.

Appendix E.2 Federated Learning

Federated learning is proposed by Google which allows the data owners to train a global model collaboratively without exposing the raw data [1]. Due to the ability to preserve the privacy of raw data, federated learning is increasingly applied in various fields, including industry [21,22], healthcare [23], and traffic [24]. Particularly, FedAvg [1] is the first algorithm for model aggregation and its performance and effectiveness has been verified in many studies. However, the straggler issue due to heterogeneity of bandwidths among devices is a bigger challenge in realistic environment. Li et al. [15] propose a synchronous aggregation algorithm where such issue is addressed by adjusting the local epoch of the straggler, while the performance has not been improved. On the other hand, a faster convergence can be achieved in asynchronous federated learning schema since the aggregation occurs when the update is received by the aggregator. Xie et al. [3] design FedAsync to improve both the flexibility and scalability of federated learning on massive devices with varying resources. Besides, asynchronous aggregation algorithms such as ASO-Fed and PAFLM are proposed to accommodate the model submissions at different rates [25, 26]. However, the extensive communication overhead brought on by the asynchronous aggregation strategy is placing additional hurdles on resource-constrained devices.

Balancing the trade-offs between the model performance, convergence rate, and communication overhead, semi-asynchronous aggregation mechanisms are developed by combining the advantages of both synchronous and asynchronous aggregations. A first-come-first-merge client selection mechanism is incorporated in SAFA to improve the efficiency of the training process and the quality of the global model [5]. FedSA is another semi-asynchronous framework proposed by Ma et al., where the number of local models for aggregation is optimized in order to minimize the training time [14]. To our knowledge, metrics that are critical to assuring model quality other than latency have not been considered in this study.

References

- 1 McMahan B, Moore E, Ramage D, et al. Communication-efficient learning of deep networks from decentralized data. In: Artificial intelligence and statistics, 2017. 1273-1282.
- 2 Xu C H, Qu Y Y, Xiang Y, et al. Asynchronous federated learning on heterogeneous devices: A survey. Computer Science Review, 2023, 50: 100595.
- 3 Xie C, Koyejo S, Gupta I. Asynchronous federated optimization. arXiv preprint, 2019, arXiv:1903.03934.
- 4 Zhang Y, Duan M, Liu D, et al. Csafl: A clustered semi-asynchronous federated learning framework. In: 2021 International Joint Conference on Neural Networks (IJCNN), 2021. 1-10.
- 5 Wu W T, He L A, Lin W W, et al. Safa: A semi-asynchronous protocol for fast federated learning with low overhead. IEEE Transactions on Computers, 2020, 70(5):655-668.
- 6 Hao J S, Zhao Y C, Zhang J L. Time efficient federated learning with semi-asynchronous communication. In: 2020 IEEE 26th International Conference on Parallel and Distributed Systems (ICPADS), 2020. 156-163.

- 7 Vepakomma P, Gupta O, Swedish T et al. Split learning for health: Distributed deep learning without sharing raw patient data. arXiv preprint arXiv:1812.00564, 2018.
- 8 Thapa C, Arachchige P C M, Camtepe S, et al. Splitfed: When federated learning meets split learning. In: Proceedings of the AAAI Conference on Artificial Intelligence, 2022. volume 36, pages 8485-8493.
- 9 Alromih A, Clark J A, Gope P. Privacy-aware split learning based energy theft detection for smart grids. In: Information and Communications Security: 24th International Conference, ICICS 2022, Canterbury, UK, 2022. 281-300.
- 10 Wang T, Liu Y, Zheng X, et al. Edge-based communication optimization for distributed federated learning. IEEE Transactions on Network Science and Engineering, 2021, 9(4):2015-2024.
- 11 Cheon J H, Kim D, Kim D, et al. Numerical method for comparison on homomorphically encrypted numbers. In: Advances in Cryptology-ASIACRYPT 2019: 25th International Conference on the Theory and Application of Cryptology and Information Security, Kobe, Japan, 2019. Part II, pages 415-445.
- 12 Finardi P, Campiotti I, Plensack G, et al. Electricity theft detection with self-attention. arXiv preprint, 2020, arXiv:2002.06219.
- 13 Li Q B, Diao Y Q, Chen Q, et al. Federated learning on non-iid data silos: An experimental study. In: 2022 IEEE 38th International Conference on Data Engineering (ICDE), 2022. 965-978.
- 14 Ma Q P, Xu Y, Xu H L, et al. Fedsa: A semi-asynchronous federated learning mechanism in heterogeneous edge computing. IEEE Journal on Selected Areas in Communications, 2021, 39(12):3654-3672.
- 15 Li T, Sahu A K, Zaheer M, et al. Federated optimization in heterogeneous networks. Proceedings of Machine learning and systems, 2020, 2:429-450.
- 16 Cheon J H, Kim A, Kim M, et al. Homomorphic encryption for arithmetic of approximate numbers. In: Advances in Cryptology-ASIACRYPT 2017: 23rd International Conference on the Theory and Applications of Cryptology and Information Security, Hong Kong, China, 2017. Part I 23, pages 409-437.
- 17 Wen M, Xie R, Lu K J, et al. Feddetect: a novel privacy-preserving federated learning framework for energy theft detection in smart grid. IEEE Internet of Things Journal, 2021, 9(8):6069-6080.
- 18 Ashraf M M, Waqas M, Abbas G, et al. Feddp: A privacy-protecting theft detection scheme in smart grids using federated learning. Energies, 2022, 15(17):6241.
- 19 Ibrahem M I, Mahmoud M, Fouda M M, et al. Privacy preserving and efficient decentralized federated learning-based energy theft detector. In: GLOBECOM 2022-2022 IEEE Global Communications Conference, 2022. 287-292.
- 20 Alromih A, Clark J A, Gope P. Privacy-aware split learning based energy theft detection for smart grids. In: Information and Communications Security: 24th International Conference, ICICS 2022, Canterbury, UK, 2022. 281-300.
- 21 Chen J B, Xue J F, Wang Y, et al. Privacy-Preserving and Traceable Federated Learning for data sharing in industrial IoT applications. Expert Systems with Applications, 2023, 213: 119036.
- 22 Wang N Y, Li X, Guan Z T, et al. FedStream: A Federated Learning Framework on Heterogeneous Streaming Data for Next-Generation Traffic Analysis. IEEE Transactions on Network Science and Engineering, 2023.
- 23 Bashir A K, Victor N, Bhattacharya S, et al. Federated learning for the healthcare metaverse: Concepts, applications, challenges, and future directions. IEEE Internet of Things Journal, 2023.
- 24 Liu L, Tian Y X, Chakraborty C, et al. Multilevel Federated Learning based Intelligent Traffic Flow Forecasting for Transportation Network Management. IEEE Transactions on Network and Service Management, 2023.
- 25 Chen Y J, Ning Y, Slawski M, et al. Asynchronous online federated learning for edge devices with non-iid data. In: 2020 IEEE International Conference on Big Data (Big Data), 2020. 15-24.
- 26 Lu X F, Liao Y Y, Lio P, et al. Privacy-preserving asynchronous federated learning mechanism for edge network computing. IEEE Access, 2020, 8:48970-48981.