

# Practical Vertical Federated Learning with Unsupervised Representation Learning

Zhaomin Wu<sup>1</sup>, Qinbin Li<sup>2</sup>, Bingsheng He<sup>3</sup>  
 {zhaomin<sup>1</sup>, qinbin<sup>2</sup>, hebs<sup>3</sup>}@comp.nus.edu.sg  
 National University of Singapore

**Abstract**—As societal concerns on data privacy recently increase, we have witnessed data silos among multiple parties in various applications. Federated learning emerges as a new learning paradigm that enables multiple parties to collaboratively train a machine learning model without sharing their raw data. *Vertical federated learning*, where each party owns different features of the same set of samples and only a single party has the label, is an important and challenging topic in federated learning. Communication costs among different parties have been a major hurdle for practical vertical learning systems. In this paper, we propose a novel communication-efficient vertical federated learning algorithm named *FedOnce*, which requires only one-shot communication among parties. To improve model accuracy and provide privacy guarantee, FedOnce features unsupervised learning representations in the federated setting and privacy-preserving techniques based on moments accountant. The comprehensive experiments on 10 datasets demonstrate that FedOnce achieves close performance compared to state-of-the-art vertical federated learning algorithms with much lower communication costs. Meanwhile, our privacy-preserving technique significantly outperforms the state-of-the-art approaches under the same privacy budget.

**Index Terms**—vertical federated learning, differential privacy

## 1 INTRODUCTION

**D**ATA silos, where physically distributed data cannot be collected on a central server due to privacy concerns and regulations, have been witnessed in many real-world applications. *Federated learning* [1], [2], [3], [4], which can be classified as *horizontal federated learning* and *vertical federated learning* according to the data partitioning among parties [5], is proposed as an emerging learning paradigm to enable collaborative training without exposing the data of each party. Different from horizontal federated learning where each party shares the same feature space but owns only a subset of samples, vertical federated learning, where each party shares the same sample space but only holds a subset of features, is a common and important scenario in practice.

*Privacy by design* (PbD), as a key concept in General Data Protection Regulation (GDPR) [6], states that personal data from different services/applications under the same company should not be directly shared with each other. This concept motivates our study to develop a collaborative machine learning framework to train a privacy-preserving model from all parties. Particularly, we consider each service/application as a party in federated learning and assume that the party holding labels, named *host party*, is trusted by other parties without labels, named *guest parties*. Such trust is practical as all the parties are under the same company. Fig. 1 presents a concrete example of this setting, which is an end-to-end procedure containing model training and model deployment. An e-commerce company provides two major services: 1) an online shopping service managed by the retailing department and 2) an e-banking service managed by the Fintech department. This company aims to train a machine learning model to recommend its users the proper investing amount. Under the principle of privacy by

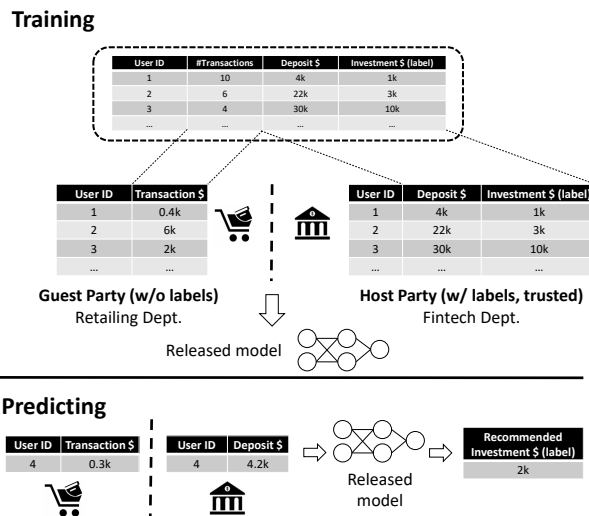


Fig. 1. Example of vertical federated learning

design, although the Fintech department can be trusted by the retailing department, direct data sharing is prohibited. Therefore, these departments have to perform training by transferring intermediate information. After the training, a model can be released to provide recommendation services for the users of this e-commerce company. Typical applications of federated learning includes Google Keyboard [7] and nationwide learning health system [8].

Considering the model training, communication issues have limited the practical adoption of vertical federated learning in many real-world applications [9]. For example,

such limitation becomes non-negligible when parties are 1) mobile devices communicating through Cellular networks, 2) IoT devices communicating through Wi-Fi networks, 3) seagoing ships communicating through satellites. These applications, featuring unstable connections and expensive communication cost, cannot be handled by existing vertical federated learning algorithms for two major reasons, including 1) **high stability requirement**: existing vertical federated learning approaches require all the parties to stay online or even synchronous during the entire training process, which is unrealistic for parties connected by unstable networks; 2) **high communication cost**: current studies in vertical federated learning usually incur large communication overhead, leading to the considerable economic cost. For example, in each training iteration, SecureBoost [10], VF<sup>2</sup>Boost [11], and Pivot [12] exchange encrypted intermediate results (e.g, gradients) between the host party and guest parties. Similarly, when training each batch of data, SplitNN [13] transfers outputs and gradients of a common layer between the host party and guest parties. All these approaches require batch/iteration-level synchronization and high communication cost.

To improve the practicability of vertical federated learning, designing a learning algorithm with a single communication round (i.e., one-shot), which can highly reduce the stability requirement and the communication cost, becomes increasingly desirable. Although some one-shot approaches [14], [15], [16] have been proposed for horizontal federated learning where each party shares the same feature space but owns only a subset of samples, these approaches cannot be directly adopted in vertical federated learning. The key challenge is that the labels only exist in a single party in vertical federated learning. These labels, which cannot be shared among parties due to privacy concerns, are necessary for the training of most models. For example, calculating gradients in neural networks or gradient boosting decision trees [17] relies on the labels. To the best of our knowledge, one-shot vertical federated learning algorithms have not been explored yet.

In this paper, we propose FedOnce, a novel one-shot vertical federated learning algorithm, the source code<sup>1</sup> of which is available online. FedOnce is designed for neural networks which are broadly used in many applications. Instead of aggregating the local models in horizontal learning, we propose the idea that the host party trains a model to aggregate the representations collected from guest parties. Specifically, guest parties first extract representative features by *unsupervised learning* [18]. These representative features, which can effectively boost the performance of federated learning, are sent to the host party. After receiving all the representative features, the host party trains an aggregation model based on these features and its own data.

Considering the model deployment after the training, these released models are threatened by membership inference attack [19] which detects the existence of training instances by inference. To protect the trained models in FedOnce from this attack, we exploit differential privacy [20] to theoretically guarantee that the training samples is indistinguishable given the trained model. Despite many

studies that apply differential privacy to horizontal federated learning [21], [22], [23], their privacy analyses cannot be directly applied to vertical federated learning. The main reason is that vertical federated learning requires analyzing inter-party privacy loss since each data record is distributed to multiple parties, whereas horizontal federated learning only requires analyzing inner-party privacy loss since each data record is held by a single party. Some existing studies [24], [25], [26] employ differential privacy in vertical federated learning. Nonetheless, these studies, focusing on inner-party privacy loss, only consider simple composition when analyzing inter-party privacy loss. In FedOnce, by refining the analysis of inter-party privacy loss based on moments accountant [27], we significantly reduce the privacy loss compared to existing approaches.

Our main contribution can be summarized as:

- We propose a practical one-shot vertical federated learning algorithm *FedOnce* by effectively taking advantage of unsupervised learning.
- We develop privacy techniques for FedOnce under the notion of differential privacy. Meanwhile, we prove a tighter bound on the privacy loss across parties compared to the current work on vertical federated learning.
- Our experiments demonstrate that FedOnce reduces the required communication cost to achieve competitive performance against state-of-the-art vertical federated learning algorithms [10], [13] by at most 94%. Meanwhile, our privacy-preserving technique significantly outperforms other approaches under the same privacy budget.

## 2 PRELIMINARIES

### 2.1 Differential Privacy

Differential privacy [28] is a technique that provides a theoretical guarantee of the privacy of individuals.

**Definition 1.** (Differential Privacy) Suppose  $\mathcal{M} : \mathcal{D} \rightarrow \mathcal{R}$  is a randomized mechanism,  $D$  and  $D'$  are two neighboring datasets in domain  $\mathcal{D}$ , i.e.,  $\|D - D'\| \leq 1$ . Then,  $\mathcal{M}$  satisfies  $(\epsilon, \delta)$ -differential privacy if for any outputs  $S \subseteq \mathcal{R}$ ,

$$\Pr[\mathcal{M}(D) \in S] \leq e^\epsilon \Pr[\mathcal{M}(D') \in S] + \delta$$

where  $\epsilon$  and  $\delta$  are privacy parameters.

In the definition of  $(\epsilon, \delta)$ -differential privacy,  $\epsilon$  bounds the difference between the outputs of two neighboring datasets, and this bound might be violated by a small probability  $\delta$ , which is typically smaller than  $1/|D|$ , where  $|D|$  is the number of samples in  $D$ .

Among many privacy mechanisms that achieve differential privacy, *Gaussian Mechanism* [20] is widely adopted in deep learning. It ensures the differential privacy of a deterministic function  $f$  by adding Gaussian noise to the output. The noise scale  $\sigma$  is determined by the privacy parameters  $\epsilon, \delta$  and the  $\ell_2$ -sensitivity of the function.

**Definition 2.** ( $\ell_2$ -sensitivity) The  $\ell_2$ -sensitivity of function  $f : \mathcal{D} \rightarrow \mathcal{R}$  is defined as

$$\Delta_2(f) \triangleq \max_{D, D'} \|f(D) - f(D')\|_2$$

1. <https://github.com/JerryLife/FedOnce>

where  $D, D'$  are neighboring datasets, i.e.,  $\|D - D'\| \leq 1$ .

**Theorem 1.** (Gaussian Mechanism) For any  $\varepsilon \in (0, 1)$ ,  $c^2 > 2 \ln(1.25/\delta)$ ,  $\sigma = c\Delta_2(f)/\varepsilon$ , Gaussian mechanism

$$\mathcal{M}(D) \triangleq f(D) + \mathcal{N}(0, \sigma^2 \mathbf{I})$$

satisfies  $(\varepsilon, \delta)$ -differential privacy.

Gaussian mechanism achieves the differential privacy of a single function. Furthermore, if multiple mechanisms are applied to one dataset, the overall differential privacy follows simple composition [20].

**Theorem 2.** (Simple Composition) Let  $\mathcal{M}_i : \mathcal{D} \rightarrow \mathcal{R}$  be a  $(\varepsilon_i, \delta_i)$ -differential privacy mechanism. Then a composition of  $k$  such mechanisms  $\mathcal{M}_k = (\mathcal{M}_1, \dots, \mathcal{M}_k)$  satisfies  $(\sum_{i=1}^k \varepsilon_i, \sum_{i=1}^k \delta_i)$ -differential privacy.

Simple composition holds regardless of the distribution of  $\mathcal{M}_i$  as long as it satisfies  $(\varepsilon_i, \delta_i)$ -differential privacy. When all the mechanisms  $\mathcal{M}_i$  are Gaussian mechanisms, this privacy bound can be improved by moments accountant [27] which first bounds the Rényi divergence [29] of the composition and then calculates the overall privacy loss  $\varepsilon$  given a specific  $\delta$ . The main theorem of moments accountant is concluded as Theorem 3.

**Theorem 3.** Given the sampling probability of SGD  $q$ , the number of epochs  $T$ , there exists constants  $c_1$  and  $c_2$  so that for  $\forall \varepsilon < c_1 q^2 T$ ,  $\forall \delta > 0$ , differentially private SGD [27] with noise scale  $\sigma$  is  $(\varepsilon, \delta)$ -differential privacy if we choose

$$\sigma \geq c_2 \frac{q \sqrt{\log(1/\delta)T}}{\varepsilon}$$

The proof of moments accountant relies on the definition of privacy loss and  $\lambda$ -th moments. Auxiliary input (e.g. hyperparameters) in the original definition is omitted for simplicity.

**Definition 3.** (Privacy Loss) For two neighboring datasets  $D$  and  $D'$ , an output  $o$  and a mechanism  $\mathcal{M}$ , the privacy loss is defined as

$$c(o; \mathcal{M}, D, D') \triangleq \ln \frac{\Pr[\mathcal{M}(D) = o]}{\Pr[\mathcal{M}(D') = o]}$$

**Definition 4.** ( $\lambda$ -th moments) The  $\lambda$ -th moments is defined as the log of the moments generating function of privacy loss evaluated at value  $\lambda$ :

$$\alpha_{\mathcal{M}}(\lambda; D, D') \triangleq \ln \mathbb{E}_{o \sim \mathcal{M}(D)}[\exp \lambda c(o; D, D')]$$

Furthermore, the maximum of all possible  $\lambda$ -th moments is defined as

$$\alpha_{\mathcal{M}}(\lambda) \triangleq \max_{D, D'} \alpha_{\mathcal{M}}(\lambda; D, D')$$

In [27], the authors first prove the moments bound  $\alpha_{\mathcal{M}}(\lambda)$  for a single Gaussian mechanism (Lemma 1). Secondly, they prove how the bound accumulates in a sequence of Gaussian mechanisms (Lemma 2). Finally, the relationship between this moments bound and differential privacy parameters are proved according to Lemma 3.

**Lemma 1.** Suppose that  $f : D \rightarrow \mathbb{R}^p$  with  $\|f(\cdot)\|_2 \leq 1$ . Suppose  $\sigma \geq 1$  and  $J$  is a sample from  $[n]$  where each  $i \in [n]$  is chosen independently with probability  $q \leq \frac{1}{16\sigma}$ . Then, for

any positive integer  $\lambda \leq \sigma^2 \ln \frac{1}{q\sigma}$ , the Gaussian mechanism  $\mathcal{M}(d) = \sum_{i \in J} f(d_i) + \mathcal{N}(0, \sigma^2 \mathbf{I})$  satisfies

$$\alpha_{\mathcal{M}}(\lambda) \leq \frac{q^2 \lambda (\lambda + 1)}{(1 - q)\sigma^2} + O(q^3 \lambda^3 / \sigma^3)$$

**Lemma 2.** Suppose  $\mathcal{M}$  contains a sequence of  $T$  adaptive mechanisms  $\mathcal{M}_1, \dots, \mathcal{M}_T$ , where  $\mathcal{M}_i : \prod_{j=1}^{i-1} \mathcal{R}_j \times \mathcal{D} \rightarrow \mathcal{R}_i$ . For  $\forall \lambda > 0$ , we have

$$\alpha_{\mathcal{M}}(\lambda) \leq \sum_{i=1}^T \alpha_{\mathcal{M}_i}(\lambda)$$

**Lemma 3.** For any  $\varepsilon > 0$ , the mechanism  $\mathcal{M}$  is  $(\varepsilon, \delta)$ -differential private if

$$\delta = \min_{\lambda} \exp(\alpha_{\mathcal{M}}(\lambda) - \lambda \varepsilon)$$

## 2.2 Unsupervised Learning

Unsupervised learning trains a model without labels. With this model, each sample can be tagged as a feature map, named the *representation* of this sample.

FedOnce does not rely on specific unsupervised learning algorithms. In this paper, we choose NAT [18], which is a simple and effective unsupervised learning algorithm suitable for both image datasets and multivariate datasets. In NAT, in order to learn a representation  $R : n \times d$  of  $n$  samples with  $d$  dimensions, a random representation  $C : n \times d$  is generated. Denoting the neural network model to generate  $R$  as  $f_{\theta}(\cdot)$ , the goal of NAT can be expressed as the Equation 1.

$$\min_{P, \theta} \frac{1}{2n} \|f_{\theta}(D) - PC\|_F^2 \quad (1)$$

where  $P$  is a permutation matrix,  $D$  is the dataset without labels, and  $PC$  (i.e., the matrix multiplication of  $P$  and  $C$ ) is a permutation of  $C$ .

In each iteration of training,  $\theta$  is then optimized by gradient descent. For every  $t_P$  iterations,  $P$  is optimized by Hungarian algorithm [30]. The updating frequency  $t_P$  can be adjusted as a hyperparameter. After the training process, the representation can be obtained by predicting with model  $\theta$ , i.e.,  $R = f_{\theta}(D)$ .

## 2.3 SplitNN

SplitNN [13] is a simple and effective vertical federated learning<sup>2</sup> model based on neural network. As explained in Fig. 2, models are split into multiple parties. Each guest party holds a local model; the host party holds a local model and an aggregation model. In forward propagation, first, all the parties forward propagate independently using their local models. Then, the outputs of local models on guest parties are sent to the host party. The host party concatenates the outputs from all the parties (including itself) and feeds the concatenated output into the aggregation model  $\theta_{agg}$ . The aggregation model continues forward propagating to produce the final prediction. In the backpropagation, first, the loss and gradients are calculated based on the final prediction. Then, the gradients w.r.t. to  $\theta_1, \theta_2, \theta_3$  are calculated

<sup>2</sup> Though some studies categorize SplitNN as split learning, we denote SplitNN as a vertical federated learning algorithm for simplicity.

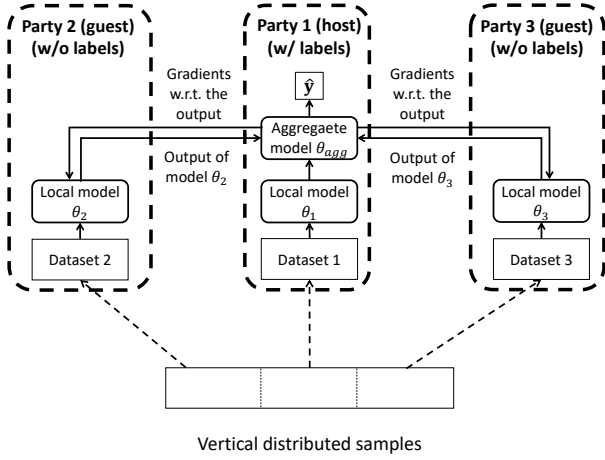


Fig. 2. Structure of SplitNN

by backpropagating through  $\theta_{agg}$ . The gradients w.r.t.  $\theta_2, \theta_3$  are sent to two guest parties, respectively. Finally, all the parties perform backpropagation on their local models. In each iteration of training, both forward propagation and backpropagation are performed for one time. Therefore, the training incurs two communication rounds between the host party and each guest party every iteration, which is impractical when parties are connected by unstable and expensive networks. This pitfall of SplitNN motivates our design of a one-shot vertical federated learning framework.

### 3 THE FEDONCE FRAMEWORK

In this section, we propose the framework of FedOnce. FedOnce with no protection on released models is denoted as FedOnce-L0, corresponding to the scenario where the trained model is used inside the company, thus free from privacy attacks. FedOnce ensuring differential privacy of the released model is denoted as FedOnce-L1, corresponding to the scenario where the trained model is released to the public (either as white-box or black-box), thus threatened by the membership inference attack [19]. Our objective is to solve the problem stated in section 3.1 by one-shot communication. The framework is presented in Figure 3.

#### 3.1 Problem Formulation

**FedOnce-L0** Given a global dataset  $\{\mathbf{x}_i, y_i\}_{i=1}^n$ , the features of each sample  $\mathbf{x}_i$  is distributed among  $k$  parties (i.e.,  $\mathbf{x}_i = [\mathbf{x}_i^1, \mathbf{x}_i^2, \dots, \mathbf{x}_i^k]$ ) while only one party has the labels  $\mathbf{y} = \{y_i\}_{i=1}^n$ . Without loss of generality, we assume party 1 is the host party and the other parties are guest parties. We denote the local dataset on party  $j$  as  $\mathbf{x}^j \triangleq \{\mathbf{x}_i^j\}_{i=1}^n$ . All the parties need to collaboratively train a model  $f(\theta; \mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^k)$  without exchanging their feature vectors or labels with each other. **One communication round with the host party is permitted for each guest party.** The training process can be formulated as

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n L(f(\theta; \mathbf{x}_i^1, \mathbf{x}_i^2, \dots, \mathbf{x}_i^k); y_i) + \lambda \omega(\theta)$$

where  $L(f(\theta; \mathbf{x}_i^1, \mathbf{x}_i^2, \dots, \mathbf{x}_i^k); y_i)$  is the loss function and  $\lambda \omega(\theta)$  is the regularization term.

Note that in the setting of FedOnce-L0, we assume that the samples have already been aligned among the parties (i.e., party  $j$  knows  $\mathbf{x}_i^j$  contains partial features of  $\mathbf{x}_i$ ). Such alignment can be implemented by privacy-preserving record linkage [31], which is an orthogonal topic to this paper. This is a common setting in vertical federated learning [10], [11], [12], [32].

**FedOnce-L1** In FedOnce-L1, we focus on the privacy risk of model deployment after the training. We assume the host party is trusted by guest parties. The attackers can obtain host and guest models including all the parameters required for inference. The goal of FedOnce-L1 is to protect the released models against membership inference attack [33] by ensuring  $(\epsilon, \delta)$ -differential privacy of the released model. The details on privacy design are presented in Section 4.

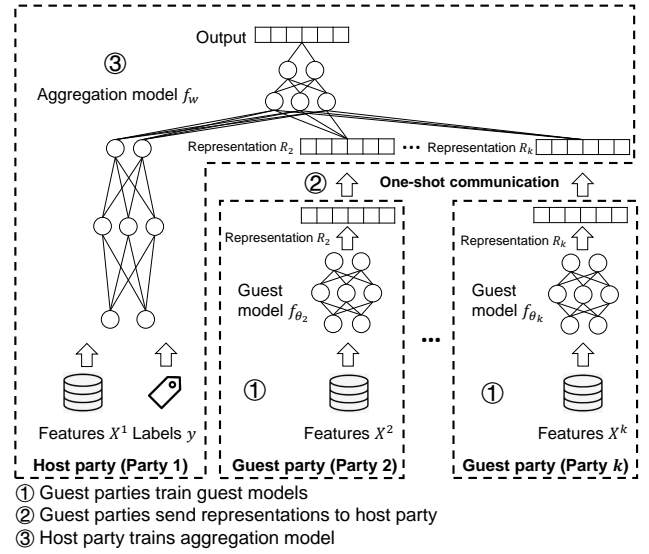


Fig. 3. Framework of FedOnce

#### 3.2 Training

Suppose among  $k$  parties  $\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_k$ ,  $\mathcal{P}_1$  is the host party and  $\mathcal{P}_2, \dots, \mathcal{P}_k$  are guest parties. In the training process, guest parties first extract representations of their data by unsupervised learning based on [18]. The representations, reflecting the high-level features of data on guest parties, are then sent to the host party, which is the only communication required in FedOnce. Finally, the host party collects all the representations and trains an aggregation model together with its own data. The training process of FedOnce is summarized as Algorithm 1, where the regularization term is omitted for simplicity. For convenience, we denote  $R_i, \dots, R_j$  as  $R_{i:j}$ . Although the training is performed by batch in practice, we present the training procedure by epoch in Algorithm 1 for simplicity.

Specifically, guest parties first initialize random targets  $C_j$  (line 3). In each epoch  $t$ , Hungarian algorithm [30] is adopted to find the best permutation of  $C_j$ , denoted as  $\tilde{P}^{(t)}$ , that minimizes the loss (line 5). Because of the high computational complexity, the Hungarian algorithm finds the best permutation in each batch rather than in the whole dataset. Later, this permutation is used to calculate



the gradient w.r.t. parameters  $\theta_j^{(t)}$  (line 6-7). If the released model requires no protection, i.e., FedOnce-L0,  $\theta_j$  is directly updated by the gradients (line 9). Otherwise, in FedOnce-L1, the gradients are clipped to bound the sensitivity (line 11). Gaussian noise is then generated and added to the gradients (line 12). Finally,  $\theta_j^{(t)}$  is updated by these gradients with noise (line 13). After each guest party  $\mathcal{P}_j (j \in [2, k])$  completes training, the calculated representation  $R_j$  is sent to the host party (line 16).

Once receiving the representations from all guest parties, the host party starts to train an aggregation model  $f_w$ , which takes representations  $R_{2:j}$  and its own data  $x_1$  as inputs and outputs the final prediction. First, gradients w.r.t.  $w$  are calculated (lines 21-22). In FedOnce-L0, the gradients are directly utilized to update  $w$  (line 24). Otherwise, in FedOnce-L1, gradients are clipped (line 26); Gaussian noise is generated and added to the gradients (lines 26-27). Finally, model parameters  $w$  are updated by these gradients with noise (line 28).

**Extension to Multi-Round Training.** While achieving impressive performance in one communication round, FedOnce can be further improved if more communication rounds are permitted. At the end of the training process of FedOnce, each party holds a guest model and the host party additionally holds an aggregation model; the combination of these models can be regarded as a pre-trained model of SplitNN. Therefore, the training process of SplitNN can be directly applied to further boost the performance of FedOnce. Specifically, in each iteration, guest parties perform forward propagation and transfer intermediate outputs to the host party. The host party receives the outputs and continues forward propagation in the aggregation model to obtain predictions  $\hat{y}$ . After calculating the loss with  $\hat{y}$  the labels  $\mathbf{y}$ , the host party performs backpropagation and transfers gradients to guest parties. Finally, the received gradients are used to perform backpropagation to update guest models. We denote this extension as *multi-round FedOnce-L0*.

### 3.3 Prediction

The prediction is conducted by all the parties collaboratively according to Algorithm 2. The features of testing data are distributed on  $k$  parties, i.e.,  $\mathbf{x} = (\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^k)$ . First, each guest party  $\mathcal{P}_j$  makes a prediction according to its guest model  $f_{\theta_j}$  (line 3). These predicted representations  $R_j$  are then sent to the host party (line 4) which makes the final prediction based on  $f_{\theta_j}(\mathbf{x}_i^j)$  and partial sample  $\mathbf{x}_i^1$  (line 7).

## 4 PRIVACY

### 4.1 Privacy of Released Model

FedOnce-L1 adopts differential privacy to ensure each sample cannot be distinguished through the released model. The noise addition method largely follows [27]. As demonstrated in Algorithm 1, gradients are first clipped by a threshold  $\Omega$  in each iteration of the training process to ensure that  $\ell_2$ -sensitivity of all the gradients are bounded by  $\Omega$ , i.e.,  $\|\mathbf{g}_{clip}\|_2 \leq \Omega$ . After gradient clipping, independent Gaussian noise is added to each gradient.

Compared to the original moments accountant, where the privacy loss is only accumulated across iterations, the

---

### Algorithm 1 Training Process of FedOnce

---

**Require:** training data  $(\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^k)$  distributed on  $k$  parties; labels  $\mathbf{y}$  on host party; number of epochs  $T_j$  of each party  $j$ ; number of samples  $n$

- 1:  $\triangleright$  Guest Parties
- 2: **for** party  $j \in [2, k]$  **do**  $\triangleright$  Parallel
- 3:     **Initialize** guest model  $f_{\theta_j}$ , random targets  $C_j$
- 4:     **for** epoch  $t \in [T_j]$  **do**
- 5:          $\tilde{P}^{(t)} \leftarrow \arg \min_{P^{(t)}} \frac{1}{2n} \|f_{\theta_j}^{(t)}(\mathbf{x}^j) - P^{(t)}C_j\|_F^2$
- 6:          $\tilde{L}^{(t)}(\theta_j, \mathbf{x}^j) \leftarrow \frac{1}{2n} \|f_{\theta_j}^{(t)}(\mathbf{x}^j) - \tilde{P}^{(t)}C_j\|_F^2$
- 7:          $\mathbf{g}^{(t)} \leftarrow \nabla_{\theta_j} L^{(t)}(\theta_j, \mathbf{x}^j)$
- 8:         **if** FedOnce-L0 **then**
- 9:              $\theta_j^{(t+1)} = \theta_j^{(t)} - \eta^{(t)}\mathbf{g}^{(t)}$
- 10:         **else if** FedOnce-L1 **then**
- 11:              $\mathbf{g}_{clip}^{(t)} = \mathbf{g}^{(t)} / \max(1, \frac{\|\mathbf{g}^{(t)}\|_2}{\Omega})$
- 12:              $\mathbf{g}_{noise}^{(t)} = \mathbf{g}_{clip}^{(t)} + \mathcal{N}(0, \sigma^2\Omega^2\mathbf{I})$
- 13:              $\theta_j^{(t+1)} = \theta_j^{(t)} - \eta^{(t)}\mathbf{g}_{noise}^{(t)}$
- 14:         **end if**
- 15:         **end for**
- 16:         Send representation  $R_j = f_{\theta_j}(\mathbf{x}^j)$  to host party
- 17:     **end for**
- 18:  $\triangleright$  Host Party
- 19: Collect  $R_j (j \in [2, k])$  from all guest parties
- 20: **for** epoch  $t \in [T_j]$  **do**
- 21:      $L^{(t)}(w, \mathbf{x}^1, R_{2:j}) \leftarrow \frac{1}{2n} \|f_w^{(t)}(\mathbf{x}^1, R_{2:j}) - \mathbf{y}\|_F^2$
- 22:      $\mathbf{g}^{(t)} \leftarrow \nabla_w L^{(t)}(w, \mathbf{x}^1, R_{2:j})$
- 23:     **if** FedOnce-L0 **then**
- 24:          $w^{(t+1)} = w^{(t)} - \eta^{(t)}\mathbf{g}^{(t)}$
- 25:     **else if** FedOnce-L1 **then**
- 26:          $\mathbf{g}_{clip}^{(t)} = \mathbf{g}^{(t)} / \max(1, \frac{\|\mathbf{g}^{(t)}\|_2}{\Omega})$
- 27:          $\mathbf{g}_{noise}^{(t)} = \mathbf{g}_{clip}^{(t)} + \mathcal{N}(0, \sigma^2\Omega^2\mathbf{I})$
- 28:          $w^{(t+1)} = w^{(t)} - \eta^{(t)}\mathbf{g}_{noise}^{(t)}$
- 29:     **end if**
- 30:     **end for**
- 31: **return** model parameters  $\theta_j (j \in [2, k]), w$

---



---

### Algorithm 2 Predicting Process of FedOnce

---

**Require:** testing data  $(\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^k)$  distributed on  $k$  parties; trained models  $f_{\theta_j}(\cdot) (j \in [1, k]), f_w$ ;

- 1:  $\triangleright$  Guest parties
- 2: **for** party  $j \in [2, k]$  **do**
- 3:      $R_j = f_{\theta_j}(\mathbf{x}^j)$
- 4:     Send  $R_j$  to host party
- 5: **end for**
- 6:  $\triangleright$  Host party
- 7:  $\hat{\mathbf{y}} = f_w(\mathbf{x}^1, R_{2:j})$
- 8: **return** prediction  $\hat{\mathbf{y}}$

---

privacy loss in differentially private federated learning also accumulates across parties. Regarding the division of overall privacy budget  $\varepsilon$  to each party, existing studies [24], [25] on differentially private vertical federated learning simply divide  $\varepsilon$  equally to each party. This method, denoted as *simple division*, satisfies differential privacy according to simple composition (Theorem 2), which is over-conservative and leads to a small privacy budget for each party. Hence,

we propose *moments division*, which is a natural extension of moments accountant. Since inner-party privacy loss is accumulated according to moments accountant, the Rényi divergence [29] for each party is calculated along with the privacy loss. Instead of accumulating  $\varepsilon$ , we accumulate the Rényi divergences of parties and calculate the overall  $\varepsilon$  given  $\delta$ . Our main result is summarized as Theorem 4.

**Theorem 4.** Given the sampling probability of SGD  $q$ , the number of epochs  $T_j$  of each party  $\mathcal{P}_j$ , and the number of parties  $k$ , there exists constants  $c_1$  and  $c_2$  so that for  $\forall \varepsilon < c_1 q^2 \sum_{j=1}^k T_j$ ,  $\forall \delta > 0$ , FedOnce-L1 is  $(\varepsilon, \delta)$ -differential privacy if we choose

$$\sigma \geq c_2 \frac{q \sqrt{\log(1/\delta) \sum_{j=1}^k T_j}}{\varepsilon}$$

The proof of Theorem 4 is provided as follows. The major difference between our proof of FedOnce-L1 and moments accountant lies in the accumulation of moments, where moments accountant only accumulates across iterations, but FedOnce-L1 also accumulates across parties. Specifically, we extend Lemma 2 to the setting of FedOnce-L1 as Lemma 4.

**Lemma 4.** Suppose  $\mathcal{M}$  contains  $k$  sequences of Gaussian mechanisms  $\mathcal{M}^1, \dots, \mathcal{M}^k$ , each sequence  $j$  contains  $T$  Gaussian mechanisms  $\mathcal{M}_1^j, \dots, \mathcal{M}_T^j$ . For  $j \in [1, k-1]$ ,  $\mathcal{M}_i^j : \mathcal{D} \rightarrow \mathcal{R}_j$  takes original dataset as input (guest parties). For  $j = k$ ,  $\mathcal{M}_i^k : \prod_{j=1}^{k-1} \mathcal{R}_j \times \mathcal{D} \rightarrow \mathcal{R}_k$  take the original dataset and the output of other mechanisms as input (host party). For  $\forall \lambda > 0$ , we have

$$\alpha_{\mathcal{M}}(\lambda) \leq \sum_{j=1}^k \sum_{i=1}^T \alpha_{\mathcal{M}_i^j}(\lambda)$$

*Proof.* Intuitively, the sequences of mechanisms on all the parties can be regarded as a large sequence of adaptive mechanisms. Hence, the overall moments is the summation of the moments of each Gaussian mechanism by simply applying Lemma 2. Formally, we provide proof of Lemma 4 as follows.

From Lemma 2, in each sequence  $j$ , we have

$$\alpha_{\mathcal{M}^j} \leq \sum_{i=1}^T \alpha_{\mathcal{M}_i^j}(\lambda) \quad (2)$$

For  $j \in [1, k-1]$ , i.e., guest parties, we define a mechanism  $\tilde{\mathcal{M}}_i^j : \prod_{t=1}^{j-1} \mathcal{R}_t \times \mathcal{D} \rightarrow \mathcal{R}_j$  as follows. We use  $o^{1:j}$  to denote  $o^1, \dots, o^j$  for simplicity.

$$\tilde{\mathcal{M}}_i^j(\mathcal{D}, o^{1:j}) = \mathcal{M}_i^j(\mathcal{D})$$

This mechanism  $\tilde{\mathcal{M}}_i^j$  satisfies the condition in Lemma 2 while it has the same output as  $\mathcal{M}_i^j$ . For  $j = k$ , i.e., host party,  $\mathcal{M}_i^k$  has already satisfied the condition in Lemma 2. Thus, by applying Lemma 2 again on all the parties, we have

$$\alpha_{\mathcal{M}}(\lambda) \leq \sum_{j=1}^k \alpha_{\mathcal{M}^j}(\lambda) \quad (3)$$

Combine (2) and (3),

$$\alpha_{\mathcal{M}}(\lambda) \leq \sum_{j=1}^k \sum_{i=1}^T \alpha_{\mathcal{M}_i^j}(\lambda)$$

□

With Lemma 4, the rest of the proof of Theorem 4 largely follows the proof of moments accountant [27] by replacing  $T$  with  $kT$ .

*Proof.* From Lemma 1, since  $\lambda \geq 1$ ,  $q < 1$ , we have

$$\begin{aligned} \alpha_{\mathcal{M}_i^j}(\lambda) &\leq \frac{q^2 \lambda (\lambda + 1)}{(1-q)\sigma^2} \\ &\leq \frac{q^2 \lambda^2}{\sigma^2} \cdot \frac{\lambda + 1}{\lambda(1-q)} \\ &\leq \frac{q^2 \lambda^2}{\sigma^2} \end{aligned}$$

From Lemma 4, the overall maximum  $\lambda$ -th moments can be calculated by a summation.

$$\alpha_{\mathcal{M}}(\lambda) \leq \sum_{j=1}^k \sum_{i=1}^T \alpha_{\mathcal{M}_i^j}(\lambda) \leq \frac{q^2 \lambda^2 \sum_{j=1}^k T_j}{\sigma^2} \quad (4)$$

Given (4), considering Lemma 3, in order to ensure the  $(\varepsilon, \delta)$ -differential privacy of FedOnce-L1, it suffices to prove

$$\frac{q^2 \lambda^2 \sum_{j=1}^k T_j}{\sigma^2} \leq \frac{\lambda \varepsilon}{2} \quad (5)$$

$$\exp\left(-\frac{\lambda \varepsilon}{2}\right) \leq \delta \quad (6)$$

From the condition of Lemma 1, we also require

$$\lambda \leq \sigma^2 \ln \frac{1}{q\sigma} \quad (7)$$

When  $\varepsilon < c_1 q^2 \sum_{j=1}^k T_j$ , it can be easily verified that there exist constants  $c_1$  and  $c_2$  that satisfy conditions (5), (6), (7) if we set

$$\sigma \geq c_2 \frac{q \sqrt{\log(1/\delta) \sum_{j=1}^k T_j}}{\varepsilon}$$

□

With this theorem, the inter-party privacy loss accumulates much slower than simple division. In our experiments, moments division reduces the overall privacy loss by 88.9% when  $k = 100$  compared to simple division.

## 4.2 Privacy of Representations

In this study, assuming the host party is trusted among all the parties, we enable collaborative training by representation sharing to comply with GDPR [6]. Nonetheless, in some applications, this setting may not be true and the representations are potentially threatened by privacy attacks. Such an attack could be similarly designed as the reconstruction attack [34] in horizontal federated learning. This privacy risk indicates merely prohibiting unnecessary data sharing is not sufficient to guarantee privacy. More strict and clear regulations above GDPR are desired to measure the permitted information leakage.

To address this privacy issue, some techniques adopted to protect gradients such as local differential privacy (LDP) [35], homomorphic encryption (HE) [10], and secure multi-party computation (MPC) [11], [12] can potentially assist.

Among these techniques, HE and MPC incur large computational cost, which is a heavy burden for devices with limited computational resources like smart Wi-Fi routers. LDP inserts noise into transferred gradients to guarantee the privacy of each party. Nevertheless, representations in FedOnce, in the format of matrices, cannot be similarly analyzed like gradients. Another approach is to regard the representations as released data and adopt privacy-preserving data release [36] to add noise to the representations. As elaborated in Section 5.8, this method is a possible solution, whereas the noise on representations should be strictly bounded to a small scale; otherwise, FedOnce could suffer considerable performance loss. In conclusion, more advanced methods are desired to protect the privacy representations in this new scenario where the host party cannot be trusted.

## 5 EXPERIMENT

### 5.1 Experiment Setup

**Datasets.** We evaluate FedOnce on eight public datasets and two real-world federated datasets, whose details are summarized in Table 1. **Public datasets:** Public datasets are used to simulate the scenario where each party has the same number of features. *gisette*, *phishing*, and *covtype* are binary classification datasets obtained from LIBSVM<sup>3</sup>. *UJIIndoorLoc* and *Superconduct* are regression datasets obtained from UCI<sup>4</sup>. *MNIST*, *KMNIST* and *Fashion-MNIST* are multi-class classification datasets obtained from TorchVision<sup>5</sup>. **Real-world federated datasets:** 1) *NUS-WIDE* [37] is a real-world multi-view image datasets, used to simulate a security company training a model based on image data collected from different surveillance cameras. In *NUS-WIDE*, each image is stored in the form of five types of low-level features. *NUS-WIDE* provides 81 labels, among which we pick the most occurring label *sky* to conduct binary classification. 2) *MovieLens* [38] is a real-world recommendation dataset, used to simulate a film production company training a recommendation model based on the data from a movie rating website and a movie streaming website. *MovieLens* contains 9992 one-hot identity features and 133 auxiliary features. We set movie ratings as labels and regard the task as regression.

**Feature Division.** To evaluate FedOnce under different feature divisions, we consider two kinds of divisions: 1) **Equal division** for public datasets: each dataset is divided into  $k$  parties equally by features. 2) **Real-world division** for real-world federated datasets: each dataset is naturally distributed on  $k$  parties by features according to real-world settings. Due to the distinct distribution of each feature, we further study the effect on the performance when the number and importance of features are biased in Section 5.5.

In **equal division**, for *gisette*, *phishing*, *covtype*, *UJIIndoorLoc* and *Superconduct* with 1D features, we equally divide the features into  $k$  parties. For *MNIST*, *KMNIST* and *Fashion-MNIST* with 2D image features, since neighboring features are correlated, we split the images with size  $28 \times 28$  both

TABLE 1  
Detailed information for datasets

(a) Public datasets			
Dataset	#samples	#features	Task
gisette	6,000	5,000	binary-classification
covtype	581,012	54	binary-classification
phishing	11,055	68	binary-classification
UJIIndoorLoc	21,048	529	regression
Superconduct	21,263	81	regression
MNIST	60,000	$28 \times 28$	multi-classification
KMNIST	60,000	$28 \times 28$	multi-classification
Fashion-MNIST	60,000	$28 \times 28$	multi-classification
(b) Real-world federated datasets			
Dataset	#samples	#features	Task
NUS-WIDE	269,648	66+144+75+128+225	bin-class
MovieLens	1,000,209	9992+133	reg

horizontally and vertically into four parts with size  $14 \times 14$  and assign each part to a party ( $k = 4$ ). This split on images is commonly adopted to evaluate the performance of vertical federated learning [39], [40]. In **real-world division**, for *NUS-WIDE*, we assign five views of each image into five parties, respectively. For *MovieLens*, we assign 9923 one-hot identity features to one party and assign 133 auxiliary features to the other party.

**Baselines.** To evaluate the communication efficiency of FedOnce-L0, we adopt five baselines: 1) Solo: each party trains locally with its own data and real labels. 2) Combine: all the data and labels are trained centrally. 3) SplitNN [13]: a vertical federated learning algorithm for neural networks. 4) SecureBoost [10]: a vertical federated learning framework for gradient boosting decision trees (GBDT). Notably, **VF<sup>2</sup>Boost [11] and Pivot [12], producing the same accuracy as SecureBoost with additional techniques on encryption, are not individually compared.** Since FedOnce-L0 provides no protection on the released model during the training process, for a fair comparison, we ignore the cryptographic overhead when calculating the communication size of SecureBoost. 5) Linear-Combine: a linear model is trained on centralized datasets. Specifically, we train logistic regression for classification tasks and train ridge regression for regression tasks on centralized datasets similarly to Combine. Among these baselines, SecureBoost is not evaluated on *MNIST*, *KMNIST* and *Fashion-MNIST* since GBDT is unsuitable for 2D features.

To evaluate FedOnce-L1, since the analyses of existing studies [24], [25], [26] are based on different models which cannot be directly applied to our algorithm based on neural networks. For a fair comparison, we extend existing analyses to a baseline approach, named *Priv-Baseline*, by only adopting their common model-agnostic inter-party analysis. The privacy analysis of *Priv-Baseline* and FedOnce-L1 are both based on moments accountant [27].

**Hardware.** We conduct the experiments on a machine with two AMD EPYC 7543 32-Core CPUs, four NVIDIA A100 GPUs, and 504GB memory. We implement FedOnce in Python 3.8 and adopt the *pytorch-dp*<sup>6</sup> (now *opacus*) library

3. <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>

4. <https://archive.ics.uci.edu/ml/datasets.php>

5. <https://pytorch.org/docs/stable/torchvision/datasets.html>

6. <https://github.com/facebookresearch/pytorch-dp>

in FedOnce-L1.

**Models.** The models used for FedOnce-L0 and FedOnce-L1 on each dataset are summarized in Table 2 and Table 3, respectively. The sizes of hidden layers are carefully tuned for each dataset.  $FC(m \times n)$  indicates fully connected layers with two hidden layers which have  $m$  and  $n$  nodes, respectively. CNN0 and CNN1 indicates two kinds of convolutional neural networks whose structures are shown in Fig. 4. NCF( $m \times n$ ) indicates neural collaborative filtering [41] with two hidden layers which have  $m$  and  $n$  nodes, respectively.

TABLE 2  
Models used in FedOnce-L0

Dataset	Guest Model	Host Model
gisette	FC(100×100×50)	FC(30×30)
covtype	FC(200×100×100)	FC(30×30)
phishing	FC(30×30)	FC(30)
UJIIndoorLoc	FC(50×50×50×30)	FC(30×10)
Superconduct	FC(30×30)	FC(30×10)
MNIST	CNN0→FC(128)	FC(128)
KMNIST	CNN0→FC(128)	FC(128)
Fashion-MNIST	CNN0→FC(128)	FC(128)
NUS-WIDE	FC(60)	FC(50)
MovieLens	NCF(32×16)/NCF(128)	FC(10)

TABLE 3  
Models used in FedOnce-L1

Dataset	Guest Model	Host Model
gisette	FC(30)	FC(10)
covtype	FC(50)	FC(30)
phishing	FC(30)	FC(10)
UJIIndoorLoc	FC(50×50)	FC(20×20)
Superconduct	FC(50×50)	FC(20×20)
MNIST	CNN1	FC(10)

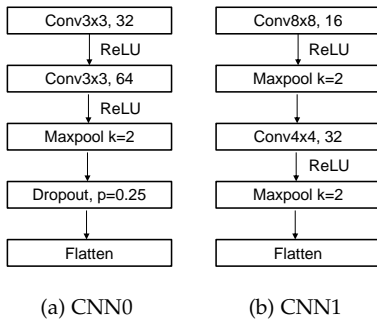


Fig. 4. Structure of CNN in experiments

**Training.** For *MNIST*, *KMNIST*, *Fashion-MNIST* and *NUS-WIDE*, public test set is used for testing. In other datasets, training and test set are split by 9:1. Five-fold cross-validation is performed for hyperparameter tuning and the test performance is presented. We report accuracy for classification tasks and root mean square error (RMSE) between the real labels and the predicted labels for regression tasks.

In Section 5.2, we present the performance of FedOnce-L0 when labels are assigned to each party. In Section 5.3, the communication efficiency of FedOnce-L0 is evaluated. In

Section 5.4, the scalability of FedOnce is evaluated on a high-dimensional dataset. In Section 5.5, we study the performance of FedOnce on biased feature split. In Section 5.6, we study the effect of different unsupervised learning methods on FedOnce. In Section 5.7, we analyze the privacy loss and present the performance of FedOnce-L1 under different privacy budget  $\epsilon$ . In Section 5.8, we demonstrate the impact of noise adding to representations on the performance.

## 5.2 Performance of Different Host Parties

In this subsection, we evaluate the performance of FedOnce-L0 when each party is chosen as the host party. For a fair comparison, we include two additional baselines: a) SplitNN-Comm: SplitNN terminates training when it costs the same amount of communication size as FedOnce-L0. b) SecureBoost-Comm: SecureBoost terminates training when it costs the same amount of communication size as FedOnce-L0. For FedOnce-L0, Solo, and SecureBoost(-Comm),  $j$  indicates the party with labels; for SplitNN(-Comm) and Combine,  $j$  is meaningless because each party is treated equally in these settings. The results are displayed in Fig. 5. Since failing to finish the first iteration under the same communication size as FedOnce, SecureBoost-Comm is missing in *gisette* and *phishing*.

Three observations can be made from Fig. 5. First, FedOnce-L0 significantly outperforms Solo, SplitNN-Comm, and SecureBoost-Comm on each party, indicating a huge performance-boosting against SplitNN and SecureBoost under the same communication size. Second, for each party  $j$ , FedOnce-L0 achieves close performance compared to SplitNN, SecureBoost, and even Combine. For example, FedOnce-L0 only incurs 0.3% accuracy loss compared to Combine on *MNIST*, and incurs 0.09% accuracy loss compared to Combine on *phishing*. Third, except on simple datasets like *gisette* and *phishing*, centralized linear models are significantly outperformed by FedOnce and most other baselines. This indicates that vertical federated learning methods with these linear models, which produce even poorer performance, are not competitive in most cases.

**Hyperparameters** For each dataset, we summarize the hyperparameters of FedOnce-L0 in Table 5 and hyperparameters of FedOnce-L1 in Appendix A. In the table,  $\eta$  refers to the learning rate,  $\lambda$  refers to weight decay,  $b$  refers to batch size,  $T$  refers to the number of epochs,  $d$  refers to the dimension of representations.  $f$  indicates the frequency of permutation matrix  $P$  to be updated. For example, if the update frequency is 3,  $P$  will be updated every three epochs.

Particularly, we study the effect of an important hyperparameter, the dimension of random targets  $d$ , in Fig. 7 in which the mean and standard variance of the performance when each party is chosen as host party is reported. Intuitively,  $d$  reflects the amount of information that can be transferred between the host party and guest parties. A small  $d$  causes a less representative  $R_j$ , leading to less contribution of the data on guest parties to the aggregation model. A large  $d$  causes a high-dimensional searching space of representations, among which it might be infeasible to locate the optimal  $R_j$ . Therefore, a moderate  $d$  leads to the best performance of FedOnce-L0. From Fig. 7, we observe that FedOnce-L0 achieves the best performance when  $d \approx 16$

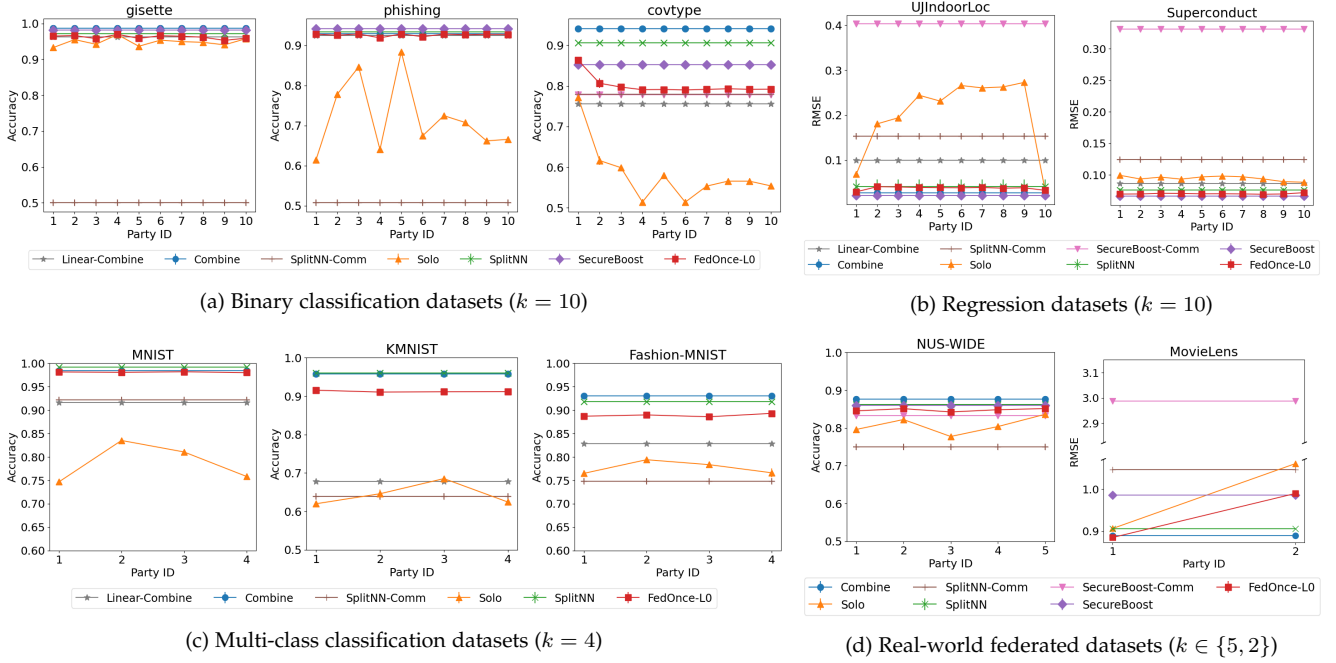


Fig. 5. Performance of FedOnce-L0 (i.e., with one communication round)

on *MNIST*. The choices of  $d$  on other datasets are similarly made.

TABLE 4  
Training time of FedOnce-L0 and SplitNN

Dataset	Training Time(min)	
	SplitNN	FedOnce-L0
gisette	0.62	0.60
covtype	171.12	52.33
phishing	4.07	1.63
UJIndoorLoc	23.15	8.96
Superconduct	16.55	7.54
MNIST	53.58	24.17
KMNIST	56.47	24.13
Fashion-MNIST	54.58	23.00
NUS-WIDE	8.86	5.17
MovieLens	39.88	34.91

**Training Time.** For each dataset, we record the training time of SplitNN and FedOnce-L0 in Table 4. As observed from the table, FedOnce-L0 generally has a lower training time than SplitNN; this is because the local models in FedOnce can be trained in parallel, while all the local models are connected with the aggregation model and trained in each iteration in SplitNN. We also highlight that the communication time is measured in a shared-memory environment of a single machine. In reality, the communication time can also be a huge burden for SplitNN which requires much more communication size as demonstrated in Fig. 5.

### 5.3 Communication Efficiency

In this subsection, we evaluate the communication efficiency of FedOnce-L0. Without loss of generality, we set party  $\mathcal{P}_1$  as the host party and report the mean and standard variance of performance across five folds. After completing one-shot

communication, we continue training the models obtained by FedOnce-L0 and track the mean and standard variance of performance across five folds in each iteration. The results are presented in Fig. 6.

Three observations can be made from the results. First, FedOnce-L0 (the first data point of multi-round FedOnce-L0) significantly outperforms SplitNN and SecureBoost at a small communication size. Second, only under considerable communication size, can SplitNN and SecureBoost achieve the same performance as FedOnce-L0. For example, on *KMNIST* dataset, SplitNN requires a 17 $\times$  communication size to achieve the same performance as FedOnce-L0. Third, multi-round FedOnce-L0 consistently outperforms SplitNN and SecureBoost as communication size increases.

### 5.4 Scalability

In this subsection, we evaluate the scalability of FedOnce on a high-dimensional dataset *gisette* with 5000 features. The performance of FedOnce-L0 and baselines when  $k$  ranges from [10, 200] is displayed in Fig. 8. The error bars indicate the standard variance of performance when labels are assigned to different parties. SecureBoost-Comm is omitted since SecureBoost cannot finish the first iteration with the same communication size as FedOnce.

From Fig. 8, we can make two observations. First, as the number of parties  $k$  increases, FedOnce-L0, whose performance remains stable while the performance of Solo degrades rapidly, is scalable. Second, even at a large number of parties, FedOnce consistently outperforms SecureBoost and SplitNN with the same communication size.

### 5.5 Performance on Biased Datasets

In this subsection, we study how the quality of features affects the performance of FedOnce. For convenience, the experiment is conducted under a two-party setting due to the

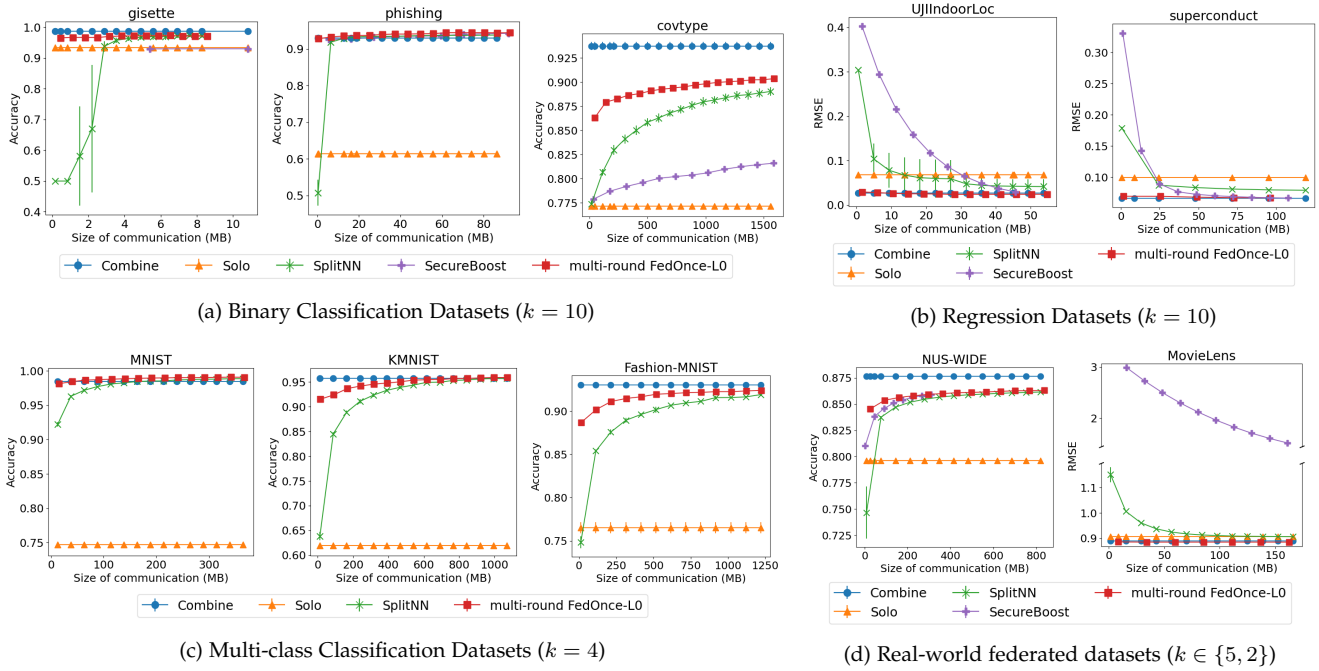


Fig. 6. Performance of multi-round FedOnce-L0 and FedOnce-L0 (i.e., the first data point of multi-round FedOnce-L0)

TABLE 5  
Hyperparameters of FedOnce-L0 on each dataset

Dataset	Host Party <sup>1</sup>	Guest Model				Host Model				$d$	$f$	Optimizer
		$\eta$	$\lambda$	$b$	$T$	$\eta$	$\lambda$	$b$	$T$			
gisette	*	1e-4	1e-5	100	100	1e-4	1e-4	100	100	3	1	Adam
covtype	*	1e-3	1e-5	100	30	1e-3	1e-3	100	100	3	1	Adam
phishing	*	1e-4	1e-5	100	100	1e-4	1e-4	100	300	3	1	Adam
UJIIndoorLoc	*	3e-4	1e-4	100	500	1e-4	1e-4	100	600	3	1	Adam
Superconduct	*	3e-4	1e-4	100	500	3e-4	1e-4	128	600	3	1	Adam
MNIST	*	1e-4	1e-5	128	100	1e-3	1e-5	128	200	16	3	Adam
KMNIST	*	1e-4	1e-5	128	100	1e-3	1e-5	128	200	16	3	Adam
Fashion-MNIST	*	1e-4	1e-5	128	100	1e-3	1e-5	128	200	16	3	Adam
NUS-WIDE	*	3e-4	1e-5	128	50	7e-5	1e-5	128	50	8	1	Adam
MovieLens	0	1e-4	1e-5	64	30	1e-4	2e-4	64	40	3	1	Adam
	1	1e-4	1e-5	64	30	5e-4	1e-5	64	100	48	1	Adam

<sup>1</sup>The party that is set as host party.  $j$ : the hyperparameters when party  $j$  is set as host party. \* means the hyperparameters remain the same when each party is set as host party.

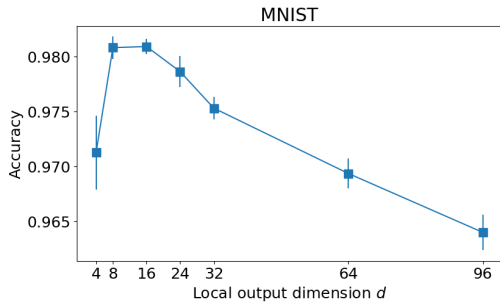


Fig. 7. Performance of FedOnce-L0 with different  $d$

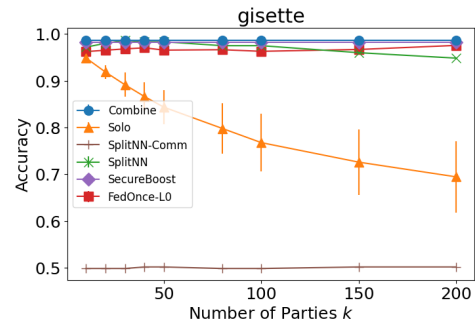


Fig. 8. Performance of FedOnce-L0 (i.e., with one communication round) when  $k$  scales from 10 to 200

constraint of our designed metric of bias. First, we leveraged XGBoost<sup>7</sup> library to calculate the importance of each feature. Then, the top- $p$  important features are marked as *dominant*

*features*. We randomly assign a ratio  $\alpha$  of dominant features to the host party and assign the remaining  $1 - \alpha$  of dominant features to the guest party. The other features are randomly

7. <https://github.com/dmlc/xgboost>

TABLE 6  
Performance of FedOnce-L1 with different  $\varepsilon$

(a) Performance on gisette, $\delta = 10^{-4}$						(b) Performance on covtype, $\delta = 10^{-5}$							
Algorithm	$\varepsilon$	Accuracy					Algorithm	$\varepsilon$	Accuracy				
		Party 1	Party 2	Party 3	Party 4	Range			Party 1	Party 2	Party 3	Party 4	Range
Priv-Baseline	2	86.20%	86.09%	75.14%	73.00%	73.00%~86.20%	Priv-Baseline	2	74.60%	65.88%	60.19%	66.53%	60.19%~74.60%
	4	89.31%	88.10%	80.67%	80.57%	80.57%~89.31%		4	75.89%	69.14%	69.88%	68.70%	68.70%~75.89%
	6	89.49%	88.46%	84.01%	82.08%	82.08%~89.49%		6	78.27%	73.75%	73.46%	73.87%	73.46%~78.27%
	8	89.71%	88.42%	83.13%	82.55%	82.55%~89.71%		8	79.25%	73.96%	73.40%	74.04%	73.40%~79.25%
FedOnce-L1	2	88.81%	88.05%	78.88%	79.51%	78.88%~88.81%	FedOnce-L1	2	76.88%	70.74%	70.11%	71.56%	70.11%~76.88%
	4	89.73%	88.54%	81.77%	82.88%	81.77%~89.73%		4	77.80%	73.33%	73.71%	72.15%	72.15%~77.80%
	6	89.25%	88.67%	83.53%	82.77%	82.77%~89.25%		6	79.23%	74.18%	73.66%	73.98%	73.98%~79.23%
	8	93.37%	93.01%	93.14%	93.64%	93.01%~93.64%		8	80.49%	75.04%	74.74%	74.80%	74.74%~80.49%
$\infty$	96.91%	97.11%	97.61%	96.42%	96.42%~97.61%	$\infty$	82.24%	74.48%	74.29%	74.06%	74.06%~82.24%		

(c) Performance on phishing, $\delta = 10^{-5}$						(d) Performance on UJIIndoorLoc, $\delta = 10^{-5}$							
Algorithm	$\varepsilon$	Accuracy					Algorithm	$\varepsilon$	RMSE				
		Party 1	Party 2	Party 3	Party 4	Range			Party 1	Party 2	Party 3	Party 4	Range
Priv-Baseline	2	87.69%	87.22%	83.98%	83.64%	83.64%~87.69%	Priv-Baseline	2	0.2326	0.2275	0.2292	0.2501	0.2275-0.2501
	4	89.43%	87.49%	88.79%	87.92%	87.49%~88.79%		4	0.1487	0.1578	0.2062	0.1943	0.1487-0.2062
	6	90.31%	89.51%	89.55%	89.29%	89.29%~90.31%		6	0.1053	0.1224	0.1475	0.1914	0.1053-0.1914
	8	90.89%	90.76%	89.69%	90.51%	89.69%~90.89%		8	0.1051	0.1198	0.1409	0.1345	0.1051-0.1345
FedOnce-L1	2	90.10%	89.36%	87.78%	88.95%	87.78%~90.10%	FedOnce-L1	2	0.1066	0.1102	0.1087	0.1489	0.1066-0.1489
	4	91.10%	90.24%	90.90%	91.16%	90.24%~91.16%		4	0.0769	0.0811	0.0843	0.1362	0.0769-0.1362
	6	91.56%	90.65%	90.36%	90.16%	90.36%~91.56%		6	0.0714	0.0746	0.0664	0.1223	0.0664-0.1223
	8	90.78%	91.12%	90.63%	91.30%	90.63%~91.30%		8	0.0542	0.0614	0.0678	0.0730	0.0542-0.0730
$\infty$	91.95%	91.28%	91.21%	90.90%	90.90%~91.95%	$\infty$	0.0355	0.0369	0.0402	0.0359	0.0355-0.0402		

(e) Performance on MNIST, $\delta = 10^{-5}$						(f) Performance on Superconduct, $\delta = 10^{-5}$							
Algorithm	$\varepsilon$	Accuracy					Algorithm	$\varepsilon$	RMSE				
		Party 1	Party 2	Party 3	Party 4	Range			Party 1	Party 2	Party 3	Party 4	Range
Priv-Baseline	2	45.29%	39.98%	50.28%	44.49%	39.98%~50.28%	Priv-Baseline	2	0.1668	0.1667	0.1783	0.1373	0.1373-0.1783
	4	70.07%	76.70%	78.00%	75.23%	70.07%~78.00%		4	0.1253	0.1387	0.1348	0.1246	0.1246-0.1387
	6	83.64%	82.92%	84.00%	82.60%	82.60%~84.00%		6	0.1173	0.1256	0.1299	0.1230	0.1173-0.1299
	8	85.40%	86.33%	86.21%	85.28%	85.28%~86.33%		8	0.1125	0.1140	0.1308	0.1148	0.1125-0.1308
FedOnce-L1	2	82.13%	82.57%	83.30%	81.43%	81.43%~83.30%	FedOnce-L1	2	0.1080	0.1043	0.1059	0.1110	0.1043-0.1110
	4	87.01%	87.51%	87.91%	87.02%	87.01%~87.91%		4	0.1071	0.1003	0.1009	0.1052	0.1003-0.1052
	6	88.05%	89.03%	88.24%	87.91%	87.91%~89.03%		6	0.0985	0.1001	0.0984	0.0973	0.0973-0.1001
	8	90.52%	90.63%	89.61%	89.32%	89.32%~90.63%		8	0.1020	0.1075	0.0985	0.0956	0.0956-0.1075
$\infty$	92.05%	92.97%	93.11%	91.74%	91.74%~93.11%	$\infty$	0.0832	0.0833	0.0841	0.0801	0.0801-0.0841		

assigned to two parties to ensure both parties hold the same number of features. Finally, we adjust the ratio  $\alpha \in [0, 1]$  to control the feature quality of the host party. The host party holds the best features when  $\alpha = 1$  while holding the worst features when  $\alpha = 0$ . We conduct experiments on *phishing* (classification) and *UJIIndoorLoc* (regression). Under five-fold cross-validation, the mean performance of the host party as well as the standard variance across five folds are reported in Fig. 9.

Two observations can be made from Fig. 9. First, FedOnce achieves the most significant performance when the quality of features is poor (i.e.,  $\alpha \rightarrow 0$ ). Second, when the feature quality of the host party is good enough (i.e.,  $\alpha \rightarrow 1$ ), federated learning is no longer necessary since even Solo provides competitive performance compared to Combine.

### 5.6 Different Unsupervised Learning Methods

FedOnce is suitable for different unsupervised learning methods. Among these unsupervised learning methods, we compare the performance of FedOnce with NAT [18]

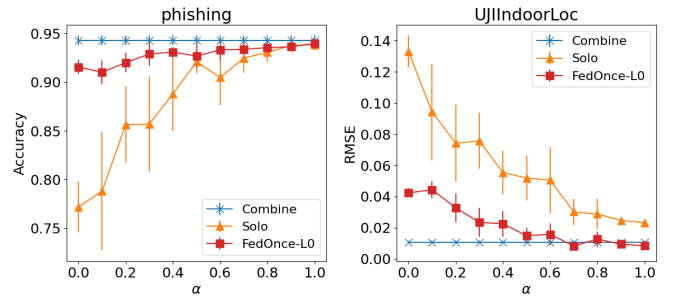


Fig. 9. Performance on biased datasets

(FedOnce-L0) and Principal Component Analysis [42] (FedOnce-PCA). Specifically, fixing both the dimensions of representative features and the number of principal components the same, we present the results of four typical datasets covering image features and multi-variant features in Table 7. Without loss of generality, party  $\mathcal{P}_1$  is selected as the host party. The choice of  $K$  is identical to that in



## Section 5.2.

TABLE 7  
Performance of FedOnce on different unsupervised learning methods

Feature Type	Dataset	Accuracy	
		FedOnce-L0	FedOnce-PCA
Multi-variant	gisetette	96.51%	94.10%
	phishing	92.78%	92.41%
Image	MNIST	98.05%	25.98%
	KMNIST	91.73%	25.59%

From Table 7, we observe that FedOnce-PCA can achieve close performance to FedOnce-L0 on multi-variant datasets, but has very poor performance on image datasets. This is because PCA is incapable to extract useful representations from complex features like images. On the contrary, NAT is a suitable unsupervised learning method for FedOnce that can handle different types of features. Generally, more “advanced” unsupervised learning algorithms tend to produce a better performance on FedOnce. The performance of unsupervised learning methods can be evaluated by commonly used metrics in the literature. For example, NAT is evaluated by the performance of a linear classifier on the learned representations.

### 5.7 Differential Privacy

In this subsection, we first compare the privacy loss of FedOnce-L1 and simple division by theoretical analysis, then present how the moments division is superior to simple division on performance.

**Privacy Loss Analysis.** Fixed the hyperparameters of each party including the number of samples  $n$ , batch size  $b$ , the number of epochs  $T$ , Gaussian noise multiplier  $\sigma$  and overall  $\delta$ , we increase the number of parties  $k$  and calculate the privacy loss of Priv-Baseline and FedOnce-L1. The results are presented in Fig. 10, from which we observe that the privacy loss of FedOnce-L1 increases much slower than Priv-Baseline when  $k > 1$ . Hence, in FedOnce-L1, each party receives a higher privacy budget, thus incurring less performance loss.

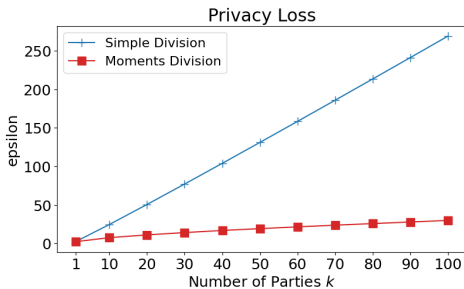


Fig. 10. Privacy loss accumulation when  $n = 60000$ ,  $b = 128$ ,  $T = 50$ ,  $\sigma = 1$ ,  $\delta = 10^{-5}$

**Performance.** Despite the much lower privacy loss compared to simple division, FedOnce-L1 still suffers significant performance loss at a large  $k$  like other vertical federated learning algorithms. Therefore, we set  $k = 4$  and report the performance under different overall privacy budgets  $\epsilon$

on six public datasets. *Fashion-MNIST* and *KMNIST*, which fail to produce reasonable accuracy, are not included in this experiment. The results are summarized in Table 6.

Two observations can be made from Table 6. First, FedOnce-L1 significantly outperforms Priv-Baseline with the same privacy budget  $\epsilon$ . Second, compared with FedOnce-L1 with  $\epsilon = \infty$  (no noise is added), FedOnce-L1 has a close performance under a modest  $\epsilon$  in most datasets. Additionally, the performance of FedOnce-L1 with a small  $\epsilon$  (e.g.,  $\epsilon = 2$ ) can be significantly affected by the noise, implying more advanced privacy mechanisms are desired.

### 5.8 Noise on Representations

In this experiment, to preserve the privacy of representations, we investigate the effect of adding noise to representations on the performance of FedOnce-L0. Specifically, independent Gaussian noise of scale  $\sigma_r$  is added to all the representations. We report the average performance of all the parties on two real-world federated datasets in Table 8.

TABLE 8  
Performance of FedOnce-L0 with Gaussian noise of scale  $\sigma_r$  on representations (bold values mean outperforming Solo)

Dataset	Accuracy/RMSE					Solo
	FedOnce-L0 w/ different $\sigma_r$					
	0.0	0.1	0.5	1.0	1.5	
NUS-WIDE	<b>84.76%</b>	<b>84.64%</b>	<b>82.75%</b>	<b>81.45%</b>	81.05%	80.71%
MovieLens	<b>0.9373</b>	<b>0.9578</b>	<b>0.9690</b>	<b>0.9679</b>	0.9969	0.9835

As can be observed from Table 8, the performance of FedOnce-L0 drops below Solo at a relatively large scale of noise (e.g.,  $\sigma_r = 1.5$ ) on both datasets. This observation indicates that techniques like privacy-preserving data releasing can potentially be used to protect representations, whereas the added noise must be restricted to a small scale (e.g.,  $\sigma_r < 1.0$ ).

## 6 RELATED WORK

In this section, we review the literature in three aspects and summarize the related work in Table 9.

**Vertical Federated Learning.** Except for a study [32] that investigates the scenario where labels are shared across all the parties, most existing studies in vertical federated learning focus on another more practical scenario where only one party holds the labels. Some tree-based approaches, including *SecureBoost* [10], *Pivot* [12], and *VF<sup>2</sup>Boost* [11], are proposed to enable vertical federated learning in gradient boosting decision trees [17]. Besides, federated random forest is also studied in [43]. *SplitNN* [13] is proposed to collaboratively train neural networks by splitting the model among parties and exchanging gradients and outputs in each epoch. These approaches that require multi-round communication suffer large communication cost, whereas our proposed FedOnce only requires one-round communication, leading to much smaller communication cost. Besides trees and neural networks, there are also studies on other machine learning algorithms such as linear regression [44] and logistic regression [45], [46]. However, these algorithms are usually incapable to handle complex tasks.



TABLE 9  
Related work in vertical federated learning

Algorithm	Label Owner <sup>1</sup>	Privacy <sup>2</sup>	Model <sup>3</sup>	Comm. Rounds <sup>4</sup>
FDML [32]	All	N/A	NN	Multi
SplitNN [13]	Single	N/A	NN	Multi
SecureBoost [10]	Single	HE	GBDT	Multi
Pivot [12] <sup>5</sup>	Single	MPC/HE	GBDT/LR	Multi
VF <sup>2</sup> Boost [11] <sup>5</sup>	Single	HE	GBDT	Multi
[24], [26]	Single	DP-S	FW	Multi
[25]	Single	DP-S	LR/HTL	Multi
<b>FedOnce</b>	<b>Single</b>	<b>DP-M</b>	<b>NN</b>	<b>Single</b>

<sup>1</sup> The number of parties that own the labels. Single: Only one party owns the labels, All: all parties own the labels

<sup>2</sup> Privacy mechanism used in the algorithm. N/A: Not specified in details, HE: homomorphic encryption, MPC: secure multi-party computation, DP-S: differential privacy with simple division among parties, DP-M: differential privacy with moments division among parties;

<sup>3</sup> Model supported by the algorithm. NN: neural networks, GBDT: gradient boosting decision trees, LR: logistic regression, FW: Floyd–Warshall algorithm, HTL: hypothesis transfer learning.

<sup>4</sup> Communication rounds required for the algorithm. Multi: multiple rounds are required, Single: single round is required.

<sup>5</sup> Pivot [12] and VF<sup>2</sup>Boost [11] have the *same* accuracy as SecureBoost with additional techniques on encryption.

**Privacy in Federated Learning.** Federated learning confronts two major privacy threats. First, during the training, intermediate results (e.g., gradients) could be vulnerable to backdoor attack [47] or reconstruction attack [34]. This privacy risk can be addressed by multiple methods including homomorphic encryption [10], secure multi-party computation [11], [12], and local differential privacy [35]. Second, after the training, the released model could be exposed to the membership inference attack [19], which can be defended by ensuring differential privacy of the released model (namely *global differential privacy*). This defense is different from local differential privacy since it protects against outside attackers instead of malicious aggregators.

Nevertheless, existing studies on global differential privacy under vertical federated learning suffer significant performance loss. One study [48] achieves differential privacy by *objective perturbation*, which is often intractable in practice compared to *gradient perturbation* according to [49]. The other two studies [24], [25] apply differential privacy based on gradient perturbation, but they both focus on inner-party privacy instead of inter-party privacy. Specifically, in both studies, the simple composition is directly applied in the analysis of privacy loss across parties, leading to excessive noise. In this paper, we apply moments accountant to analyze the inter-party privacy loss, thus effectively reducing the overall privacy loss compared to simple composition.

**Communication-Efficient Federated Learning.** Most existing studies in communication-efficient federated learning [14], [15], [50], [51], [52] focus on horizontal federated learning. These approaches, requiring each party to train independently, cannot be applied to vertical federated learning where only one party holds the labels. Though some approaches [10], [40], [53] study the communication efficiency in vertical federated learning, they all require multiple communication rounds and a certain level of synchronization. As observed from Table 9, vertical federated learning with one-shot communication remains unexplored.

## 7 CONCLUSION

In this paper, we propose FedOnce, a novel one-shot algorithm for vertical federated learning. Moreover, we develop a privacy mechanism for FedOnce (FedOnce-L1) under the

notion of differential privacy and refine inter-party privacy loss. Our experiments demonstrate that FedOnce-L0 achieves impressive performance compared to state-of-the-art vertical federated learning algorithms with only one-shot communication. Besides, FedOnce-L1 significantly outperforms the baseline in our evaluation.

## ACKNOWLEDGMENTS

This research is supported by the National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No: AISG2-RP-2020-018). Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not reflect the views of National Research Foundation, Singapore. Qinbin is also in part supported by a Google PhD Fellowship.

## REFERENCES

- [1] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings *et al.*, “Advances and open problems in federated learning,” *arXiv preprint*, 2019.
- [2] Q. Li, Z. Wen, Z. Wu, S. Hu, N. Wang, and B. He, “Federated learning systems: vision, hype and reality for data privacy and protection,” *TKDE*, 2021.
- [3] L. Zhang, T. Zhu, P. Xiong, W. Zhou, and P. Yu, “A robust game-theoretical federated learning framework with joint differential privacy,” *TKDE*, 2022.
- [4] P. Zhou, K. Wang, L. Guo, S. Gong, and B. Zheng, “A privacy-preserving distributed contextual federated online learning framework with big data support in social recommender systems,” *TKDE*, 2019.
- [5] Q. Yang, Y. Liu, T. Chen, and Y. Tong, “Federated machine learning: Concept and applications,” *TIST*, vol. 10, no. 2, 2019.
- [6] P. Voigt and A. Von dem Bussche, “The eu general data protection regulation (gdpr),” *A Practical Guide, 1st Ed.*, Cham: Springer International Publishing, vol. 10, p. 3152676, 2017.
- [7] T. Yang, G. Andrew, H. Eichner, H. Sun, W. Li, N. Kong, D. Ramage, and F. Beaufays, “Applied federated learning: Improving google keyboard query suggestions,” *arXiv preprint arXiv:1812.02903*, 2018.
- [8] C. P. Friedman, A. K. Wong, and D. Blumenthal, “Achieving a nationwide learning health system,” *Science translational medicine*, 2010.
- [9] D. Loghin, S. Cai, G. Chen, T. T. A. Dinh, F. Fan, Q. Lin, J. Ng, B. C. Ooi, X. Sun, Q.-T. Ta *et al.*, “The disruptions of 5g on data-driven technologies and applications,” *TKDE*, 2020.

- [10] K. Cheng, T. Fan, Y. Jin, Y. Liu, T. Chen, and Q. Yang, "Secureboost: A lossless federated learning framework," *arXiv preprint arXiv:1901.08755*, 2019.
- [11] F. Fu, Y. Shao, L. Yu, J. Jiang, H. Xue, Y. Tao, and B. Cui, "Vf2boost: Very fast vertical federated gradient boosting for cross-enterprise learning," ser. SIGMOD, 2021.
- [12] Y. Wu, S. Cai, X. Xiao, G. Chen, and B. C. Ooi, "Privacy preserving vertical federated learning for tree-based models," *VLDB*, 2020.
- [13] P. Vepakomma, O. Gupta, T. Swedish, and R. Raskar, "Split learning for health: Distributed deep learning without sharing raw patient data," *ICLR 2019 Workshop on AI for social good*, 2019.
- [14] N. Guha, A. Talwalkar, and V. Smith, "One-shot federated learning," *arXiv preprint arXiv:1902.11175*, 2019.
- [15] Y. Zhou, G. Pu, X. Ma, X. Li, and D. Wu, "Distilled one-shot federated learning," *arXiv preprint*, 2020.
- [16] A. Sharifnassab, S. Salehkaleybar, and S. J. Golestani, "Order optimal one-shot federated learning for non-convex loss functions," *arXiv preprint arXiv:2108.08677*, 2021.
- [17] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *KDD*, 2016.
- [18] P. Bojanowski and A. Joulin, "Unsupervised learning by predicting noise," in *ICML*, 2017.
- [19] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, "Membership inference attacks against machine learning models," *SP*, 2017.
- [20] C. Dwork, A. Roth *et al.*, "The algorithmic foundations of differential privacy." *FTTCS*, 2014.
- [21] L. Zhao, L. Ni, S. Hu, Y. Chen, P. Zhou, F. Xiao, and L. Wu, "Inprivate digging: Enabling tree-based distributed data mining with differential privacy," in *INFOCOM*. IEEE, 2018.
- [22] B. Jayaraman, L. Wang, D. Evans, and Q. Gu, "Distributed learning without distress: Privacy-preserving empirical risk minimization," in *NeurIPS*, 2018.
- [23] R. Liu, Y. Cao, H. Chen, R. Guo, and M. Yoshikawa, "Flame: Differentially private federated learning in the shuffle model," in *AAAI*, 2020.
- [24] J. Lou and Y.-m. Cheung, "Uplink communication efficient differentially private sparse optimization with feature-wise distributed data," in *AAAI*, 2018.
- [25] Q. Yao, X. Guo, J. T. Kwok, W.-W. Tu *et al.*, "Privacy-preserving stacking with application to cross-organizational diabetes prediction." in *IJCAI*, 2019.
- [26] J. Lou and Y.-m. Cheung, "An uplink communication-efficient approach to featurewise distributed sparse optimization with differential privacy," *IEEE Trans Neural Netw Learn Syst*, 2020.
- [27] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep learning with differential privacy," in *CCS*, 2016.
- [28] C. Dwork, "Differential privacy: A survey of results," in *TAMC*. Springer, 2008.
- [29] I. Mironov, "Rényi differential privacy," in *CSF*. IEEE, 2017.
- [30] H. W. Kuhn, "The hungarian method for the assignment problem," *Nav. Res. Logist. Q.*, 1955.
- [31] A. Vidanage, T. Ranbaduge, P. Christen, and R. Schnell, "Efficient pattern mining based cryptanalysis for privacy-preserving record linkage," in *ICDE*. IEEE, 2019.
- [32] Y. Hu, D. Niu, J. Yang, and S. Zhou, "Fdml: A collaborative machine learning framework for distributed features," in *KDD*, 2019.
- [33] M. Nasr, R. Shokri, and A. Houmansadr, "Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning," in *SP*, 2019.
- [34] J. Geiping, H. Bauermeister, H. Dröge, and M. Moeller, "Inverting gradients—how easy is it to break privacy in federated learning?" *arXiv preprint*, 2020.
- [35] P. Kairouz, S. Oh, and P. Viswanath, "Extremal mechanisms for local differential privacy," *Advances in neural information processing systems*, vol. 27, pp. 2879–2887, 2014.
- [36] I. Kotsogiannis, Y. Tao, X. He, M. Fanaeepour, A. Machanavajhala, M. Hay, and G. Miklau, "Privatesql: a differentially private sql query engine," *VLDB*, 2019.
- [37] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y.-T. Zheng, "Nus-wide: A real-world web image database from national university of singapore," in *CIVR*, 2009.
- [38] F. M. Harper and J. A. Konstan, "The movielens datasets: History and context," *TIIS*, 2015.
- [39] Y. Liu, X. Zhang, and L. Wang, "Asymmetrical vertical federated learning," *arXiv preprint arXiv:2004.07427*, 2020.
- [40] T. Chen, X. Jin, Y. Sun, and W. Yin, "Vaf1: a method of vertical asynchronous federated learning," *arXiv preprint arXiv:2007.06081*, 2020.
- [41] X. He, L. Liao, H. Zhang, L. Nie, X. Hu, and T.-S. Chua, "Neural collaborative filtering," in *WWW*, 2017.
- [42] L. Xu, "Theories for unsupervised learning: Pca and its nonlinear extensions," in *Proceedings of 1994 IEEE International Conference on Neural Networks (ICNN'94)*. IEEE, 1994.
- [43] Y. Liu, Y. Liu, Z. Liu, Y. Liang, C. Meng, J. Zhang, and Y. Zheng, "Federated forest," *IEEE Transactions on Big Data*, 2020.
- [44] Q. Zhang, B. Gu, C. Deng, and H. Huang, "Secure bilevel asynchronous vertical federated learning with backward updating," in *AAAI*, 2021.
- [45] Y. Hu, P. Liu, L. Kong, and D. Niu, "Learning privately over distributed features: An admm sharing approach," *arXiv preprint arXiv:1907.07735*, 2019.
- [46] Y. Liu, Y. Kang, X. Zhang, L. Li, Y. Cheng, T. Chen, M. Hong, and Q. Yang, "A communication efficient collaborative learning framework for distributed features," *arXiv preprint arXiv:1912.11187*, 2020.
- [47] H. Wang, K. Sreenivasan, S. Rajput, H. Vishwakarma, S. Agarwal, J.-y. Sohn, K. Lee, and D. Papailiopoulos, "Attack of the tails: Yes, you really can backdoor federated learning," *NeurIPS 2020*, 2020.
- [48] D. Xu, S. Yuan, and X. Wu, "Achieving differential privacy in vertically partitioned multiparty learning," *arXiv preprint arXiv:1911.04587*, 2019.
- [49] D. Wang, M. Gaboardi, and J. Xu, "Empirical risk minimization in non-interactive local differential privacy revisited," in *NeurIPS*, 2018.
- [50] J. Bernstein, Y.-X. Wang, K. Azizzadenesheli, and A. Anandkumar, "signsgd: Compressed optimisation for non-convex problems," *arXiv preprint arXiv:1802.04434*, 2018.
- [51] R. Jin, Y. Huang, X. He, H. Dai, and T. Wu, "Stochastic-sign sgd for federated learning with theoretical guarantees," *arXiv preprint arXiv:2002.10940*, 2020.
- [52] J. Hamer, M. Mohri, and A. T. Suresh, "Fedboost: A communication-efficient algorithm for federated learning," in *ICML*. PMLR, 2020.
- [53] S. Feng and H. Yu, "Multi-participant multi-class vertical federated learning," *arXiv preprint arXiv:2001.11154*, 2020.



**Zhaomin Wu** is a Ph.D. candidate in School of Computing of National University of Singapore. He received the bachelor degree in computer science from Huazhong University of Science and Technology (2015-2019). His current research interests include federated learning and privacy.



**Qinbin Li** is currently a Ph.D. candidate in School of Computing of National University of Singapore. His current research interests include machine learning, federated learning, and privacy.



**Bingsheng He** is an Associate Professor in School of Computing of National University of Singapore. He received the bachelor degree in computer science from Shanghai Jiao Tong University (1999-2003), and the PhD degree in computer science in Hong Kong University of Science and Technology (2003-2008). His research interests are high performance computing, distributed and parallel systems, and database systems.

**APPENDIX A  
CHOICE OF HYPERPARAMETERS**

For each dataset, we summarize the hyperparameters of FedOnce-L1 in Table 10. In the table,  $\eta$  refers to the learning rate,  $\lambda$  refers to weight decay,  $b$  refers to batch size,  $T$  refers to the number of epochs,  $d$  refers to the dimension of representations.  $f$  indicates the frequency of permutation matrix  $P$  to be updated. For example, if the update frequency is 3,  $P$  will be updated every three epochs.  $\epsilon$  refers to the overall privacy budget.  $\Omega$  refers to the clipping norm. SGD refers to stochastic gradient descent without momentum, i.e.,  $momentum = 0$ . We adopt the SGD optimizer in FedOnce-L1 because our analysis of differential privacy is based on SGD.

TABLE 10  
Hyperparameters of FedOnce-L1 on each dataset

Dataset	$\epsilon$	Guest Model				Host Model				$d$	$f$	$\Omega$
		$\eta$	$\lambda$	$b$	$T$	$\eta$	$\lambda$	$b$	$T$			
gisette	2	0.3	0	32	10	0.3	0	32	10	3	1	1.0
	4	0.2	0	32	6	0.2	0	32	30	3	1	1.0
	6	0.3	0	128	10	0.6	0	128	15	3	1	1.5
	8	0.3	0	128	10	0.6	0	128	40	3	1	1.5
phishing	2	0.3	0	32	10	0.3	0	32	30	3	1	1.0
	4	0.2	0	32	10	0.2	0	32	30	3	1	1.0
	6	0.3	0	32	10	0.1	0	32	30	3	1	1.0
	8	0.3	0	128	10	0.3	0	128	40	3	1	1.5
MNIST	2	0.5	0	256	6	0.8	0	256	15	8	1	1.5
	4	0.5	0	256	6	0.9	0	256	15	8	1	1.5
Fashion-MNIST	6	0.5	0	256	7	0.9	0	256	15	8	1	1.5
	8	0.5	0	512	10	0.9	0	512	30	8	1	1.5
UJIIndoorLoc	2	0.4	0	128	10	0.5	0	128	10	6	1	1.5
	4	0.4	0	128	12	0.5	0	128	20	6	1	1.5
	6	0.4	0	128	12	0.5	0	128	30	6	1	1.5
	8	0.4	0	128	12	0.3	0	128	60	6	1	1.5

**APPENDIX B  
SUPPLEMENT EXPERIMENT ON CIFAR-10**

To evaluate the performance of FedOnce on complex image datasets, additional experiments are conducted on CIFAR-10 which consists of 60000 32x32 color images in 10 classes. We divide the features both horizontally and vertically into four parties; each party holds a 16x16 partial image with full three channels. Since it is an open problem to fit more advanced deep learning models in unsupervised learning contexts, we use CNN0 in Fig. 4 (the same model we use for other image datasets) as a starting point to demonstrate our proposal on more complex datasets. The experimental setting is similar to that in Section 5.3. We present the performance of each approach as the communication cost increases in Fig. 11.

As can be observed from Fig. 11, multi-round FedOnce 1) significantly outperforms SplitNN when the communication size is small and 2) consistently outperforms SplitNN as the communication size increases.

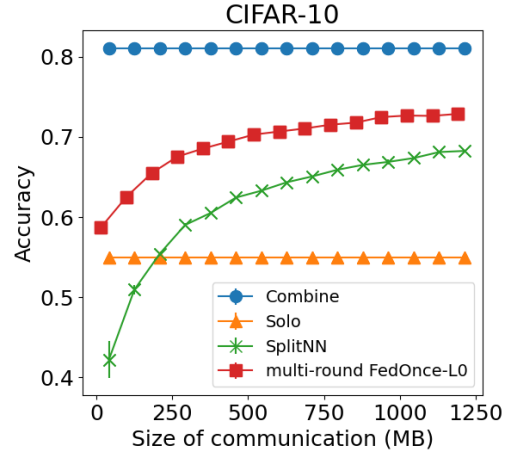


Fig. 11. Performance of multi-round FedOnce-L0 and FedOnce-L0 (i.e., the first data point of multi-round FedOnce-L0) on CIFAR-10 ( $k = 4$ )

**APPENDIX C  
DISCUSSION OF LAGGY PARTIES**

During the distributed machine learning process, some parties may respond slowly due to communication bottlenecks. The existence of these parties, denoted as *laggy parties*, is one of the major hurdles of federated learning. Although FedOnce significantly reduces the effect of laggy parties by communicating for only one round, the impact of laggy parties has not been eliminated. Therefore, in this section, we present an intuitive solution to tackle this issue in FedOnce.

Setting a threshold  $t_r$  for the responding time, all the parties that respond later than  $t_r$  are regarded as laggy parties which do not participate in the training of FedOnce. Specifically, all the representations from laggy parties are set to zero. Thus, FedOnce can still learn from the host party and the representations from other parties.