

# Practical One-Shot Federated Learning for Cross-Silo Setting

Qinbin Li<sup>1</sup>, Bingsheng He<sup>1</sup>, Dawn Song<sup>2</sup>

<sup>1</sup>National University of Singapore

<sup>2</sup>University of California, Berkeley

{qinbin, hebs}@comp.nus.edu.sg, dawnsong@cs.berkeley.edu

## Abstract

Federated learning enables multiple parties to collaboratively learn a model without exchanging their data. While most existing federated learning algorithms need many rounds to converge, one-shot federated learning (i.e., federated learning with a single communication round) is a promising approach to make federated learning applicable in cross-silo setting in practice. However, existing one-shot algorithms only support specific models and do not provide any privacy guarantees, which significantly limit the applications in practice. In this paper, we propose a practical one-shot federated learning algorithm named FedKT. By utilizing the knowledge transfer technique, FedKT can be applied to any classification models and can flexibly achieve differential privacy guarantees. Our experiments on various tasks show that FedKT can significantly outperform the other state-of-the-art federated learning algorithms with a single communication round.

## 1 Introduction

While the size of training data can influence the machine learning model quality a lot, the data are often dispersed over different parties in reality. Due to regulations on data privacy, the data cannot be centralized to a single party for training. A popular solution is federated learning [Kairouz *et al.*, 2019; Li *et al.*, 2019; Yang *et al.*, 2019], which enables multiple parties to collaboratively learn a model without exchanging their local data.

A typical and widely used federated learning algorithm is FedAvg [McMahan and others, 2016]. Its training is an iterative process with four steps in each iteration. First, the server sends the global model to the selected parties. Second, each of the selected parties updates its model with their local data. Third, the updated models are sent to the server. Last, the server averages all the received models to update the global model. There are many variants of FedAvg [Li *et al.*, 2020b; Karimireddy and others, 2020; Lin *et al.*, 2020], which have similar frameworks to FedAvg.

The above iterative algorithms are mainly designed for the cross-device setting, where the parties are mobile devices. In

such a setting, the server is usually managed by the federated learning service provider (e.g., Google) and the parties are the users that are willing to improve their service quality (e.g., Google Keyboard users). The server can sample part of devices to conduct federated learning in each round and there are always users available. However, in cross-silo settings, where the parties are usually organizations, approaches like FedAvg may not work in practice due to the following reasons. First, the algorithm requires parties to participate multi-round training, which is not practical in some scenarios such as model markets [Vartak and others, 2016; Baylor and others, 2017]. Second, federated learning across rounds may suffer from attacks (e.g., inference attacks [Shokri *et al.*, 2017]) from curious parties. Last, it is hard to find a fair and trusted server to lead the training process.

A promising solution is one-shot federated learning (i.e., federated learning with only a single communication round). With one-shot federated learning, the parties can simply sell or upload their local models to a model market (or a model management platform) [Vartak and others, 2016; Baylor and others, 2017]. Then, a buyer or the market can use these models collectively to learn a final model, which is very suitable for cross-silo settings. Such a process largely reduces the multi-round requirements on the stability of the parties. It is natural that sellers put their models into the model market (in return for incentives and benefits, whose design is interesting but out of scope of this paper). Usually, sellers are not engaged in purchasing and consuming the models, and one-shot federated learning algorithms are a must here.

There have been several studies on one-shot federated learning [Guha *et al.*, 2019; Zhou *et al.*, 2020; Yurochkin *et al.*, 2019]. However, existing one-shot federated learning studies have the following obvious shortcomings. First, they usually are specially designed for a special model architecture (i.e., support vector machines or multi-layer perceptrons), which significantly limit the applications in the real-world. Second, they do not provide any privacy guarantees. This is important in the model market scenario, since the models may be sold to anyone including attackers.

In this paper, we propose a new one-shot federated learning algorithm named FedKT (Federated learning via Knowledge Transfer). Inspired by the success of the usage of unlabelled public data in many studies [Papernot *et al.*, 2017; Papernot and others, 2018; Chang *et al.*, 2019; Lin *et al.*,

2020], which often exists such as text and images and can be obtained by public repositories or synthetic data generator or various data markets, we design a two-tier knowledge transfer framework to achieve effective and private one-shot federated learning. As such, unlike most existing studies that only work on either differentiable models (e.g., neural networks [McMahan and others, 2016]) or non-differentiable models (e.g., decision trees [Li *et al.*, 2020a]), FedKT is able to learn any classification model. Moreover, we develop differentially private versions and theoretically analyze the privacy loss of FedKT in order to provide different differential privacy guarantees. Our experiments on various tasks and models show that FedKT significantly outperforms the other state-of-the-art federated learning algorithms with a single communication round.

Our main contributions are as follows.

- Based on the knowledge transfer approach, we propose a new federated learning algorithm named FedKT. To the best of our knowledge, FedKT is the first one-shot federated learning algorithm which can be applied to any classification models.
- We consider comprehensive privacy requirements and show that FedKT is easy to achieve both example-level and party-level differential privacy and theoretically analyze the bound of its privacy cost.
- We conduct experiments on various models and tasks and show that FedKT can achieve much better accuracy compared with the other federated learning algorithms with a single communication round.

## 2 Background and Related Work

### 2.1 Knowledge Transfer

Knowledge transfer has been successfully used in previous studies [Hinton *et al.*, 2015; Papernot *et al.*, 2017; Papernot and others, 2018]. Through knowledge transfer, an ensemble of models can be compressed into a single model. A typical example is the PATE (Private Aggregation of Teacher Ensembles) [Papernot *et al.*, 2017] framework. In this framework, PATE first divides the original dataset into multiple disjoint subsets. A teacher model is trained separately on each subset. Then, the max voting method is used to make predictions on the public unlabelled datasets with the teacher ensemble, i.e., choosing the majority class among the teachers as the label. Last, a student model is trained on the public dataset. A good feature of PATE is that it can easily satisfy differential privacy guarantees by adding noises to the vote counts. Moreover, PATE can be applied to any classification model regardless of the training algorithm. However, PATE is not designed for federated learning.

### 2.2 Federated Learning with Knowledge Transfer

There have been several studies [Li and Wang, 2019; Chang *et al.*, 2019; He *et al.*, 2020; Lin *et al.*, 2020; Zhu *et al.*, 2021] using knowledge transfer in federated learning. [Li and Wang, 2019] needs a public labeled dataset to conduct initial transfer learning, while FedKT only needs a public unlabeled dataset. [Chang *et al.*, 2019] and [He *et al.*, 2020]

have different objectives from FedKT. Specifically, [Chang *et al.*, 2019] designs a robust federated learning algorithm to protect against poisoning attacks. [He *et al.*, 2020] considers cross-device setting with limited computation resources and uses group knowledge transfer to reduce the overload of each edge device. [Lin *et al.*, 2020] has a similar setting with FedKT. They use a public dataset to improve the global model in the server side. As we will show in the experiment, FedKT has a much better accuracy than [Lin *et al.*, 2020] with a single round.

All the above studies conduct in an iterative way, which require many communication rounds to converge and cannot be applied in the model market scenario. Moreover, all existing studies transfer the prediction vectors (i.g., logits) on the public dataset between clients and the server. As we will show in Section 3, FedKT transfers the voting counts and can easily satisfy differential privacy guarantees with a tight theoretical bound on the privacy loss.

We notice that there is a contemporary work [Zhu *et al.*, 2021] which also utilizes PATE in federated learning. While they simply extend PATE to a federated setting, we design a two-tier PATE structure and provide more flexible differential privacy guarantees.

### 2.3 One-Shot Federated Learning

There have been several studies [Yurochkin *et al.*, 2019; Guha *et al.*, 2019; Zhou *et al.*, 2020; Kasturi *et al.*, 2020] on one-shot federated learning. Instead of simply averaging all the model weights in FedAvg, [Yurochkin *et al.*, 2019] propose PFNM by adopting a Bayesian nonparametric model to aggregate the local models when they are multilayer perceptrons (MLPs). Their method shows a good performance in a single communication round and can also be applied in multiple communication rounds. [Guha *et al.*, 2019] propose an one-shot federated learning algorithm to train support vector machines (SVMs) in both supervised and semi-supervised settings. [Zhou *et al.*, 2020] and [Kasturi *et al.*, 2020] transfer the synthetic data or data distribution to the server, which trains the final model using generated dataset. Such data sharing approaches do not fit the mainstream model sharing schemes. Moreover, all existing one-shot federated learning studies do not provide privacy guarantees, which is important especially in the model market scenario where the models are sold and may be bought by anyone including attackers.

### 2.4 Differential Privacy

Differential privacy [Dwork *et al.*, 2014] is a popular standard of privacy protection. It guarantees that the probability of producing a given output does not depend much on whether a particular data record is included in the input dataset or not. It has been widely used to protect the machine learning models [Abadi *et al.*, 2016; Choquette-Choo *et al.*, 2021].

**Definition 1.** ( $(\epsilon, \delta)$ -Differential Privacy) Let  $\mathcal{M}: \mathcal{D} \rightarrow \mathcal{R}$  be a randomized mechanism with domain  $\mathcal{D}$  and range  $\mathcal{R}$ .  $\mathcal{M}$  satisfies  $(\epsilon, \delta)$ -differential privacy if for any two adjacent inputs  $d, d' \in \mathcal{D}$  and any subset of outputs  $S \subseteq \mathcal{R}$  it holds that:

$$\Pr[\mathcal{M}(d) \in S] \leq e^\epsilon \Pr[\mathcal{M}(d') \in S] + \delta. \quad (1)$$

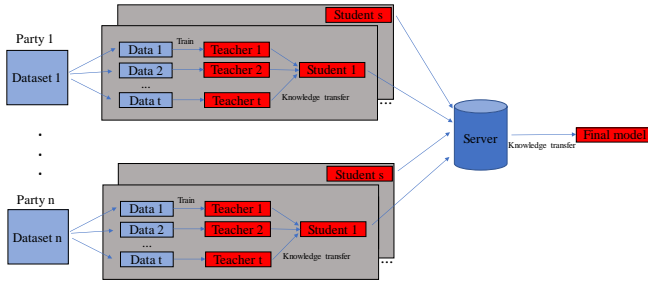


Figure 1: The framework of FedKT

The moments accountant method [Abadi *et al.*, 2016] is a state-of-the-art approach to track the privacy loss. We briefly introduce the key concept, and refer readers to the previous paper [Abadi *et al.*, 2016] for more details.

**Definition 2.** (Privacy Loss) Let  $\mathcal{M}: \mathcal{D} \rightarrow \mathcal{R}$  be a randomized mechanism. Let  $\text{aux}$  denote an auxiliary input. For two adjacent inputs  $d, d' \in \mathcal{D}$ , an outcome  $o \in \mathcal{R}$ , the privacy loss at  $o$  is defined as:

$$c(o; \mathcal{M}, \text{aux}, d, d') \triangleq \log \frac{\Pr[\mathcal{M}(\text{aux}, d) = o]}{\Pr[\mathcal{M}(\text{aux}, d') = o]}. \quad (2)$$

**Definition 3.** (Moments Accountant) Let  $\mathcal{M}: \mathcal{D} \rightarrow \mathcal{R}$  be a randomized mechanism. Let  $\text{aux}$  denote an auxiliary input. For two adjacent inputs  $d, d'$ , the moments accountant is defined as:

$$\alpha_{\mathcal{M}}(\lambda) \triangleq \max_{\text{aux}, d, d'} \alpha_{\mathcal{M}}(\lambda; \text{aux}, d, d') \quad (3)$$

where  $\alpha_{\mathcal{M}}(\lambda; \text{aux}, d, d') \triangleq \log \mathbb{E}_o[\exp(\lambda c(o; \mathcal{M}, \text{aux}, d, d'))]$  is the log of moment generating function.

The moments have good composability and can be easily converted to  $(\epsilon, \delta)$ -differential privacy [Abadi *et al.*, 2016].

**Party-level Differential Privacy** In addition to the standard example-level differential privacy, party-level differential privacy [Geyer *et al.*, 2017; McMahan *et al.*, 2018] is more strict and attractive in the federated setting. Instead of aiming to protect a single record, party-level differential privacy ensures that the model does not reveal whether a party participated in federated learning or not.

**Definition 4.** (Party-adjacent Datasets) Let  $d$  and  $d'$  be two datasets of training examples, where each example is associated with a party. Then,  $d$  and  $d'$  are party-adjacent if  $d'$  can be formed by changing the examples associated with a single party from  $d$ .

### 3 Our Approach

**Problem Statement** Suppose there are  $n$  parties  $P_1, \dots, P_n$ . We use  $\mathcal{D}^i$  to denote the dataset of  $P_i$ . With the help of a central server and a public unlabelled dataset  $\mathcal{D}_{aux}$ , our objective is to build a machine learning model over the datasets  $\bigcup_{i \in [n]} \mathcal{D}^i$  without exchanging the raw data. Moreover, the learning process should be able to support three different privacy level settings: (1)  $L0$ : like most federated learning

studies [McMahan and others, 2016; Yurochkin *et al.*, 2019; Lin *et al.*, 2020], the possible inference attacks [Shokri *et al.*, 2017; Fredrikson *et al.*, 2015] on the models are not considered. At  $L0$ , we do not enforce any privacy mechanism on the model. (2)  $L1$  (server-noise): in the case where the final model has to be sent back to the parties or even published, it should satisfy differential privacy guarantees to protect against potential inference attacks. (3)  $L2$  (party-noise): in the case where the server is curious and all the models transferred from the parties to the server during training should satisfy differential privacy guarantees. These three levels have their own application scenarios in practice. For example, in the model market,  $L0$  is for a model market within an enterprise where different departments can share model but not the data due to data privacy regulations, and departments are trusted.  $L1$  is similar to  $L0$  but additionally the final model will be published in the market.  $L2$  is for the public market where sellers need to protect their own data privacy.

**The Overall Framework** The framework of FedKT is shown in Figure 1. It is designed to be a simple and practical one-shot algorithm for cross-silo settings. Specifically, FedKT adopts a **two-tier knowledge transfer structure**. On the party side, each party uses knowledge transfer to learn student models and sends them to the server. On the server side, the server takes the received models as teachers to learn a final model using knowledge transfer again. The final model is sent back to the parties and used for predictions. Two techniques, multi-partitioning and consistent voting, are proposed to improve the performance.

**Learning Student Models on the Parties (Multi-Partitioning)** Locally, each party has to create  $s$  ( $s \geq 1$ ) partitions and learn a student model on each partition. Each partition handles the entire local dataset. Since the operations in each partition are similar, here we describe the process in one partition for ease of presentation. Inside a partition, the local dataset is divided into  $t$  disjoint subsets. We train a teacher model separately on each subset, denoted as  $T_1, \dots, T_t$ . Then, the ensemble of teacher models is used to make predictions on the public dataset  $\mathcal{D}_{aux}$ . For an example  $\mathbf{x} \in \mathcal{D}_{aux}$ , the vote count of class  $m$  is the number of teachers that predicts  $m$ , i.e.,  $v_m(\mathbf{x}) = |\{i : i \in [t], T_i(\mathbf{x}) = m\}|$ . The prediction result of the ensemble is the class that has the maximum vote counts, i.e.,  $f(\mathbf{x}) = \arg \max_m v_m(\mathbf{x})$ . Then, we use the public dataset  $\mathcal{D}_{aux}$  with the predicted labels to train a student model. For each partition, we get a student model with the above steps. After all the student models are trained, the parties send their student models to the server for further processing.

**Learning the Final Model on the Server** Suppose the student models of party  $i$  are denoted as  $U_1^i, \dots, U_s^i$  ( $i \in [n]$ ). After receiving all the student models, like the steps on the party side, the server can use these student models as an ensemble to make predictions on the public dataset  $\mathcal{D}_{aux}$ . The public dataset with the predicted labels is used to train the final model.

**Consistent voting** Here we introduce a technique named *consistent voting* for computing the vote counts of each class. If the student models of a party make the same prediction

---

**Algorithm 1: The FedKT algorithm**

---

**Input:** local datasets  $\mathcal{D}^1, \dots, \mathcal{D}^n$ , number of partitions  $s$  in each party, number of subsets  $t$  in each partition, number of classes  $u$ , public dataset  $\mathcal{D}_{aux}$ , privacy parameter  $\gamma$ , privacy level  $l$

**Output:** The final model  $F$ .

```
1 for  $i = 1, \dots, n$  do
  /* Conduct on party  $P_i$  */
2  Create  $s$  partitions (i.e.,  $D_1^i, \dots, D_s^i$ ) on dataset  $\mathcal{D}^i$ 
   such that  $D_j^i = \bigcup_{k \in [t]} D_{j,k}^i$  for all  $j \in [s]$ , where
    $D_{j,k}^i$  is a subset.
3  for  $j = 1, \dots, s$  do
4    for  $k = 1, \dots, t$  do
5      Train a teacher model  $T_{j,k}^i$  on subset  $D_{j,k}^i$ .
6      for all  $\mathbf{x} \in \mathcal{D}_{aux}$  do
7        for all  $m \in [u]$  do Vote-1
8           $v_m(\mathbf{x}) \leftarrow |\{k : k \in [t], T_{j,k}^i(\mathbf{x}) = m\}|$ 
9          if  $l == L2$  then
10              $v_m(\mathbf{x}) \leftarrow v_m(\mathbf{x}) + Lap(1/\gamma)$ 
11           $f(\mathbf{x}) = \arg \max_m v_m(\mathbf{x})$ 
12      Train a student model  $U_j^i$  on dataset
        $\{(\mathbf{x}, f(\mathbf{x}))\}_{\mathbf{x} \in \mathcal{D}_{aux}}$ .
13  Send the student models  $\{U_j^i : j \in [s]\}$  to server.
  /* Conduct on the server */
14 for all  $\mathbf{x} \in \mathcal{D}_{aux}$  do
15   for all  $i \in [n]$  do
16     for all  $m \in [u]$  do
17        $v_m^i(\mathbf{x}) \leftarrow |\{k : k \in [s], U_k^i(\mathbf{x}) = m\}|$ 
18   for all  $m \in [u]$  do Vote-2
19      $v_m(\mathbf{x}) \leftarrow s \cdot |\{i : i \in [n], v_m^i(\mathbf{x}) = s\}|$ 
20     if  $l == L1$  then
21        $v_m(\mathbf{x}) \leftarrow v_m(\mathbf{x}) + Lap(1/\gamma)$ 
22    $f(\mathbf{x}) = \arg \max_m v_m(\mathbf{x})$ 
23 Train the final model  $F$  on dataset  $\{(\mathbf{x}, f(\mathbf{x}))\}_{\mathbf{x} \in \mathcal{D}_{aux}}$ .
```

KD2

on an example, we take their predictions into account. Otherwise, the party is not confident at predicting this example and thus we ignore the predictions of its student models. Formally, given an example  $\mathbf{x} \in \mathcal{D}_{aux}$ , we first compute the vote count of class  $m$  on the student models of party  $i$  as  $v_m^i(\mathbf{x}) = |\{k : k \in [s], U_k^i(\mathbf{x}) = m\}|$ . Next, with consistent voting, the final vote count of class  $m$  on all parties is computed as  $v_m(\mathbf{x}) = s \cdot |\{i : i \in [n], v_m^i(\mathbf{x}) = s\}|$ .

**Differentially Private Versions of FedKT** FedKT can easily satisfy differential privacy guarantees by providing differentially private prediction results to the query dataset. Given the privacy parameter  $\gamma$ , we can add noises to the vote count histogram such that  $f(\mathbf{x}) = \arg \max_m \{v_m(\mathbf{x}) + Lap(\frac{1}{\gamma})\}$ , where  $Lap(\frac{1}{\gamma})$  is the noises generated from Laplace distribution with location 0 and scale  $\frac{1}{\gamma}$ . Note that we do not need to add noises on both the parties and the server. For the  $L1$

setting, we only need to add noises on the server side. The parties can train and send non-differentially private student models to the server. For the  $L2$  setting, we only need to add noises on the party side so that the student models are differentially private. Then, the final model naturally satisfies differential privacy guarantees. More analysis on privacy loss will be presented in Section 4.

Algorithm 1 shows the whole training process of FedKT. In the algorithm, for each party (i.e., Line 1) and its each partition (i.e., Line 3), we train a student model using knowledge transfer (i.e., Lines 4-12). The student models are sent to the server. Then, the server trains the final model using knowledge transfer again (i.e., Lines 14-23). For different privacy level settings, we have the corresponding noises injection operations on the server side (i.e., Lines 20-21) or the party side (i.e., Lines 9-10).

**Communication/Computation Overhead of FedKT** Suppose the size of each model is  $M$ . Then, the total communication size of FedKT is  $nsM$  for sending the student models to the server. Suppose the number of communication rounds in FedAvg is  $r$  and all the parties participate in the training in every iteration. Then the total communication size of FedAvg is  $2nMr$  including the server sends the global model to the parties<sup>1</sup> and the parties send the local models to the server. Thus, when  $r > \frac{s}{2}$ , the communication cost of FedAvg is higher than FedKT. This value can be quite small, e.g.,  $r = 2$  if we set  $s = 2$ . Moreover, when  $s = 2$ , FedAvg has the same communication cost with FedKT in the first round. The computation overhead of FedKT is usually larger than FedAvg in each round since FedKT needs to train multiple teacher and student models. However, the computation overhead is acceptable in the cross-silo setting, where the parties (e.g., companies, data centers) usually have a relatively large computation power.

## 4 Data-Dependent Privacy Analysis of FedKT

In this section, we use the moments accountant method [Abadi *et al.*, 2016] to track the privacy loss in the training process. For  $L1$  setting, we mainly consider the party-level differential privacy [Geyer *et al.*, 2017; McMahan *et al.*, 2018], which is more attractive in the federated setting. Instead of aiming to protect a single record, party-level differential privacy ensures that the learned model does not reveal whether a party participated in federated learning or not. For  $L2$  setting, we mainly consider the example-level differential privacy since it is not practical to require each local model to satisfy party-level differential privacy. For proofs of the theorems in this section, please refer to Section 4 of Appendix.

**FedKT-L1** Considering we change the whole dataset of a party, then at most  $s$  student models will be influenced. Thus, on the server side, the sensitivity of the vote count histogram is  $2s$  (i.e., the vote count of a class increases by  $s$  and the vote count of another class decreases by  $s$ ). According to the Laplace mechanism, we have the following theorem.

<sup>1</sup>It also happens in the first round to ensure all parties have the same initialized model [McMahan and others, 2016].



**Theorem 1.** Let  $\mathcal{M}$  be the  $f$  function executed on the server side. Given the number of partitions  $s$  and the privacy parameter  $\gamma$ ,  $\mathcal{M}$  satisfies  $(2s\gamma, 0)$  party-level differential privacy.

Similar with [Papernot *et al.*, 2017], we conduct a data-dependent privacy analysis for FedKT with the moments accountant method.

**Theorem 2.** Let  $\mathcal{M}$  be  $(2s\gamma, 0)$  party-level differentially private. Let  $q \geq \Pr[\mathcal{M}(d) \neq o^*]$  for some outcome  $o^*$ . Let  $l, \gamma \geq 0$  and  $q < \frac{e^{2s\gamma}-1}{e^{4s\gamma}-1}$ . For any  $\mathbf{aux}$  and any two party-adjacent datasets  $d$  and  $d'$ ,  $\mathcal{M}$  satisfies

$$\alpha_{\mathcal{M}}(l; \mathbf{aux}, d, d') \leq \min(\log((1-q)\left(\frac{1-q}{1-e^{2s\gamma}q}\right)^l + qe^{2s\gamma l}), 2s^2\gamma^2 l(l+1)).$$

Here  $\Pr[\mathcal{M}(d) \neq o^*]$  can be bounded by Lemma 4 of [Papernot *et al.*, 2017].

With Theorem 2, we can track the privacy loss of each query [Abadi *et al.*, 2016].

**FedKT-L2** For each partition on the party side, we add Laplace noises to the vote counts, which is same with the PATE approach. Thus, we have the following theorem.

**Theorem 3.** Let  $\mathcal{M}$  be the  $f$  function executed on each partition of a party. Let  $q \geq \Pr[\mathcal{M}(d) \neq o^*]$  for some outcome  $o^*$ . Let  $l, \gamma \geq 0$  and  $q < \frac{e^{2\gamma}-1}{e^{4\gamma}-1}$ . For any  $\mathbf{aux}$  and any two adjacent datasets  $d$  and  $d'$ ,  $\mathcal{M}$  satisfies

$$\alpha_{\mathcal{M}}(l; \mathbf{aux}, d, d') \leq \min(\log((1-q)\left(\frac{1-q}{1-e^{2\gamma}q}\right)^l + qe^{2\gamma l}), 2\gamma^2 l(l+1)).$$

After computing the privacy loss of each party, we can use the parallel composition to compute the privacy loss of the final model.

**Theorem 4.** Suppose the student models of party  $P_i$  satisfy  $(\varepsilon_i, \delta)$ -differential privacy. Then, the final model satisfies  $(\max_i \varepsilon_i, \delta)$ -differential privacy.

Note that the above privacy analysis is data-dependent. Thus, the final privacy budget is also data-dependent and may have potential privacy breaches if we publish the budget. Like previous studies [Papernot *et al.*, 2017; Jordon *et al.*, 2019], we report the data-dependent privacy budgets in the experiments. As future work, we plan to use the smooth sensitivity algorithm [Nissim *et al.*, 2007] to add noises to the privacy losses. Also, we may get a tighter bound of the privacy loss if adopting the Gaussian noises [Papernot and others, 2018].

## 5 Experiments

To evaluate FedKT, we conduct experiments on four public datasets: (1) A random forest on *Adult* dataset. The number of trees is set to 100 and the maximum tree depth is set to 6. (2) A gradient boosting decision tree (GBDT) model on *cod-rna* dataset. The maximum tree depth is set to 6. (3)

A multilayer perceptron (MLP) with two hidden layers on *MNIST* dataset. Each hidden layer has 100 units using ReLU activations. (4) A CNN on extended *SVHN* dataset. The CNN has two 5x5 convolution layers followed with 2x2 max pooling (the first with 6 channels and the second with 16 channels), two fully connected layers with ReLU activation (the first with 120 units and the second with 84 units), and a final softmax output layer. For the first two datasets, we split the original dataset at random into train/test/public sets with a 75%/12.5%/12.5% proportion. For MNIST and SVHN, we use one half of the original test dataset as the public dataset and the remaining as the final test dataset. Like many existing studies [Yurochkin *et al.*, 2019; Lin *et al.*, 2020; Li *et al.*, 2021], we use the Dirichlet distribution to simulate the heterogeneous data partition among the parties. Suppose there are  $n$  parties. We sample  $p_k \sim \text{Dir}_n(\beta)$  and allocate a  $p_{k,j}$  proportion of the instances of class  $k$  to party  $j$ , where  $\text{Dir}(\beta)$  is the Dirichlet distribution with a concentration parameter  $\beta$  (0.5 by default). By default, we set the number of parties to 50 for *Adult* and *cod-rna* and to 10 for MNIST and SVHN. We set  $s$  to 2 and  $t$  to 5 by default for all datasets. For more details in the experimental settings and study of the hyper-parameters, please refer to Section B.1 and Section B.2 of Appendix. The code is publicly available <sup>2</sup>.

We compare FedKT with eight baselines: (1) SOLO: each party trains its model locally and does not participate in federated learning. (2) FedAvg [McMahan and others, 2016]; (3) FedProx [Li *et al.*, 2020b]; (4) SCAFFOLD [Karimireddy and others, 2020]; (5) FedDF [Lin *et al.*, 2020]; (6) PNFM [Yurochkin *et al.*, 2019]; (7) PATE [Papernot *et al.*, 2017]: we use the PATE framework to train a final model on all data in a centralized setting (i.e., only a single party with the whole dataset) without adding noises. This method defines an upper bound of learning a final model using knowledge transfer with public unlabelled data. (8) XGBoost [Chen and Guestrin, 2016]: the XGBoost algorithm for the GBDT model on the whole dataset in a centralized setting. This method defines an upper bound of learning the GBDT model. Here approaches (2)-(5) are popular or state-of-the-art federated learning algorithms and approach (6) is an one-shot algorithm. Moreover, same as FedKT, approaches (5) and (7) also utilize the unlabeled public dataset.

### 5.1 Effectiveness

Table 1 shows the accuracy of FedKT<sup>3</sup> compared with the other baselines. For fair comparison and practical usage in model markets, we run all approaches for a single communication round. For SOLO, we report the average accuracy of the parties. From this table, we have the following observations. First, except for FedKT, the other federated learning algorithms can only learn specific models. FedKT is able to learn all the studied models including trees and neural networks. Second, FedKT can achieve much better performance than the other federated learning algorithms running with a single round. FedKT can achieve about 6.5% higher accu-

<sup>2</sup><https://github.com/QinbinLi/FedKT>

<sup>3</sup>For simplicity, we use FedKT to denote FedKT-L0, unless specified otherwise.

Table 1: The test accuracy comparison between FedKT and the other baselines in a single round.

| Datasets | FedKT                   | SOLO  | FedAvg | FedProx | SCAFFOLD | FedDF | PNFM  | PATE  | XGBOOST |
|----------|-------------------------|-------|--------|---------|----------|-------|-------|-------|---------|
| Adult    | <b>82.2%</b> $\pm$ 0.6% | 68.6% | \      |         |          |       |       | 83.5% | \       |
| cod-rna  | <b>88.3%</b> $\pm$ 0.6% | 65.0% |        |         |          |       |       | 91.1% | 91.2%   |
| MNIST    | <b>90.5%</b> $\pm$ 0.3% | 69.0% | 62.8%  | 44.3%   | 51.7%    | 83.8% | 65.9% | 92.7% | \       |
| SVHN     | <b>83.2%</b> $\pm$ 0.4% | 62.8% | 26.8%  | 20.1%   | 16.2%    | 77.2% | \     | 86.6% | \       |

Table 2: The privacy loss and test accuracy of FedKT-L1 and FedKT-L2 given different  $\gamma$  and number of queries. L0 acc is the test accuracy of FedKT-L0. The failure probability  $\delta$  is set to  $10^{-5}$ .

| datasets | FedKT-L1 |          |            |       |                 | FedKT-L2 |          |            |       |                 |
|----------|----------|----------|------------|-------|-----------------|----------|----------|------------|-------|-----------------|
|          | $\gamma$ | #queries | $\epsilon$ | acc   | non-private acc | $\gamma$ | #queries | $\epsilon$ | acc   | non-private acc |
| Adult    | 0.04     | 0.5%     | 2.56       | 76.8% | 82.2%           | 0.04     | 0.5%     | 2.59       | 77.6% | 82.4%           |
|          | 0.04     | 1.0%     | 4.73       | 80.2% |                 | 0.04     | 1.0%     | 3.72       | 78.6% |                 |
| cod-rna  | 0.06     | 0.50%    | 5.48       | 82.6% | 88.0%           | 0.05     | 1.0%     | 4.51       | 82.7% | 89.7%           |
|          | 0.1      | 0.5%     | 6.89       | 84.7% |                 | 0.05     | 2.0%     | 9.78       | 84.7% |                 |

racy than FedDF, which also utilizes the public dataset. Third, although PNFM is an one-shot algorithm specially designed for MLPs, FedKT outperforms PNFM about 25% accuracy on MNIST. Last, the accuracy of FedKT is close to PATE and XGBoost, which means our design has little accuracy loss compared with the centralized setting. The gap between FedKT and the upper bound is very small.

## 5.2 Privacy

We run FedKT with different  $\gamma$  and number of queries. When running FedKT-L1, we tune the percentage of number of queries on the server side. When running FedKT-L2, on the contrary, we tune the percentage of number of queries on the party side. The selected results on Adult and cod-rna are reported in Table 2. While differentially private FedKT does not need any knowledge on the model architecture, the accuracy is still comparable to the non-private version given a privacy budget less than 10. For more results, please refer to Section B.4 of Appendix.

## 5.3 Size of the Public Dataset

We show the performance of FedKT with different size of public dataset. The results are shown in Table 3. We can observe that FedKT is stable even though reducing the size of the public dataset. The accuracy decreases no more than 1% and 2% using only 20% of the public dataset on cod-rna and MNIST (i.e., 1807 examples on cod-rna and 1000 examples on MNIST), respectively. Moreover, the accuracy is almost unchanged on Adult.

## 5.4 Extend to Multiple Rounds

While FedKT is an one-shot algorithm, it is still applicable in the scenarios where multiple rounds are allowed. FedKT can be used as an initialization step to learn a global model in the first round. Then the parties can use the global model to conduct iterative federated learning algorithms. By combining FedKT with FedProx (denoted as FedKT-Prox), FedKT-Prox is much more communication-efficient than the other algorithms. FedKT-Prox needs about 11 rounds to achieve 87% accuracy, while the other approaches need at least 32 rounds. For more details, please refer to Section B.3 of Appendix.

Table 3: The test accuracy of FedKT with different size of the public dataset. We vary the portion of the public dataset used in the training from 20% to 100%.

| Datasets | Portion of the public dataset used in training |       |       |       |       |
|----------|--|-------|-------|-------|-------|
|          | 20%  | 40%   | 60%   | 80%   | 100%  |
| Adult    | 82.1%  | 82.1% | 82.1% | 82.3% | 82.2% |
| cod-rna  | 87.3%  | 87.6% | 87.8% | 87.8% | 88.3% |
| MNIST    | 89.3%  | 89.7% | 90.2% | 90.3% | 90.5% |
| SVHN     | 80.1%  | 81.3% | 82.1% | 82.9% | 83.2% |

## 6 Conclusions

Motivated by the rigid multi-round training of current federated learning algorithms and emerging applications like model markets, we propose FedKT, a one-shot federated learning algorithm for the cross-silo setting. Our experiments show that FedKT can learn different models with a much better accuracy than the other state-of-the-art algorithms with a single communication round. Moreover, the accuracy of differentially private FedKT is comparable to the non-differentially private version with a modest privacy budget. Overall, FedKT is a practical one-shot solution for model-based sharing in cross-silo federated learning.

## Acknowledgements

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.

## References

[Abadi *et al.*, 2016] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *ACM CCS*, 2016.

- [Baylor and others, 2017] Denis Baylor et al. Tfx: A tensorflow-based production-scale machine learning platform. In *ACM SIGKDD*, 2017.
- [Bun and Steinke, 2016] Mark Bun and Thomas Steinke. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *TCC*. Springer, 2016.
- [Chang et al., 2019] Hongyan Chang, Virat Shejwalkar, Reza Shokri, and Amir Houmansadr. Cronus: Robust and heterogeneous collaborative learning with black-box knowledge transfer. *arXiv:1912.11279*, 2019.
- [Chen and Guestrin, 2016] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *ACM SIGKDD*, 2016.
- [Choquette-Choo et al., 2021] Christopher A. Choquette-Choo, Natalie Dullerud, Adam Dziedzic, Yunxiang Zhang, Somesh Jha, Nicolas Papernot, and Xiao Wang. Capc learning: Confidential and private collaborative learning. In *ICLR*, 2021.
- [Dwork et al., 2014] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 2014.
- [Fredrikson et al., 2015] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. In *ACM CCS*, 2015.
- [Geyer et al., 2017] Robin C Geyer, Tassilo Klein, and Moin Nabi. Differentially private federated learning: A client level perspective. *arXiv preprint arXiv:1712.07557*, 2017.
- [Guha et al., 2019] Neel Guha, Ameet Talwalkar, and Virginia Smith. One-shot federated learning. *arXiv preprint arXiv:1902.11175*, 2019.
- [He et al., 2020] Chaoyang He, Murali Annavaram, and Salman Avestimehr. Group knowledge transfer: Federated learning of large cnns at the edge. *NeurIPS*, 33, 2020.
- [Hinton et al., 2015] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [Jordon et al., 2019] James Jordon, Jinsung Yoon, and Michaela van der Schaar. Differentially private bagging: Improved utility and cheaper privacy than subsample-and-aggregate. In *NeurIPS*, 2019.
- [Kairouz et al., 2019] Peter Kairouz, H Brendan McMahan, et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- [Karimireddy and others, 2020] Sai Praneeth Karimireddy et al. Scaffold: Stochastic controlled averaging for on-device federated learning. In *ICML*. PMLR, 2020.
- [Kasturi et al., 2020] Anirudh Kasturi, Anish Reddy Ellore, and Chittaranjan Hota. Fusion learning: A one shot federated learning. In *ICCS*. Springer, 2020.
- [Li and Wang, 2019] Daliang Li and Junpu Wang. Fedmd: Heterogenous federated learning via model distillation. *arXiv preprint arXiv:1910.03581*, 2019.
- [Li et al., 2019] Qinbin Li, Zeyi Wen, Zhaomin Wu, Sixu Hu, Naibo Wang, Yuan Li, Xu Liu, and Bingsheng He. A survey on federated learning systems: vision, hype and reality for data privacy and protection. *arXiv preprint arXiv:1907.09693*, 2019.
- [Li et al., 2020a] Qinbin Li, Zeyi Wen, and Bingsheng He. Practical federated gradient boosting decision trees. In *AAAI*, 2020.
- [Li et al., 2020b] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. In *ML-Sys*, 2020.
- [Li et al., 2021] Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In *CVPR*, 2021.
- [Lin et al., 2020] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. *NeurIPS*, 33, 2020.
- [McMahan and others, 2016] H Brendan McMahan et al. Communication-efficient learning of deep networks from decentralized data. *arXiv:1602.05629*, 2016.
- [McMahan et al., 2018] H. Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially private recurrent language models. In *ICLR*, 2018.
- [Nissim et al., 2007] Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. Smooth sensitivity and sampling in private data analysis. In *ACM STOC*, 2007.
- [Papernot and others, 2018] Nicolas Papernot et al. Scalable private learning with pate. In *ICLR*, 2018.
- [Papernot et al., 2017] Nicolas Papernot, Martín Abadi, Úlfar Erlingsson, Ian Goodfellow, and Kunal Talwar. Semi-supervised knowledge transfer for deep learning from private training data. In *ICLR*, 2017.
- [Shokri et al., 2017] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *IEEE SP*, 2017.
- [Vartak and others, 2016] Manasi Vartak et al. Modeldb: a system for machine learning model management. In *Proceedings of the Workshop on Human-In-the-Loop Data Analytics*, 2016.
- [Yang et al., 2019] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept and applications. *ACM TIST*, 2019.
- [Yurochkin et al., 2019] Mikhail Yurochkin, Mayank Agarwal, Soumya Ghosh, Kristjan Greenewald, Nghia Hoang, and Yasaman Khazaeni. Bayesian nonparametric federated learning of neural networks. In *ICML*. PMLR, 2019.
- [Zhou et al., 2020] Yanlin Zhou, George Pu, Xiyao Ma, Xiaolin Li, and Dapeng Wu. Distilled one-shot federated learning. *arXiv preprint arXiv:2009.07999*, 2020.
- [Zhu et al., 2021] Yuqing Zhu, Xiang Yu, Yi-Hsuan Tsai, Francesco Pittaluga, Masoud Faraki, Manmohan Chandraker, and Yu-Xiang Wang. Voting-based approaches for differentially private federated learning, 2021.

## Appendix

In this appendix, we first present the data-dependent privacy analysis of FedKT-L1 and FedKT-L2 in Section A. Next, in Section B.1, we show the details of our experimental settings. Then, we show additional experimental results in Section B.2 to Section B.4. Specifically, in Section B.2, we study the performance of FedKT with different hyper-parameters. In Section 5.3, we study the effect of the size of the public dataset. In Section B.3, we extend FedKT to multiple rounds. Last, in section B.4, we show the performance of FedKT-L1 and FedKT-L2.

### A Privacy Analysis of FedKT

In this section, we analyze the privacy loss of FedKT-L1 and FedKT-L2 using the moments accountant method [Abadi *et al.*, 2016].

#### A.1 Proof of Theorem 2

*Proof.* We first introduce two theorems from the previous studies, which will be used in our analysis. The first theorem is from [Bun and Steinke, 2016] and the second theorem is from [Papernot *et al.*, 2017].

**Lemma 5.** Let  $\mathcal{M}$  be  $(2\gamma, 0)$ -differentially private. For any  $l$ ,  $\mathbf{aux}$ , neighboring inputs  $d$  and  $d'$ , we have

$$\alpha_{\mathcal{M}}(l; \mathbf{aux}, d, d') \leq 2\gamma^2 l(l+1)$$

**Lemma 6.** Let  $\mathcal{M}$  be  $(2\gamma, 0)$ -differentially private and  $q \geq \Pr[\mathcal{M}(d) \neq o^*]$  for some outcome  $o^*$ . Let  $l, \gamma \geq 0$  and  $q < \frac{e^{2\gamma} - 1}{e^{4\gamma} - 1}$ . Then for any  $\mathbf{aux}$  and any neighbor  $d'$  of  $d$ ,  $\mathcal{M}$  satisfies

$$\alpha_{\mathcal{M}}(l; \mathbf{aux}, d, d') \leq \log((1-q) \left( \frac{1-q}{1-e^{2\gamma}q} \right)^l + qe^{2\gamma l})$$

$\Pr[\mathcal{M}(d) \neq o^*]$  can be bounded by the following lemma.

**Lemma 7.** Let  $\mathbf{v}$  be the label score vector for an instance  $d$  with  $v_{o^*} \geq v_o$  for all  $o$ . Then

$$\Pr[\mathcal{M}(d) \neq o^*] \leq \sum_{o \neq o^*} \frac{2 + \gamma(v_{o^*} - v_o)}{4 \exp(\gamma(v_{o^*} - v_o))}$$

According to Theorem 1 of the paper,  $\mathcal{M}$  satisfies  $(2s\gamma, 0)$  party-level differential privacy. Thus, substituting  $s\gamma$  into Lemma 5 and 6, we can get

$$\alpha_{\mathcal{M}}(l; \mathbf{aux}, d, d') \leq \min(\log((1-q) \left( \frac{1-q}{1-e^{2s\gamma}q} \right)^l + qe^{2s\gamma l}), 2s^2\gamma^2 l(l+1)).$$

□

#### A.2 Proof of Theorem 3

*Proof.* The noises injection on the party side is similar to the PATE approach. The sensitivity of the  $f$  function executed on the party side is 2. Thus, we can get the theorem by combining Lemma 5 and Lemma 6. We have

$$\alpha_{\mathcal{M}}(l; \mathbf{aux}, d, d') \leq \min(\log((1-q) \left( \frac{1-q}{1-e^{2\gamma}q} \right)^l + qe^{2\gamma l}), 2\gamma^2 l(l+1)).$$

□

#### A.3 Example-Level Differential Privacy Analysis of FedKT-L1

Here we analyze the example-level differential privacy of FedKT-L1. If we change a single example of the original dataset (i.e., the union of all the local datasets), only a single party will be influenced. More precisely, for each partition of the party, only a single teacher model will be influenced. Then, even though changing a single record, the student model is still unchanged if the top-2 vote counts of the teachers differ at least 2. Thus, if not applying consistent voting, we have the following theorem.

**Theorem 8.** Let  $\mathcal{M}$  be the  $f$  function executed in the server. Let  $q \geq \Pr[\mathcal{M}(d) \neq o^*]$  for some outcome  $o^*$ . Let  $\mathcal{D}_{aux}$  denotes the query dataset. Given a query  $i$ , suppose the top-2 vote counts are  $v_1^i$  and  $v_2^i$ . In party  $P_i$ , let  $z_i$  denotes the number of partitions that there  $\exists q \in \mathcal{D}_{aux}$  such that  $v_1^q - v_2^q \leq 1$  when training the student model. Let  $z = \max_i z_i$ . Let  $l, \gamma \geq 0$  and  $q < \frac{e^{2z\gamma} - 1}{e^{4z\gamma} - 1}$ . We have

$$\alpha_{\mathcal{M}}(l; \mathbf{aux}, d, d') \leq \min(\log((1-q) \left( \frac{1-q}{1-e^{2z\gamma}q} \right)^l + qe^{2z\gamma l}), 2z^2\gamma^2 l(l+1))$$

*Proof.* Given a query dataset  $\mathcal{D}_{aux}$ ,  $z$  is the number of partitions such that there exists a query that the top-2 vote counts differ at most 1. In other words, there are at most  $z$  student models will be changed if we change a single record of the original dataset. Thus, the vote counts change by at most  $2z$  on the server side and  $\mathcal{M}$  is  $(2z\gamma, 0)$ -differentially private with respect to  $d$  and  $\mathcal{D}_{aux}$ . Then, we can get this theorem by substituting  $\gamma$  of Theorem 5 and Theorem 6 to  $z\gamma$ . □

Note that the example-level differential privacy is same as party-level differential privacy when  $z = s$ . Also, if we applied consistent voting in FedKT-L1, the vote counts may change by  $s$  even if only a single student model is affected. Thus, the example-level differential privacy of FedKT-L1 is usually same as party-level differential privacy.

#### A.4 Party-Level Differential Privacy Analysis of FedKT-L2

Here we study party-level differential privacy of FedKT-L2.

**Theorem 9.** Let  $\mathcal{M}$  be the  $f$  function executed on each partition of a party. Given the number of subsets in each partition  $t$  and the privacy parameter  $\gamma$ ,  $\mathcal{M}$  satisfies  $(2t\gamma, 0)$  party-level differential privacy.



*Proof.* Considering changing the whole local dataset, then  $t$  teachers will be influenced and the vote counts change by at most  $2t$ . Thus, from a party-level perspective, the sensitivity of the vote counts is  $2t$  and  $\mathcal{M}$  satisfies  $(2t\gamma, 0)$ -differential privacy.  $\square$

Like Theorem 8, we have the following theorem to track the moments.

**Theorem 10.** Let  $\mathcal{M}$  be the  $f$  function executed on the party side and  $q \geq \Pr[\mathcal{M}(d) \neq o^*]$  for some outcome  $o^*$ . Suppose the number of subsets in a partition is  $t$ . Let  $l, \gamma \geq 0$  and  $q < \frac{e^{2t\gamma}-1}{e^{4t\gamma}-1}$ . Then for any  $\mathbf{aux}$  and any two party-adjacent datasets  $d$  and  $d'$ ,  $\mathcal{M}$  satisfies

$$\alpha_{\mathcal{M}}(l; \mathbf{aux}, d, d') \leq \min(\log((1-q)\left(\frac{1-q}{1-e^{2t\gamma}q}\right)^l + qe^{2t\gamma l}), 2t^2\gamma^2 l(l+1)).$$

Note that the party-level privacy loss of FedKT-L2 can be quite large. To mitigate the impact of the noises, we usually expect  $t$  to be large to have a tighter bound of  $\Pr[\mathcal{M}(d) \neq o^*]$ . However, when  $t$  is large, the bound in Theorem 10 is also large. In fact, since every student model satisfies differential privacy, it is not necessary to apply party-level differential privacy in FedKT-L2.

## B Experiments

### B.1 Additional Details of Experimental Settings

We use *Adult*, *cod-rna*, *MNIST*, and *SVHN* for our experiments, where *Adult* and *cod-rna* are downloaded from this link<sup>4</sup>. The details of the datasets are shown in Table 4. We run experiments on a Linux cluster with 8 RTX 2080 Ti GPUs and 12 Intel Xeon W-2133 CPUs.

We compare FedKT with the other eight baselines. For the federated learning algorithms, the learning rate is tuned from  $\{0.001, 0.01\}$  and the number of local epochs is tuned from  $\{10, 20, 40\}$ . We found that the baselines can get best performance when the learning rate is set to 0.001 and the number of local epochs is set to 10. For FedProx, the regularization term  $\mu$  is tuned from  $\{0.1, 1\}$ . For SCAFFOLD, same as the experiments of the paper [Karimireddy and others, 2020], we use option II to calculate the control variates. For FedKT, SOLO, and PATE, the number of local epochs is simply set to 100. For PATE, the number of teacher models is set to be same as the number of parties of FedKT. We use the Adam optimizer and the  $L_2$  regularization is set to  $10^{-6}$ . The final parameters are summarized in Table 5.

### B.2 Hyper-parameters Study

#### Number of partitions in each party

Here we study the impact of the number of partitions (i.e., the parameter  $s$ ) on FedKT. Table 6 shows the test accuracy of FedKT with different  $s$ . From Table 6, we can see that the accuracy can be improved if we increase  $s$  from 1 to 2. However, if we further increase  $s$ , there is little or no improvement

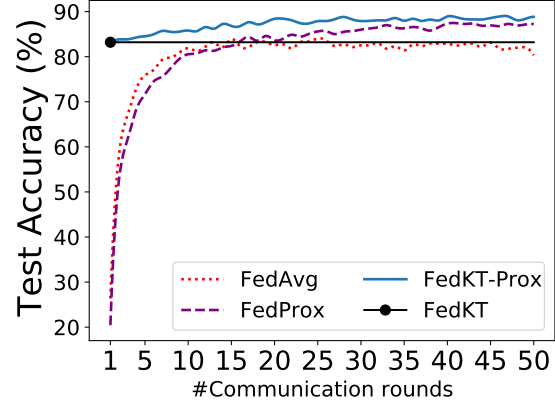


Figure 2: The test accuracy with increasing communication rounds/communication size on SVHN.

on the accuracy while the communication and computation overhead is larger. Thus, from our empirical study, we suggest users to simply set  $s$  to 2 for FedKT-L0 if they do not want to tune the parameters. For FedKT-L1 and FedKT-L2, since the privacy loss increases as  $s$  increases, we suggest to set  $s$  to small values. Users can simply set  $s$  to 1 or tune  $s$  from small values (i.e., 1 or 2) to find the best accuracy-privacy trade-off.

#### Number of teachers in each partition

Here we study the impact of number of teachers in each partition (i.e., the parameter  $t$ ). Table 7 shows the test accuracy of FedKT with different  $t$ . As we can see, FedKT can always get the best performance if setting  $t$  to 5. If  $t$  is large, the size of each data subset is small and the teacher models may not be good at predicting the public dataset. From our empirical study, users can simply set  $t$  to 5 if they do not want to tune the parameter. However, if the student models need to satisfy differential privacy (i.e., in FedKT-L2), the privacy loss may potentially be smaller if we increase  $t$  according to Lemma 7. Users need to tune  $t$  to find the best trade-off between the performance and the privacy loss.

### B.3 Extend to Multiple Rounds

FedKT can be used as an initialization step if applied to iterative training process. Here we combine FedKT and FedProx (denoted as FedKT-Prox) and compare it with FedAvg and FedProx. The results are shown in Figure 2. We can observe that FedKT-Prox can always achieve much higher accuracy than FedAvg and FedProx. Overall, FedKT-Prox is much more communication-efficient than FedAvg and FedProx.

### B.4 Differential Privacy

Table 8 and Table 9 present the results of FedKT-L1 and FedKT-L2. Besides the heterogeneous partition ( $\beta = 0.5$ ), we also try homogeneous partition (i.e., the dataset is randomly and equally partitioned into the parties). From the tables, we can see that the accuracy of FedKT-L1 and FedKT-L2 are comparable to the non-private version with a modest privacy budget. Moreover, the moments accountant method

<sup>4</sup><https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

Table 4: The datasets used and their learning models. The detailed model structures are shown in the paper.

| Datasets              | Adult               | cod-rna       | MNIST         | SVHN          |
|-----------------------|---------------------|---------------|---------------|---------------|
| #Training examples    | 24421               | 54231         | 50000         | 604288        |
| #Public examples      | 4070                | 9039          | 5000          | 13016         |
| #Test examples        | 4070                | 9039          | 5000          | 13016         |
| #Classes              | 2                   | 2             | 10            | 10            |
| #Parties (by default) | 50                  | 50            | 10            | 10            |
| Model                 | Random Forest       | GBDT          | MLP           | CNN           |
| Implemented library   | scikit-learn 0.22.1 | XGBoost 1.0.2 | PyTorch 1.6.0 | PyTorch 1.6.0 |

Table 5: The default parameters used in our experiments.

| Parameters |                                  | Adult | cod-rna | MNIST | SVHN  |
|------------|----------------------------------|-------|---------|-------|-------|
| common     | #parties                         | 50    | 50      | 10    | 10    |
|            | tree depth                       | 6     | 6       | —     | —     |
|            | learning rate                    | —     | 0.05    | 0.001 | 0.001 |
|            | batch size                       | —     | —       | 32    | 64    |
|            | #epochs                          | —     | —       | 10    | 10    |
| FedKT      | number of partitions in a party  | 2     | 2       | 2     | 2     |
|            | number of subsets in a partition | 5     | 5       | 5     | 5     |
| FedProx    | regularization term $\mu$        | —     | —       | 0.1   | 0.1   |

Table 6: The test accuracy of FedKT with number of partitions ranging between 1 and 5. We run 5 trials and report the mean and standard deviation. The number of subsets in each partition is set to 5 by default.

| #partitions | 1                | 2                       | 3                       | 4                | 5                |
|-------------|------------------|-------------------------|-------------------------|------------------|------------------|
| Adult       | 80.8% $\pm$ 1.4% | <b>82.2%</b> $\pm$ 0.6% | 81.5% $\pm$ 0.6%        | 81.2% $\pm$ 0.5% | 81.1% $\pm$ 0.1% |
| cod-rna     | 87.7% $\pm$ 0.6% | <b>88.3%</b> $\pm$ 0.6% | <b>88.3%</b> $\pm$ 0.5% | 88.1% $\pm$ 0.5% | 88.2% $\pm$ 0.5% |
| MNIST       | 89.2% $\pm$ 0.4% | <b>90.5%</b> $\pm$ 0.3% | 89.9% $\pm$ 0.2%        | 90.1% $\pm$ 0.2% | 90.2% $\pm$ 0.2% |
| SVHN        | 81.5% $\pm$ 0.6% | 83.2% $\pm$ 0.4%        | <b>83.5%</b> $\pm$ 0.4% | 83.5% $\pm$ 0.4% | 83.4% $\pm$ 0.3% |

Table 7: The test accuracy of FedKT with number of subsets in each partition ranging between 5 and 20. We run 5 trials and report the mean and standard deviation. For Adult, since there is a party with less than 15 examples, the experiment cannot successfully run when the number of subsets is not smaller than 15. For cod-rna, since there is a party with less than 20 examples, the experiment cannot successfully run when setting the number of subsets to 20. The number of partitions in each party is set to 2 by default.

| #subsets | 5                       | 10               | 15               | 20               |
|----------|-------------------------|------------------|------------------|------------------|
| Adult    | <b>82.0%</b> $\pm$ 0.6% | 81.1% $\pm$ 0.7% | —                | —                |
| cod-rna  | <b>88.3%</b> $\pm$ 0.6% | 87.4% $\pm$ 0.6% | 83.1% $\pm$ 0.6% | —                |
| MNIST    | <b>90.5%</b> $\pm$ 0.3% | 89.4% $\pm$ 0.2% | 89% $\pm$ 0.4%   | 88.4% $\pm$ 0.3% |
| SVHN     | <b>83.2%</b> $\pm$ 0.4% | 81.3% $\pm$ 0.5% | 80.0% $\pm$ 0.5% | 79.1% $\pm$ 0.6% |

usually can achieve a tighter privacy loss than the advanced composition [Dwork *et al.*, 2014]. For example, if we run cod-rna under homogeneous data partition setting  $\gamma = 0.1$  and the fraction of queries to 1%, the advanced composition gives us  $\varepsilon \approx 20.2$  and our analysis gives  $\varepsilon \approx 11.2$ . Note that the techniques in [Papernot and others, 2018] can also be applied to FedKT. For example, we may get a smaller privacy loss if adopting Gaussian noises instead of Laplace noises. Generally, our framework can also benefit from the state-of-the-art approaches on the privacy analysis of PATE, which we may investigate in the future.

Table 8: The accuracy and party-level  $\varepsilon$  of FedKT-L1 on Adult and cod-rna given different  $\gamma$  values and number of queries. For each setting, we run 3 trials and report the median accuracy and the corresponding  $\varepsilon$ . The number of partitions is set to 1 and the number of subsets in each partition is set to 5. The failure probability  $\delta$  is set to  $10^{-5}$ .

|         | data partitioning | $\gamma$ | #queries | $\varepsilon$ | acc          | FedKT-L0 acc |
|---------|-------------------|----------|----------|---------------|--------------|--------------|
| Adult   | Heterogeneous     | 0.04     | 0.1%     | 0.64          | 71.3%        | 82.2%        |
|         |                   |          | 0.5%     | 2.56          | 76.8%        |              |
|         |                   |          | 1%       | <b>4.73</b>   | <b>80.2%</b> |              |
|         |                   | 0.06     | 0.1%     | 0.96          | 75.7%        |              |
|         |                   |          | 0.5%     | 3.64          | 77.6%        |              |
|         |                   |          | 1%       | 5.78          | 80.2%        |              |
|         |                   | 0.08     | 0.1%     | 1.23          | 76.0%        |              |
|         |                   |          | 0.5%     | 4.25          | 76.3%        |              |
|         |                   |          | 1%       | 7             | 80.3%        |              |
|         | Homogeneous       | 0.02     | 0.1%     | 0.32          | 72.2%        | 82.4%        |
|         |                   |          | 0.5%     | 1.25          | 76.1%        |              |
|         |                   |          | 1%       | 1.79          | 80.1%        |              |
|         |                   | 0.04     | 0.1%     | 0.6           | 76.0%        |              |
|         |                   |          | 0.5%     | 1.9           | 80.4%        |              |
|         |                   |          | 1%       | 3.32          | 81.5%        |              |
|         |                   | 0.06     | 0.1%     | 0.75          | 76.1%        |              |
|         |                   |          | 0.5%     | 2.03          | 81.7%        |              |
|         |                   |          | 1%       | <b>3.36</b>   | <b>82.1%</b> |              |
| cod-rna | Heterogeneous     | 0.04     | 0.1%     | 1.09          | 66.8%        | 88.3%        |
|         |                   |          | 0.5%     | 3.54          | 72.5%        |              |
|         |                   |          | 1%       | 5.14          | 75.2%        |              |
|         |                   | 0.06     | 0.1%     | 1.52          | 69%          |              |
|         |                   |          | 0.5%     | 5.48          | 82.6%        |              |
|         |                   |          | 1%       | 8.1           | 79.5%        |              |
|         |                   | 0.1      | 0.1%     | 2.12          | 69%          |              |
|         |                   |          | 0.5%     | <b>6.89</b>   | <b>84.7%</b> |              |
|         |                   |          | 1%       | 11.2          | 85.3%        |              |
|         | Homogeneous       | 0.02     | 0.2%     | 0.53          | 67.0%        | 88.6%        |
|         |                   |          | 0.5%     | 1.71          | 73.2%        |              |
|         |                   |          | 1%       | 2.45          | 75.3%        |              |
|         |                   | 0.04     | 0.2%     | 1.5           | 73%          |              |
|         |                   |          | 0.5%     | 3.06          | 84.1%        |              |
|         |                   |          | 1%       | 5.10          | 85%          |              |
|         |                   | 0.06     | 0.2%     | 1.63          | 80.2%        |              |
|         |                   |          | 0.5%     | 3.10          | 84.1%        |              |
|         |                   |          | 1%       | <b>5.14</b>   | <b>86.1%</b> |              |

Table 9: The accuracy and example-level  $\varepsilon$  of FedKT-L2 on Adult and cod-rna given different  $\gamma$  values and number of queries. For each setting, we run 3 trials and report the median accuracy and the corresponding  $\varepsilon$ . The number of parties is set to 20 to ensure that FedKT has enough data to train each teacher model. The number of partitions is set to 1 and the number of subsets in each partition is set to 25. The failure probability  $\delta$  is set to  $10^{-5}$ .

|         | Data Partition | $\gamma$ | #queries | $\varepsilon$ | acc          | FedKT-L0 acc |
|---------|----------------|----------|----------|---------------|--------------|--------------|
| Adult   | Heterogeneous  | 0.04     | 0.1%     | 1.13          | 76.1%        | 82.4%        |
|         |                |          | 0.5%     | 2.56          | 76.5%        |              |
|         |                |          | 1%       | 3.72          | 78.5%        |              |
|         |                | 0.05     | 0.1%     | 1.32          | 76.1%        |              |
|         |                |          | 0.5%     | 3.24          | 79.0%        |              |
|         |                |          | 1%       | 4.76          | 79.2%        |              |
|         |                | 0.06     | 0.1%     | 1.96          | 76.2%        |              |
|         |                |          | 0.5%     | 3.93          | 78.5%        |              |
|         |                |          | 1%       | <b>5.79</b>   | <b>79.4%</b> |              |
|         | Homogeneous    | 0.04     | 0.3%     | 2.13          | 76.1%        | 82.6%        |
|         |                |          | 0.5%     | 2.59          | 78.7%        |              |
|         |                |          | 1%       | <b>3.72</b>   | <b>81.7%</b> |              |
|         |                | 0.06     | 0.3%     | 2.97          | 76.3%        |              |
|         |                |          | 0.5%     | 3.93          | 79.9%        |              |
|         |                |          | 1%       | 5.79          | 81.8%        |              |
|         |                | 0.08     | 0.3%     | 3.77          | 76.3%        |              |
|         |                |          | 0.5%     | 5.04          | 80.4%        |              |
|         |                |          | 1%       | 7.89          | 82.0%        |              |
| cod-rna | Heterogeneous  | 0.04     | 0.5%     | 3.54          | 77.7%        | 89.7%        |
|         |                |          | 1%       | 5.14          | 79.8%        |              |
|         |                |          | 2%       | 7.63          | 82.0%        |              |
|         |                | 0.05     | 0.5%     | 4.51          | 81.4%        |              |
|         |                |          | 1%       | 6.58          | 82.0%        |              |
|         |                |          | 2%       | <b>9.78</b>   | <b>84.7%</b> |              |
|         |                | 0.06     | 0.5%     | 5.50          | 81.2%        |              |
|         |                |          | 1%       | 8.10          | 83.2%        |              |
|         |                |          | 2%       | 12.2          | 85.9%        |              |
|         | Homogeneous    | 0.03     | 0.5%     | 2.64          | 79.2%        | 90.6%        |
|         |                |          | 1%       | 3.78          | 80.5%        |              |
|         |                |          | 2%       | 5.51          | 83.1%        |              |
|         |                | 0.04     | 0.5%     | 3.54          | 80.1%        |              |
|         |                |          | 1%       | 5.14          | 83.7%        |              |
|         |                |          | 1.5%     | 6.45          | 84.0%        |              |
|         |                | 0.05     | 0.5%     | 4.51          | 80.8%        |              |
|         |                |          | 1%       | 6.58          | 84.2%        |              |
|         |                |          | 1.5%     | <b>8.3</b>    | <b>84.7%</b> |              |