Secure Federated Submodel Learning

Chaoyue Niu¹, Fan Wu¹, Shaojie Tang², Lifeng Hua³, Rongfei Jia³, Chengfei Lv³, Zhihua Wu³, and Guihai Chen¹ Shanghai Key Laboratory of Scalable Computing and Systems, Shanghai Jiao Tong University, China

²Naveen Jindal School of Management, University of Texas at Dallas, USA

³Alibaba Group, Hangzhou and Beijing, China

Email: ¹{rvince, wu-fan}@sjtu.edu.cn; gchen@cs.sjtu.edu.cn; ²shaojie.tang@utdallas.edu;

³{issac.hlf, rongfei.jrf, chengfei.lcf, zhihua.wzh}@alibaba-inc.com

Abstract—Federated learning was proposed with an intriguing vision of achieving collaborative machine learning among numerous clients without uploading their private data to a cloud server. However, the conventional framework requires each client to leverage the full model for learning, which can be prohibitively inefficient for resource-constrained clients and large-scale deep learning tasks. We thus propose a new framework, called federated submodel learning, where clients download only the needed parts of the full model, namely submodels, and then upload the submodel updates. Nevertheless, the "position" of a client's truly required submodel corresponds to her private data, and its disclosure to the cloud server during interactions inevitably breaks the tenet of federated learning. To integrate efficiency and privacy, we have designed a secure federated submodel learning scheme coupled with a private set union protocol as a cornerstone. Our secure scheme features the properties of randomized response, secure aggregation, and Bloom filter, and endows each client with a customized plausible deniability, in terms of local differential privacy, against the position of her desired submodel, thus protecting her private data. We further instantiated our scheme with the e-commerce recommendation scenario in Alibaba, implemented a prototype system, and extensively evaluated its performance over 30-day Taobao user data. The analysis and evaluation results demonstrate the feasibility and scalability of our scheme from model accuracy and convergency, practical communication, computation, and storage overheads, as well as manifest its remarkable advantages over the conventional federated learning framework.

Index Terms—federated submodel learning, private set union, randomized response, local differential privacy, secure aggregation, Bloom filter, e-commerce recommendation

I. INTRODUCTION

A. Motivating Industrial Scenario in Alibaba

The industrial scenario in Alibaba that motivated federated submodel learning is the desire to provide customized and accurate e-commerce recommendations for billion-scale clients while keeping user data on local devices.

This work was supported in part by Science and Technology Innovation 2030 – "New Generation Artificial Intelligence" Major Project No. 2018AAA0100905, in part by China NSF grant 61972252, 61972254, 61672348, and 61672353, in part by the Open Project Program of the State Key Laboratory of Mathematical Engineering and Advanced Computing 2018A09, and in part by Alibaba Group through Alibaba Innovation Research (AIR) Program. The opinions, findings, conclusions, and recommendations expressed in this paper are those of the authors and do not necessarily reflect the views of the funding agencies or the government.

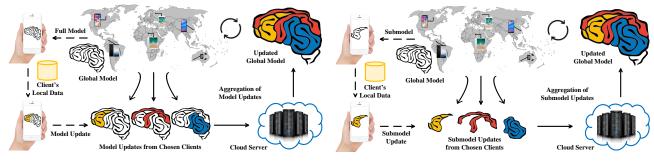
F. Wu is the corresponding author.

Currently, the recommendation systems in Alibaba are cloud based and require the server cluster to collect, process, and store numerous user data. In addition, the deployed recommendation models follow a golden paradigm of embedding¹ and Multi-Layer Perceptron (MLP) [3]: user data are first encoded into high-dimensional sparse feature vectors, then embedded into low-dimensional dense vectors, and finally fed into fully connected layers. To improve accuracy, Deep Interest Network (DIN) [4] introduces the attention mechanism to activate the user's historical behaviors, namely relative interests, with respect to the target item; Deep Interest Evolution Network (DIEN) [5] further extracts latent interests and monitors interest evolution through Gated Recurrent Unit (GRU) coupled with attention update gate; and Behavior Sequence Transformer (BST) [6] incorporates transformer to capture the sequential signals underlying the user's behavior sequence.

However, typical fields of user data involved in recommendation include user profile (e.g., user ID, gender, and age), user behavior (e.g., the list of visited goods IDs and relevant information, such as category IDs and shop IDs), and context (e.g., time, page number, and display position). More or less, these data fields are sensitive, and some clients who value security and privacy highly may refuse to share their data. In addition, according to the General Data Protection Regulation (GDPR), which was legislated by the European Commission and took effect on May 25, 2018, any institution or company is prohibited from uploading user data and storing it in the cloud without the explicit permissions from the European Union users [7], [8]. Under such circumstances, refining the recommendation models and further providing accurate recommendations become urgent demands as well as thorny challenges in practice.

Federated learning, which decouples the ability to do machine learning from the need to upload and store data in the cloud, is a potential solution. However, the original framework of federated learning, proposed by Google researchers in [9], requires each client to download the full machine learning

¹In general, deep learning with a huge and sparse input space (e.g., e-commerce goods IDs, natural language texts, and locations) requires an embedding layer to first transform inputs into a lower-dimensional space [1]. Additionally, the full embedding matrix tends to occupy a large proportion of the whole model parameters (e.g., 98.22% in our evaluated DIN model and more than two-thirds in Gboard's CIFG language model [2]).



(a) Conventional Federated Learning

(b) Federated Submodel Learning

Fig. 1. A visual comparison of conventional federated learning and our new federated submodel learning.

model for training and then to upload the update of the full model, which is impractical for resource-constrained clients in the context of complex deep learning tasks. For example, as the largest online consumer-to-consumer platform in China, Taobao (owned by Alibaba) has roughly two billion goods in total [10], which is far larger than the 10,000 word vocabulary in the natural language scenario of Google's Gboard [2], [11], [12]. This implies that the full embedding matrix of goods has roughly two billion rows and roughly occupies 134GB of space, when the embedding dimension is 18, and each element adopts 32-bit representation. If each client directly uses the full matrix for learning, it inevitably incurs huge overheads, which are unacceptable and unaffordable for one billion Taobao users with smart devices. To improve efficiency, we observe that a certain user tends to browse, click, and buy a small number of goods, and thus just needs a tailored model, which can sharply reduce the overheads and is more practical for mobile clients. Continuing with the above example, if a Taobao user's historical data involve 100 goods, she only requires the corresponding 100 rows, rather than the entire two billion rows, of the embedding matrix. Based on this key observation, we propose a new framework of federated learning, called federated submodel learning, as follows.

B. Framework of Federated Submodel Learning

We plot the workflow of federated submodel learning in Fig. 1(b) and also provide the traditional federated learning in Fig. 1(a) for an intuitive comparison.

In the beginning of one communication round, a cloud server first selects a certain number of eligible clients, typically end users whose mobile devices are idle, charging, and connected to an unmetered Wi-Fi network. This eligibility criteria is used to avoid a negative effect on the user experience, data usage, or battery life. Then, each chosen client downloads part of the global model as she requires, namely a submodel, from the cloud server. For example, in the e-commerce scenario above, a client's submodel mainly consists of the embedding parameters for the displayed and clicked goods in her historical data, as well as the parameters of the other network layers. Afterwards, the client trains the submodel over her private data locally. At the end of one round, the cloud server lets those chosen clients who are still alive upload the updates of

their submodels and further aggregates the submodel updates to form a consensus update to the global model. Considering the convergencies of the global model at the cloud server and the submodels on clients, the above process is iterated for several rounds.

If each client leverages the full model rather than her required submodel for learning, federated submodel learning will degenerate to conventional federated learning. Compared with the conventional one, our new framework further decouples the ability to accomplish federated learning from the need to use the prohibitively large full model, which can dramatically improve efficiency. For example, in our evaluation, the size of a client's desired submodel is only 1.99% of the full model's size. Thus, our framework is more practical for resource-constrained clients and deep learning tasks.

C. Newly Introduced Privacy Risks

Just as every coin has two sides, federated submodel learning not only brings in efficiency but also introduces extra privacy risks. On one hand, compared with using the public full model in conventional federated learning, the download of a submodel and the upload of the submodel update would require each client to provide an index set as auxiliary information, specifying the "position" of her submodel. However, the index set normally corresponds to the client's private data. For example, to specify the required rows of the embedding matrix in the e-commerce scenario, a client mainly needs to provide the goods IDs in her user data as the index set. Thus, the disclosure of a client's real index set to the cloud server can still be regarded as the leakage of the client's private data, breaking the tenet of federated learning. On the other hand, compared with the aligned full model in federated learning, each client only submits the update of her customized and highly differentiated submodel in federated submodel learning. As a result, the aggregation of updates with respect to a certain index can come from a unique client (e.g., with probability 86.7% for 100 chosen clients in our evaluated Taobao dataset), which indicates that the cloud server not only can ascertain that the client has a certain index but also can learn her detailed update. These two kinds of knowledge both breach the client's private data. Further, such a privacy risk in e-commerce is more severe than that in natural language because compared with the vocabularies of different Gboard users, the goods IDs of different Taobao users are more differentiated. We will detail and visualize the preceding privacy risks in Section III-A and Fig. 2.

D. Fundamental Problems and Challenges

In essence, to mitigate the above privacy risks, we need to jointly solve two fundamental problems modeled from the processes of downloading a submodel and uploading a submodel update, respectively. One is how a client can download a row from a matrix, maintained by an untrusted cloud server, without revealing which row, alternatively the row index, to the cloud server. The other is how a client can modify a row of the matrix, still without revealing which row was modified and the altered content to the cloud server. Using the terminology from file system permissions, the first problem has a "readonly" attribute, where the client only reads the file. In contrast, the second problem is in a "write" mode, where the client can edit the file. Further incorporating the obscure requirement of two operations, the second problem appears more challenging than the first one. We now analyze these two problems in detail as follows.

We start with the first problem. One trivial method is that the client downloads the full matrix, as in conventional federated learning, and then extracts the required row locally. Although this method perfectly hides the fetched row index, it incurs significant communication cost, which can be unaffordable for resource-constrained mobile devices, especially when the matrix is huge, e.g., representing a deep neural network. To avoid downloading the full matrix, Private Information Retrieval (PIR) [13]–[15] can be applied, which exactly matches our problem settings, including the read-only mode and the privacy preservation requirement of the retrieved elements. The stateof-the-art constructions of private information retrieval include Microsoft's SealPIR [13] and Labeled PSI [14] and Goolge's PSIR [15], where two Microsoft protocols have been deployed in its Pung private communication system [16]. We note that another celebrated cryptographic primitive, called Oblivious Transfer (OT) [17], is stronger than private information retrieval. It not only guarantees that the cloud server does not know which row the client has downloaded, as in private information retrieval, but also ensures that the client does not know the other rows of the matrix, which is instead not needed in practical federated submodel learning. Therefore, if we consider the first problem independently, private information retrieval may be a good choice.

We next dissect the second problem. For a concrete row of the full matrix, if clients modify this row one by one, the cloud server definitely knows those clients who modified this row and their detailed contents of modification. Thus, one feasible way is to first securely aggregate all the modifications without revealing any individual modification, and then apply the aggregate modification to the row of the full matrix once. In particular, such a guarantee can be provided by the secure aggregation protocol in [18] and some other schemes for oblivious addition, e.g., based on additively homomorphic

cryptosystems [19]-[21]. With the secure aggregation guarantee, if more than one client participates in aggregation and at least one of their modifications is nonzero, then the cloud server cannot reveal which client(s) truly intend to modify this row and their detailed modifications. Further, a larger number of involved clients implies a stronger privacy guarantee. One extreme case is in conventional federated learning, which harshly lets all chosen clients in one communication round be involved, no matter whether they truly intend to modify this row or not. Thus, it can offer the best privacy guarantee. Nevertheless, considering each client needs to be involved for each row of the full matrix, it is too inefficient to be applicable in the large-scale deep learning context. Another extreme case is in federated submodel learning, which simply lets those clients who really intend to modify this row be involved. Hence, each chosen client only needs to be involved for those rows that she truly intends to modify, implying the best efficiency. However, different clients tend to modify highly differentiated or even mutually exclusive rows. For the joint modification with respect to some row, chances are high that only one client is involved. Under such a circumstance, the secure aggregation guarantee no longer works, which leaks the client's real intention and her detailed modification. In a nutshell, trivial solutions to the second problem cannot well balance or support tuning privacy and efficiency.

E. Our Solution Overview and Major Contributions

Jointly considering the above two fundamental problems and several practical issues, we propose a secure scheme for federated submodel learning. In our scheme, each chosen client generates three types of index sets locally; real, perturbed, and succinct. First, the real index is extracted from a client's private data and is kept secret from the other system participants, including the cloud server and any other chosen client. Second, the perturbed index set is used to interact with others in the download and upload phases. It is generated by applying randomized response twice with one memoization step between. Such a design, together with secure aggregation, allows the client to hold a self-controllable deniability against whether she really intends or does not intend to download some row and to upload the modification of this row, even if the client may be chosen to participate in multiple communication rounds. The strength of deniability is rigorously quantified using local differential privacy. Further, rather than trivially using the prohibitively large-scale full index set as the questionnaire of randomize response in every communication round, we identify a necessary and sufficient index set, namely the union of the chosen clients' real index sets. Considering the secrecy of each client's real index set, we propose an efficient and scalable Private Set Union (PSU) protocol based on Bloom filter, secure aggregation, and randomization, allowing clients to obtain the union under the mediation of an untrusted cloud server without revealing any individual real index set. In particular, private set union promises a wide range of applications but receives little attention. Due to unaffordable overheads, none of the existing protocols can be deployed in

practice yet. Last, the succinct index set is derived from the intersection between the real and perturbed index sets, and it is used to prepare the data and submodel for local training.

We summarize our key contributions in this work as follows:

- To the best of our knowledge, we are the first to propose the framework of federated submodel learning and further to identify and remedy new privacy risks.
- Our proposed secure scheme mainly features the properties of randomized response and secure aggregation to empower each client with a tunable deniability against her real intention of downloading the desired submodel and uploading its update, thus protecting her private data. As a moat, we designed an efficient and scalable private set union protocol based on Bloom filter and secure aggregation, which can be of independent and significant value in practice.
- We instantiated with Taobao's e-commerce scenario, adopted Deep Interest Network (DIN) for recommendation, and implemented a prototype system. Additionally, we extensively evaluated over one month of Taobao data. The evaluation and analysis results demonstrate the practical feasibility of our scheme, as well as its remarkable advantages over the conventional federated learning framework in terms of model accuracy and convergency, communication, computation, and storage overheads. Specifically, when the number of chosen clients in one round is 100, compared with conventional federated learning, which diverges in the end, our scheme improves the highest Area Under the Curve (AUC) by 0.072. In addition, at the same security and privacy levels as conventional federated learning with secure aggregation, our scheme reduces 80.05% of communication overhead on both sides of the client and the cloud server. Moreover, our scheme reduces 85.02% (resp., 45.43%) and 72.51% (resp., 63.77%) of computation (resp., memory) overheads on the sides of the client and the cloud server, respectively. Furthermore, when the size of the full model scales further, it does not incur additional overhead to our scheme, but it prohibits conventional federated learning from being applied. Finally, for our private set union, the communication overhead per client is less than 1MB, and the computation overheads of the client and the cloud server are both less than 40s, even if the dropout ratio of the chosen clients reaches 20%.

II. RELATED WORK

In recent years, federated learning has become an active topic in both academic and industrial fields. In this section, we briefly review some major focuses and relevant work as follows. For more related work, we direct interested readers to the surveys written by Li et al. [22] and Yang et al. [23].

First and most important is to identify and address security and privacy issues of federated learning. Bonawitz et al. [18] proposed a secure, communication-efficient, and failure-robust aggregation protocol in both honest-but-curious and active adversary settings. It can ensure that the untrusted cloud server

learns nothing but the aggregate (or mathematically, the sum) of the model updates contributed by chosen clients, even if part of clients drop out during the aggregation process. To bound the leakage of a certain client's training data from her individual model update, several differentially private mechanisms were proposed. McMahan et al. [24] offered clientlevel differential privacy for recurrent language models based on the celebrated moments account scheme in [25]. Here, the moments account allows the release of all intermediate results during the training process, particularly the gradients per iteration; keeps track of privacy loss in every iteration; and provides a tighter compositive/cumulative privacy guarantee. However, in the practical federated learning scenario, only the model update after multiple iterations/epochs is revealed. whereas all intermediate gradients are hidden. Specific to this case, Feldman et al. [26] analyzed the detailed amplification effect of hiding intermediate results on differential privacy. In contrast to these defense mechanisms, Bagdasaryan et al. [27] developed a model replacement attack launched by malicious clients to backdoor the global model at the cloud server. Melis et al. [1] exploited membership and property inference attacks to uncover features of the clients' training data from model updates.

Second is to improve the communication efficiency, especially the expensive and limited up-link bandwidth for mobile clients. To overcome this bottleneck, two types of solution methods have been proposed in general. One is to reduce the total number of communication rounds between the cloud server and the clients. A pioneering work is the federated averaging algorithm proposed by McMahan et al. [9]. Its key principle is to let each client locally train the global model for multiple epochs, and then upload the model update. Thus, it is more communication efficient than the common practice of conventional distributed learning to exchange gradients per iteration in datacenter-based scenarios. The other complementary way is to further reduce the size of the transmitted message in each communication round, particularly through compressing model updates. Typical compression techniques include sparsification, subsampling, and probabilistic quantization coupled with random rotation. For example, after quantization, the original float-type elements of the update of the global model can be encoded as integer-type values with a few bits [28], [29]. Considering the compressed model updates are discrete, while classic differentially private deep learning mechanisms, hinging on the Gaussian mechanism, only support continuous inputs, Agarwal et al. [30] proposed a Binomial mechanism to guarantee differential privacy for one iteration while enjoying communication efficiency. Another effective approach to improving communication efficiency is to first apply dropout strategies to the global model, and then let clients train over the same reduced model architecture [31]. As a result, the downloaded model and the uploaded model update can be compressed in terms of dimension.

Third is from learning theory. The federated learning framework has several atypical characteristics: non independent and identically distributed (non-iid) and unbalanced data dis-

tributed over numerous clients with limited and unstable network connections. Such statistical heterogeneity and existence of stragglers make most existing analysis techniques for iid data infeasible and pose significant challenges for designing theoretically robust and efficient learning algorithms. The federated averaging algorithm mentioned above, as a cornerstone of federated learning, empirically shows its effectiveness in some tasks, but was observed to diverge for a large number of local epochs in [9]. More specifically, it lets multiple chosen clients run mini-batch Stochastic Gradient Descent (SGD) in parallel, and then lets the cloud server periodically aggregate the model updates in a weighted manner, where weights are proportional to the sizes of the clients' training sets. Recently, Yu et al. [32] and Li et al. [33] advanced the convergency analysis of federated averaging by imposing smooth and bounded assumptions on the loss function. The follow-up work [34] further presented a momentum extension of parallel restarted SGD, which is compatible with federated learning. Different from the above work, Smith et al. [35] focused on learning separate but related personalized models for distinct clients by leveraging multitask learning for shared representation. Chen et al. [36] instead adopted meta-learning to enable client-specific modeling, where clients contribute information at the algorithm level rather than the model level to help train the meta-learner. Mohri et al. [37] considered an unfairness issue that the global model can be unevenly biased toward different clients. They thus proposed a new agnostic federated learning framework where the global model can be optimized for any possible target distribution, which is formed via a mixture of client distributions. Eichner et al. [38] captured data heterogeneity in federated learning, particularly cyclic patterns, and offered a pluralistic solution for convex objectives and sequential SGD.

Fourth is regarding production and standardization. Google has deployed federated learning in its Android keyboard, called Gboard, to polish several language tasks, including nextword prediction [2], query suggestion [39], out-of-vocabulary words learning [11], and emoji prediction [12]. In particular, the query suggestion used logistic regression as the triggering model for on-device training to determine whether the candidate suggestion should be shown or not. In addition, the other three tasks leveraged a tailored Long Short-Term Memory (LSTM) Recurrent Neural Network (RNN), called Coupled Input and Forget Gate (CIFG). In [40], Google's team also detailed their initial system design and summarized practical deployment issues, such as irregular device availability, unreliable network connectivity and interrupted execution, orchestration of lock-step execution across heterogonous devices, and limited device storage and computation resources. They also pointed out some future optimization directions in bias, convergence time, device scheduling, and bandwidth. To facilitate open research, Google integrated a federated learning simulation interface into its deep learning framework, called TensorFlow Federated [41]. However, this open-source module lacks several core functionalities, e.g., secure and privacy preserving mechanisms, on-device training, socket commu-

TABLE I
FREQUENTLY USED NOTATIONS AND ABBREVIATIONS.

Notation	Remark
$\mathbf{W}^{m imes d}$	Global/Full model at the cloud server, denoted by a
	matrix with m rows and d columns
$S = \{1, 2, \dots, m\}$	Full row index set of W
C, C = n	The set of n clients chosen by the cloud server in one
<u>^</u>	communication round, the cardinality of ${\mathcal C}$
$\hat{\mathcal{C}} \subset \mathcal{C}$	The up-to-date set of clients who are alive throughout
	the communication round
$i \in \mathcal{C}$	A chosen client i
$\mathcal{S}^{(i)} \subset \mathcal{S}$	Client <i>i</i> 's real index set that corresponds to local data and specifies truly required rows of W
$S''^{(i)}$	A perturbed index set of client i, to download the
	submodel from the cloud server and to securely upload
	the update of the submodel to the cloud server
$\mathbf{W}_{\mathcal{S}''(i)}$	Client i's downloaded submodel
$\mathbf{W}_{\mathcal{S}^{(i)} \cap \mathcal{S}^{\prime\prime(i)}}^{\mathcal{S}^{\prime\prime(i)}}$	Client i's succinct submodel for local training
$\Delta \mathbf{W}_{\mathcal{S}^{(i)} \cap \mathcal{S}^{\prime\prime(i)}}$	Client i's succinct submodel update
$\Delta \mathbf{W}_{S''(i)}$	Client i's uploaded submodel update by padding zero
3 . ,	vectors to the succinct submodel update
ϵ	A privacy level/budget of local differential privacy
$\mathbf{b}, \beta, h, \phi$	A Bloom filter with β bits and h hash functions,
	representing/accommodating a set of ϕ elements
	Dimension of vector in the secure aggregation protocol
$p_1^{(i)}, p_2^{(i)}, p_3^{(i)}, p_4^{(i)} \\ p_5^{(i)} = p_1^{(i)}(p_3^{(i)} - p_4^{(i)}) + p_4^{(i)}$	Client i's probability parameters to generate $S''^{(i)}$
$p_5^{(i)} = p_1^{(i)}(p_3^{(i)} - p_4^{(i)}) + p_4^{(i)}$	The probability that an index in client i's real index
	set will fall into her perturbed index set
$p_6^{(i)} = p_2^{(i)}(p_3^{(i)} - p_4^{(i)}) + p_4^{(i)}$	The probability that an index not in client i 's real index
	set will fall into her perturbed index set
p_7	The probability of the cloud server ascertaining that
	an index belongs to some client's real index set and
	also learning her detailed update with respect to this
	index from the securely aggregated submodel update
p_8	The probability of the cloud server ascertaining that
	an index does not belong to some client's real index set from the securely aggregated submodel update
0	The expected cardinality of each client's real index set
$\mathbb{Z}_R = \{0, 1, \dots, R - 1\}$	The least residue system modulo R
γ	A level of stochastic quantization mechanism
r FL	Conventional federated learning
SFL	Secure federated learning, namely conventional feder-
	ated learning with secure aggregation
SFSL	Secure federated submodel learning
CPP	Choice of probability parameters
SA	Secure aggregation

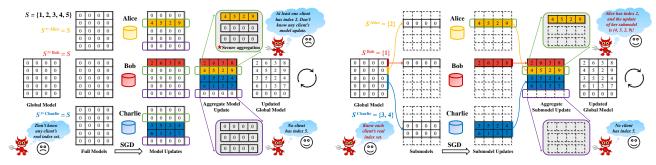
nication between cloud server and clients, task scheduling, and dropout/exception handling. These significantly suppress federated learning related productions in other commercial companies. Caldas et al. [42] released a benchmark for federated learning, called LEAF. Currently, LEAF comprises some representative datasets, evaluation metrics, and a referenced implementation of federated averaging.

Parallel to existing work, where clients use the same (simplified) global model for learning, we propose a novel federated submodel learning framework for the sake of scalability. Under this framework, we identify and remedy new security and privacy issues, due to the dependence between the position of a client's desired submodel and her private data as well as the misalignment of clients' submodel updates in aggregation.

III. PRELIMINARIES

In this section, we elaborate on the privacy risks sketched in Section I-C and formally define the corresponding security requirements. We also review some existing building blocks.

We first introduce some necessary notations. For the sake of clarity, frequently used notations and abbreviations throughout this paper are also listed in Table I. We use a two-dimensional



(a) Federated Learning with Secure Aggregation

(b) Federated Submodel Learning with Secure Aggregation

Fig. 2. Federated submodel learning with secure aggregation can still leak a concrete client's real indices and her detailed updates to the cloud server, compared with conventional federated learning with secure aggregation. The rounded rectangle colored in dark gray denotes the process of secure aggregation [18], where the cloud server, as the aggregator, only obtains the sum of vectors from multiple clients but does not know any individual client's vector. Additionally, the row in a table with dashed lines indicates that the client does not download, train over, or upload this row.

matrix with m rows and d columns to represent the global/full model, denoted as \mathbf{W} . Such a matrix-based representation not only suffices for the recommendation models used in Alibaba but also can easily degenerate to a widely used vector-based representation [18], [30], by setting the number of columns d to 1. Additionally, we let $\mathcal{S} = \{1, 2, \ldots, m\}$ denote the entire row index set of \mathbf{W} . Moreover, we let \mathcal{C} denote those clients who are selected by the cloud server to participate in one communication round of federated submodel learning. For a chosen client $i \in \mathcal{C}$, we let $\mathcal{S}^{(i)} \subset \mathcal{S}$ denote her real index set, which implies that the user data of client i involves the rows in \mathbf{W} with indices $\mathcal{S}^{(i)}$.

A. Details on Privacy Risks and Security Requirements

We now expand on two kinds of privacy leakages that the federated submodel learning brings in, compared with conventional federated learning. We provide Fig. 2 for illustration. We here adopt an *honest-but-curious* security model, in which the cloud server and all clients follow the designed protocol, but try to glean sensitive information about others.

The first kind of privacy leakage is the disclosure of a client's real index set, which specifies the position of a submodel and implies the client's private data, to the cloud server. For example, each row of the embedding matrix for goods in the recommendation model is linked with a certain goods ID, which indicates that a client's real index set, specifying her required rows of the embedding matrix, is in fact the goods IDs in her private data. Similarly, when federated submodel learning is applied to the natural language scenario (e.g., nextword prediction in Gboard), a client's real index set to locate her wanted parameters of word embedding is actually the vocabulary extracted from her typed texts. Thus, the disclosure of a client's real index set to the cloud server is still regarded as the leakage of the client's private data. In contrast, for conventional federated learning, each client essentially uses the full index set, which is public to the cloud server and all other clients, and does not reveal any private information.

The second kind of privacy leakage is from the aggregation of misaligned submodel updates, where the cloud server may not only know that a certain client has a concrete index but also learn her detailed update with respect to this index. In addition to the fact that the real index reveals a client's private data, the client's individual submodel update can still memorize or even allow reconstruction of her private data, namely "model inversion" attack [1], [43]-[46]. To conceal a client's individual update in conventional federated learning, the secure aggregation protocol [18] can be applied, which allows the cloud server to obtain the sum of multiple vectors without learning any individual vector. As shown in Fig. 2(a), with respect to index 2, Alice submits the update, denoted by the vector (4,5,2,9), whereas Bob and Charlie submit two zero vectors. The secure aggregation protocol can guarantee that the cloud server only obtains the sum of three vectors, i.e., (4,5,2,9), but does not know the content of any individual vector. This further implies that from the aggregate result, the cloud server can merely infer that at least one client has index 2, but cannot identify which client(s). Such a functionality is essentially analogous to anonymization. In a nutshell, the zero updates from Bob and Charlie function as two shields of Alice. However, in federated submodel learning, due to the differentiation and misalignment of clients' submodels, the "zero" shields from other clients vanish, and the aggregation of updates with respect to a certain index can come from one unique client, making secure aggregation ineffective. For example, in Fig. 2(b), only Alice who has index 2 submits the update (4, 5, 2, 9), whereas Bob and Charlie submit nothing. Without the blindings from Bob and Charlie, the cloud server not only knows that Alice has index 2 while Bob and Charlie do not have but also learns Alice's detailed update (4, 5, 2, 9).

Given the two kinds of privacy leakages above, we define the corresponding security requirements. First, for the disclosure of real index sets when clients interact with the cloud server, we consider that each client should have *plausible deniability* of whether a certain index is or is not in her real index set. To measure the strength of plausible deniability, we adopt local differential privacy, which is a variant of standard differential privacy in the local setting. Specifically, the perturbation in local differential privacy is performed by clients

in a distributed manner, rather than relying on a data curator, as a trusted authority to conduct centralized perturbation in differential privacy. Thus, the privacy of an individual client's data is not only preserved from external attackers but also from the untrusted data curator, e.g., the cloud server in our context. Due to its intriguing security properties, local differential privacy for various population statistics has recently received significant industrial deployments (e.g., in Google [47], [48], Apple [49], and Microsoft [50]), as well as lasting academic attention [51]–[61]. We now present the formal definition of local differential privacy as follows:

Definition 1 (Local Differential Privacy). A randomized mechanism M satisfies ϵ -local differential privacy, if for any pair of inputs from a client, denoted as x and y, and for any possible output of M, denoted as z, we have

$$\frac{\Pr\left(M\left(x\right)=z\right)}{\Pr\left(M\left(y\right)=z\right)} \le \exp\left(\epsilon\right),\,$$

where ϵ is a privacy budget controlled by the client. A smaller ϵ offers a better privacy guarantee.

Intuitively, the above definition says that the output distribution of the randomized mechanism does not change too much, given distinct inputs from the client. Thus, local differential privacy formalizes a sort of plausible deniability: no matter what output is revealed, it is approximately equally as likely to have come from one input as any other input. In addition, when local differential privacy applies to obscure the membership of a certain index in federated submodel learning, the inputs and the outputs are boolean values, where possible inputs (resp., outputs) are two states: a certain index "in" or "not in" a client's real (resp., revealed) index set. Moreover, we can check that conventional federated learning provides the strongest deniability, where the level of local differential privacy is $\epsilon = \ln(1/1) = 0$ for each client. The reason is that no matter whether an index is or is not in a client's real index set (different inputs), this index will definitely be revealed (the same output). In contrast, federated submodel learning provides the weakest deniability, where the level of local differential privacy is $\epsilon = \ln(1/0) = \infty$ for each client, because if an index is in (resp., not in) a client's real index set (different inputs), this index will definitely (resp., definitely not) be revealed, i.e., the output with probability 1 (resp., 0).

Second, direct secure aggregation of submodel updates is the most efficient but insecure case, which can leak whether some client has a certain index as well as her detailed update. In contrast, the other extreme case is conventional federated learning with secure aggregation, which is most secure but inefficient. Specifically, all participating clients upload the full model updates, which can perfectly prevent privacy leakages due to the misalignment of customized submodels. To enable clients to tune privacy and efficiency in a fine-grained manner, we define a client-controllable privacy protection mechanism for submodel updates aggregation.

Definition 2. A privacy protection mechanism for submodel updates aggregation is client controllable, if it enables partic-

ipating clients to determine the probabilities of the following two complementary events: From the securely aggregated submodel update,

- Event 1: the cloud server ascertains that an index belongs to some client's real index set and also learns her detailed update with respect to this index;
- Event 2: the cloud server ascertains that an index does not belong to some client's real index set.

We note that revealing the states of some clients having and not having a certain index should both be regarded as privacy leakages. Furthermore, when the above definition applies to federated learning, and if at least two clients participate in aggregation, the probability of Event 1 is 0, and the probability of Event 2 is still 0 for those indices within the union of the chosen clients' real index sets. For an index outside the union, e.g., index 5 shown in Fig. 2(a), the probability of Event 2 is approaching 1. The reason is that from the aggregate zero vector, the cloud server almost ascertains that all clients do not have this index, despite of some rare cases (e.g., Alice and Bob submit two vectors of elements differing in signs, and Charlie submits a zero vector).

B. Building Blocks

We review randomized response, secure aggregation, and Bloom filter underlying our design.

1) Randomized Response: Randomized response, due to Warner in 1965 [62], is a survey technique in the social sciences to collect statistical information about illegal, embarrassing, or sensitive topics, where the respondents want to preserve privacies of their answers. A classical example for illustrating this technique is the "Are you a member of the communist party?" question. For this question, each respondent flips a fair coin in secret and tells the truth if it comes up tails; otherwise, she flips a second coin and responds "Yes" if heads and "No" if tails. Thus, a communist (resp., non-communist) will answer "Yes" with probability 75% (resp., 25%) and "No" with probability 25% (resp., 75%).

The intuition behind randomized response is that it provides plausible deniability for both "Yes" and "No" answers. In particular, a communist can contribute her response of "Yes" to the event that the first and second coin flips were both heads, which occurs with probability 25%. Meanwhile, a noncommunist can also contribute her response of "No" to the event that the first coin is heads and the second coin is tails, which still occurs with probability 25%. Furthermore, the plausible deniability of randomized response can be rigorously quantified by local differential privacy. As analyzed in [47], [63], for a one-time response, each respondent has local differential privacy at the level $\epsilon = \ln(75\%/25\%) = \ln 3$, irrespective of any attacker's prior knowledge.

2) Secure Aggregation: An individual model update may leak a client's private data under the notorious model inversion attack. Nevertheless, to update the global model in federated learning, the cloud server does not need to access any individual model update and only requires the aggregate,

basically the sum, of multiple model updates. For example, if n clients participate in the aggregation, denoted as C, where client $i \in \mathcal{C}$ holds a vector $\Delta \mathbf{w}^{(i)} \in \mathbb{Z}^l$ of dimension l, the cloud server should just obtain the sum $\sum_{i\in\mathcal{C}}\Delta\mathbf{w}^{(i)}$, while maintaining each individual $\Delta\mathbf{w}^{(i)}$ in secret. For this purpose and the characteristics of mobile devices, particularly limited and unstable network connections and common dropouts, Google researchers proposed a secure aggregation protocol in [18]. Implied by the functionality of oblivious addition, secure aggregation in federated learning can further ensure that even if the model inversion attack succeeds, the attacker (e.g., the honest-but-curious or active adversary cloud server, or an external intruder) may only infer that a group of clients has a certain data item but cannot identify which concrete client. This functionality is similar to anonymization. In what follows, we briefly review the secure aggregation protocol from communication settings, technical intuitions, scalability, and efficiency.

First, we introduce its communication settings. During the aggregation process, a client can neither establish direct communication channels with other clients nor natively authenticate other clients. However, each client has a secure (private and authenticated) channel with the cloud server. Thus, if one client intends to exchange messages with other clients, she needs to hinge on the cloud server as a relay. In addition, to guarantee confidentiality and integrity against the mediate cloud server, client-to-client messages should be encrypted with symmetric authenticated encryption, where the secret key is set up through Diffie-Hellman key exchange between two clients. Moreover, to defend active adversaries, a digital signature scheme is required for consistency checks. These basic settings make the secure aggregation protocol different from other relevant work about oblivious addition [19]–[21], or, more generally, secure multiparty computation [64]-[67], which requires direct peer-to-peer communication between clients; assumes the availability of multiple noncolluding cloud servers; or resorts to a trusted third party for key generation and distribution.

Second, we outline the technical intuitions behind secure aggregation. Each client doubly masks her private data, including a self mask and a mutual mask. Here, the self mask is chosen by the client, whereas the mutual mask is agreed on with the other clients through Diffie-Hellman key exchange and is additively cancelable when summed with others. Considering that some clients may drop out at any point, their masks cannot be canceled. To handle this problem, each client uses a threshold secret sharing scheme to split her private seed of a Pseudo-Random Number Generator (PRNG) for generating the self mask as well as her private key for generating the mutual mask, and then distribute the shares to the other clients. As long as some minimum (no less than the threshold) number of clients remain alive, they can jointly help the cloud server remove the self masks of live clients and the mutual masks between dropped and live clients.

Third, we present the scalability and efficiency of the secure aggregation protocol. We list its communication, computation,

TABLE II
COMPLEXITIES OF THE SECURE AGGREGATION PROTOCOL IN THE
HONEST-BUT-CURIOUS SETTING.

	Communication	Computation	Storage
Client Server	$ \begin{array}{ c c } O(n+l) \\ O(n^2+nl) \end{array} $	$O(n^2 + nl)$ $O(n^2l)$	$O(n+l)$ $O(n^2+l)$

and storage complexities in Table II, where n is the number of clients involved in the aggregation, and l denotes the number of data items held by each client or the dimension of her data vector. We can see that this protocol is quite efficient for large-scale data vectors, especially from communication overhead, and thus can apply to mobile applications. In particular, as reported in [18], when 2^{14} clients are involved in the aggregation, and each client has 2^{24} 16-bit values, the communication overhead of the secure aggregation protocol expands $1.98\times$ over sending data in the clear.

3) Bloom Filter: Bloom filter, conceived by Bloom in 1970 [68], is a space-efficient probabilistic data structure to represent a set whose elements come from a huge domain. In addition, when testing whether an element is a member of the set, a false positive is possible, but a false negative is impossible. In other words, an element that is diagnosed to be present in the set possibly does not belong to the set in reality, and an element that is judged to be not present definitively does not belong to the set. We describe its technical details and properties as follows.

A Bloom filter is a β -length bit vector initially set to 0, denoted as b. In addition, it requires h different independent hash functions. The output range of these hash functions is $\{1, 2, \dots, \beta\}$, which corresponds to the β positions of the Bloom filter. To represent a set of ϕ elements, we apply h hash functions to each element and set the Bloom filter at the positions of h hash values to 1. In the membership test phase, to check whether an element belongs to the set, we simply check the Bloom filter at the positions of its h hash values. If any of the bits at these positions is 0, the element is definitely not in the set. If all are 1, then either the element is in the set, or the bits have by chance been set to 1 during the insertion of other elements, resulting in a false positive. Specifically, the false positive rate (FPR) of a Bloom filter depends on the length of Bloom filter β ; the number of hash functions h; and the cardinality of set ϕ . According to [69], [70], its detailed formula is given as

$$\mathrm{FPR} = \left(1 - \left(1 - \frac{1}{\beta}\right)^{h\phi}\right)^h \approx \left(1 - \exp\left(-\frac{h\phi}{\beta}\right)\right)^h.$$

Given β and ϕ , to minimize the false positive rate, the optimal number of hash functions is

$$h = \ln 2 \frac{\beta}{\phi}.$$

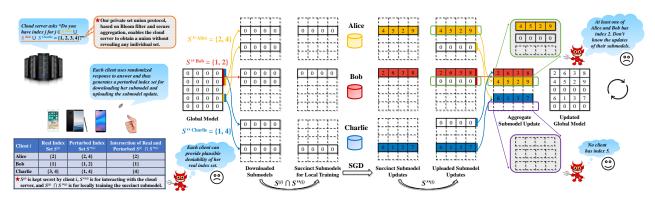


Fig. 3. An illustration of the design rationale of secure federated submodel learning.

In addition, given ϕ and assuming the optimal number of hash functions h is used, to achieve a desired false positive rate FPR, the optimal length of Bloom filter should be

$$\beta = -\frac{\phi \ln \text{FPR}}{\left(\ln 2\right)^2}.\tag{1}$$

Thus, the optimal number of bits per the set's element is

$$\frac{\beta}{\phi} = -\frac{\ln \text{FPR}}{\left(\ln 2\right)^2},$$

and the corresponding number of hash functions is

$$h = -\frac{\ln \text{FPR}}{\ln 2}.$$

The above deductions mean that for a given false positive rate, the length of a Bloom filter is proportional to the size of the set being filtered, while the required number of hash functions only relies on the target false positive rate.

We further introduce an appealing property of Bloom filter when performing union over the underlying sets. To represent sets, denoted as $\forall i \in \mathcal{C}, \mathcal{S}^{(i)}$, we use Bloom filters with the same length and the same hash functions, denoted as $\forall i \in \mathcal{C}, \mathbf{b}^{(i)}$. Then, the union of these sets, i.e., $\bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)}$, can be represented by a Bloom filter, which performs bitwise OR operations over the Bloom filters, i.e., $\bigvee_{i \in \mathcal{C}} \mathbf{b}^{(i)}$. Such a union operation is lossless (implying that the false positive rate remains unchanged) in the sense that the resulting Bloom filter is the same as the Bloom filter created from scratch using the union of these sets. In addition, because the Bloom filter needs to accommodate the union of sets, the parameter ϕ , denoting the cardinality of the set, should be determined by estimating the cardinality of the union $\bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)}$ rather than that of each individual set $\mathcal{S}^{(i)}$.

We finally present a generalized version of Bloom filter, called counting Bloom filter [71]. It consists of β counters/integers rather than β bits and can represent a multiset, where an element can occur more than once. The main difference between counting Bloom filter and Bloom filter lies in that when representing an element, we increment the counters by one at the positions of its h hash values. Thus, compared with the membership test in Bloom filter, counting Bloom filter further supports more general counting queries

of a given element, i.e., whether the count number of a given element in a multiset is smaller than a certain threshold.

IV. DESIGN OF SECURE FEDERATED SUBMODEL LEARNING

In this section, we present the design rationale and the design details of Secure Federated Submodel Learning (SFSL).

A. Design Rationale

We illustrate our key design principles mainly through demonstrating how to handle two fundamental problems raised in Section I-D and how to resolve several practical issues.

As shown in Fig. 3, we handle two fundamental problems in a unified manner rather than in separate ways. During both download and upload phases, a client consistently uses a perturbed index set in place of her real index set. In contrast, during the local training phase, the client leverages the intersection of her real index set and perturbed index set to prepare the succinct submodel and involved user data. With the blinding of the perturbed index set to interact with the outside world, the client holds a plausible deniability of some index being in or not in her real index set. Specifically, the client generates her perturbed index set locally with randomized response as follows. First, the sensitive question asked by the cloud server here is "Do you have a certain index?". Then, the client answers "Yes" with two customized probabilities, conditioned on whether the index is or is not in her real index set. These two probabilities allow the client to fine-tune balance between privacy and utility.

We further carefully examine the feasibility of our index set perturbation method in handling two fundamental problems. For the first problem in the download phase, if a client intends to download a certain row of the full matrix, and she actually downloads, she can blame her action to randomization, i.e., the occurrence of the event that the index not in a client's real index set returns a "Yes" answer. Similarly, if a client does not intend to download the row, and she actually does not, she can still attribute her action to randomization, i.e., the occurrence of the event that the index in a client's real index returns a "No" answer. Regarding the second problem in the upload phase, the usage of the perturbed index set

still empowers a client to deny her underlying intention of modifying or not modifying some row of the full matrix, even if the cloud server observes her binary action of modifying or not modifying. Additionally, for a concrete row, there are two different groups of clients involved in the joint modification: (1) One group consists of those clients who intend to modify the row and contribute nonzero modifications; and (2) the other group comprises those clients who do not intend to modify the row and pretend to modify by submitting zero modifications. Under the secure aggregation guarantee, even though the cloud server observes the aggregate modification, it is hard for the cloud server, as an adversary, to identify any individual modification and further to infer whether some client originally intends to perform a modification or not. The hardness is controlled by the sizes of two groups, alternatively the probabilities of an index in and not in the real index set returning a "Yes" answer, which are fully tunable by clients.

In addition to these two basic problems, there still exist two practical issues to be solved before the above index set perturbation method can apply to federated submodel learning. The first issue regards practical efficiency, i.e., whether it is practical and necessary for the cloud server to ask "Do you have a certain index?" for each index in the full index set. In our context, the number of rows of the matrix, representing the deep learning model, is in the magnitude of billions. Thus, it is impractical for a client to answer billion-scale questions and further download and securely upload those rows with "Yes" answers of the full model. We thus turn to narrowing down the size of the questions and identify a sufficient and necessary index set, namely the union of the chosen clients' real index sets. Our optimization is inspired by an example: if a client's real index set is of size 100, and the full index set is of size 1 billion, using the probability parameters in the survey of party membership, her expected number of "Yes" answers is $10^2 \times 75\% + (10^9 - 10^2) \times 25\% \approx 2.5 \times 10^8$. Such a calculation implies that the dominant "Yes" answers are those with "No" in reality but "Yes" due to randomness. Nevertheless, most of the "No"-to-"Yes" answers are useless. More specifically, for those indices that do not belong to any client's real index set, e.g., index 5 in Fig. 2 and Fig. 3, although part (25% in expectation) of clients upload zero vectors for randomization, the cloud server can still infer from the aggregate zero vectors that these clients do not actually have the indices. Thus, it is not necessary to cover any index outside the union.

Accompanied with the first issue, another fundamental and thorny problem that arises is how multiple clients can obtain the union of their real index sets under the mediation of an untrusted cloud server without revealing any individual client's real index set to others, i.e., the need of a private set union protocol. Considering any existing scheme doest not satisfy the atypical setting and the stringent requirement of federated submodel learning, we design a novel private set union scheme based on Bloom filter, secure aggregation, and randomization. We first let each chosen client represent her real index set as a Bloom filter. Then, different from the common practice to derive the union of sets by performing bit-wise OR operations

over their Bloom filters, which naturally requires both addition and multiplication operations, we let the cloud server directly "sum" the Bloom filters. Here, the sum operation can be performed obliviously and efficiently under the coordination of an untrusted cloud server with the secure aggregation protocol. The aggregate Bloom filter is actually a counting Bloom filter, equivalent to constructing it from scratch by inserting each set in sequence. Besides the membership information, the counting Bloom filter still contains the count number of each element in the union. To prevent such an undesirable leakage, we let each client replace bit 1s in her Bloom filter with random integers, while keeping each bit 0 unchanged. When recovering the union of real index sets, one naive method for the cloud server is to do membership tests for the full index set. which is prohibitively time consuming, and can also introduce a large number of false positives. To handle these problems, we let the cloud server first divide the full index set into a certain number of partitions, and then let each client fill in a bit vector to indicate whether there exists an element of her real index set falling into these partitions. Then, just like computing the union with Bloom filters, the cloud server can securely determine those partitions that contain clients' real indices to further facilitate efficient union reconstruction.

The second issue regards the longitudinal privacy guarantee when a client is chosen to participate in multiple communication rounds. However, the initial version of randomized response only provides a rigorous privacy guarantee when an audience answers the same question once, facilitating only a one-time response to "Do you have a certain index?" in our context. Thus, we need to extend the original randomized response mechanism to allow repeated responses from the same client to those already answered indices in a privacypreserving manner. Our extension leverages key principles from Randomized Aggregatable Privacy-Preserving Ordinal Response (RAPPOR) [47], [48] and plays the randomized response game twice with a memoization step between. Specifically, the noisy answers generated by the inner randomized response will be memoized and permanently replace real answers in the outer randomized response. This ensures that even though a client responds to the membership of a concrete index for an infinite number of times, she can still hold a plausible deniability of her real answer, where the level of deniability is lower bounded by the memoized noisy answer.

B. Design Details

Following the guidelines in Section IV-A, we propose a secure scheme for federated submodel learning. We introduce the scheme in a top-down manner, where we first give an overview of its top-level architecture and then show two underlying modules, namely index set perturbation and private set union. For the sake of clarity, we outline our design in Algorithm 1, Algorithm 2, and Algorithm 3.

1) Secure Federated Submodel Learning: Before presenting our secure federated submodel learning framework, we first briefly review the federated averaging algorithm [9], which is the cornerstone and core of conventional federated learning.

Algorithm 1: Secure Federated Submodel Learning

```
/* Cloud server's process */
 1 Initializes the global model W;
2 foreach communication round do
         Randomly selects n clients, denoted as C, where |C| = n;
3
         Launches private set union (Algorithm 3), gets the union of
           n clients' real index sets, namely \bigcup_{i\in\mathcal{C}}\mathcal{S}^{(i)}, and delivers the union result to the up-to-date set of clients who are
           alive, denoted as \hat{\mathcal{C}} \subset \mathcal{C};
         foreach Client i \in \hat{\mathcal{C}} do
5
               Receives and stores the perturbed index set S''^{(i)} from
                 client i, and returns the submodel \mathbf{W}_{\mathcal{S}''(i)} and
                 training hyperparameters to i;
               Securely aggregates weighted submodel
               updates and count vectors
         foreach j \in \bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)} do
7
               Determines the live clients involving index j, denoted
                 as \hat{C}_j = \{i | i \in \hat{C} \land j \in \mathcal{S}^{\prime\prime(i)}\}, lets them submit
                 materials for secure aggregation, and obtains the sum
                of (weighted) updates \sum_{i \in \hat{C}_i} \Delta \mathbf{w}_j^{(i)} and the total
                count number of relevant samples \sum_{i \in \hat{C}_i} v_j^{(i)};
               Updates the j-th row of the global model \mathbf{W} by
 9
                adding \sum_{i \in \hat{\mathcal{C}}_i} \Delta \mathbf{w}_j^{(i)} / \sum_{i \in \hat{\mathcal{C}}_i} v_j^{(i)};
    /* Client i's process */
10 Determines her real index set S^{(i)} based on local data;
11 Participates in private set union (Algorithm 3);
12 Generates a perturbed index set S''^{(i)} (Algorithm 2);
13 Uses S''(i) to download a submodel, denoted as \mathbf{W}_{S''(i)};
14 Depending on the succinct index set S''^{(i)} \cap S^{(i)}, locally
     extracts the succinct submodel \mathbf{W}_{\mathcal{S}''(i)\bigcap\mathcal{S}^{(i)}} from \mathbf{W}_{\mathcal{S}''(i)}
     and prepares involved data as the succinct training set;
15 Locally trains \mathbf{W}_{\mathcal{S}''(i) \bigcap \mathcal{S}^{(i)}} using the hyperparameters and
     obtains the update of succinct submodel \Delta \mathbf{W}_{\mathcal{S}''(i) \bigcap \mathcal{S}^{(i)}};
16 Initializes the submodel update to be uploaded, denoted as
      \Delta \mathbf{W}_{\mathcal{S}''(i)}, all to 0, and then adds \Delta \mathbf{W}_{\mathcal{S}''(i) \cap \mathcal{S}^{(i)}};
17 Counts the number of samples involving each index j \in \mathcal{S}''^{(i)}
     and stores the results to the vector \mathbf{v}_{S''(i)};
18 Updates \Delta \mathbf{W}_{\mathcal{S}''(i)} by multiplying each row with the
     corresponding count number in \mathbf{v}_{\mathcal{S}''(i)};
19 Uploads materials for securely aggregating \Delta \mathbf{W}_{S''(i)}, \mathbf{v}_{S''(i)}.
```

In particular, federated averaging is a synchronous distributed learning method, for non-iid and unbalanced training data distributed at massive communication-constrained clients, under the coordination of a cloud server. At the beginning of one communication round, the cloud server sends the up-to-date parameters of the global model and the training hyperparameters to some clients. Here, the training hyperparameters include the optimization algorithm, typically mini-batch SGD, the local batch size (the number of training samples used to locally update the global model once, namely per iteration), the local epochs (the number of passes over a client's entire training data), and the learning rate. Then, each chosen client trains the global model on her data and uploads the update of the global model together with the size of her training data to the cloud server. The cloud server takes a weighted average of all updates, where one client's weight is proportional to the size of her local data, and finally adds the aggregate update to the global model.

We now present secure federated submodel learning in Algorithm 1, which generalizes the federated averaging algorithm to support effective and efficient submodel learning and preserves desirable security and privacy properties while incorporating the unstable and limited network connections of mobile devices. At the initial stage, the cloud server randomly initializes the global model (Line 1). For each communication round, the cloud server first selects some clients to participate (Line 3) and also maintains an up-to-date set of clients who are alive throughout the whole round. A chosen client determines her real index set based on her local data, which can specify the "position" of her truly required submodel (Line 10). For example, if the visited goods IDs of a Taobao user include $\{1,2,4\}$, then she requires the first, second, and fourth rows of the embedding matrix for goods IDs, which further implies that her real index set should contain $\{1, 2, 4\}$. Then, the cloud server launches the private set union protocol to obtain the union of all chosen clients' real index sets while keeping each individual client's real index set in secret (Lines 4 and 11). The union result will be further delivered to live clients, based on which each client can perturb her real index set with a customized local differential privacy guarantee (Line 12). In addition, each client will use the perturbed index set, rather than the real index set, to download her submodel and upload the submodel update (Lines 13 and 19). In other words, when interacting with the cloud server, a client's real index set is replaced with her perturbed index set, which provides deniability of her real index set and thus obscures her training data. Upon receiving the perturbed index set from a client, the cloud server stores it for later usage and returns the corresponding submodel and the training hyperparameters to the client (Line 6). Depending on the intersection of the real index set and the perturbed index set, called the succinct index set, the client extracts a succinct submodel and prepares involved data as the succinct training set (Line 14). For example, if a Taobao user's real index set is $\{1, 2, 4\}$ and her perturbed index set is $\{2, 4, 6, 9\}$, she receives a submodel with row indices $\{2,4,6,9\}$ from the cloud server, but just needs to train the succinct submodel with row indices $\{2,4\}$ over her local data involving goods IDs $\{2,4\}$. After training under the preset hyperparameters, the client derives the update of the succinct submodel (Line 15) and further prepares the submodel update to be uploaded with the perturbed index set by adding the update of the succinct submodel to the rows with the succinct indices and padding zero vectors to the other rows (Line 16). Additionally, to facilitate the cloud server in averaging submodel updates according to the sizes of relevant local training data, each client also needs to count the number of her samples involving every index in the perturbed index set (Line 17). In particular, the numbers of samples involving the indices outside the succinct index set are all zeros. Furthermore, each client prepares the submodel update to be uploaded by multiplying each row with the corresponding count number, namely the weight, in advance (Line 18).

Algorithm 2: Client i's Index Set Perturbation

```
Input: Client i's real index set S^{(i)}, Union of the chosen
               clients' real index sets \bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)}, Client i's memoized
               index set \mathcal{Y}^{(i)} (resp., \mathcal{N}^{(i)}) initialized to \emptyset at very
               beginning with "Yes" (resp., "No") permanent answers
               to the question "Do you have a certain index?", Client
               i's customized probability parameters
    0 \leq p_1^{(i)}, p_2^{(i)}, p_3^{(i)}, p_4^{(i)} \leq 1. Output: Client i's doubly perturbed index set \mathcal{S}''^{(i)}
 1 \mathcal{S}'^{(i)} = \emptyset, \mathcal{S}''^{(i)} = \emptyset;
     // Permanent randomized response
 2 foreach j \in \bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)} \bigwedge j \notin \mathcal{Y}^{(i)} \bigcup \mathcal{N}^{(i)} do
          if j \in \mathcal{S}^{(i)} then
               Adds j to \mathcal{S}^{\prime(i)} with probability p_1^{(i)};
 5
                Adds j to \mathcal{S}'^{(i)} with probability p_2^{(i)};
                Memoization of permanent responses
          if j \in \mathcal{S}^{\prime(i)} then
 7
             | \mathcal{Y}^{(i)} = \mathcal{Y}^{(i)} \bigcup j; 
 8
                \mathcal{N}^{(i)} = \mathcal{N}^{(i)} \cup j;
10
          Instantaneous randomized response
11 foreach j \in \bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)} do
          if j \in \mathcal{Y}^{(i)} then
12
               Adds j to S''^{(i)} with probability p_3^{(i)};
13
14
                Adds j to S''^{(i)} with probability p_4^{(i)};
15
16 return S''^{(i)}
```

Finally, the weighted submodel updates and the count vectors from live clients are securely aggregated under the guidance of the cloud server (Lines 7–9 and 19). Specifically, the cloud server guides the secure aggregation by enumerating every index in the union of the chosen clients' real index sets. For each index, the cloud server first determines the set of live clients whose perturbed index sets contain this index and then lets these clients submit the materials for securely adding up the weighted updates and the count numbers with respect to this index (Line 8). The cloud server finally applies the update to the global model in this row by adding the quotient of the sum of the weighted updates and the total count number, namely the weighted average (Line 9). Considering that the weighted submodel updates and the count numbers are aggregated side by side, each client can augment the matrix, denoting her weighted submodel update, with the transposed count vector in the last column, when preparing materials for secure aggregation (Line 19). In addition, to reduce the interactions between the cloud server and a client, they can package all the materials supporting secure aggregation, rather than exchange the materials for one index each time (Lines 7– 9), i.e., the cloud server executes Lines 7 and 8 for each live client $i \in \mathcal{C}$ in parallel and then executes Line 9 for each index in the union $j \in \bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)}$.

2) Index Set Perturbation: We now present how a client can generate a perturbed index set to download her submodel and to upload the update of the submodel, with a customized local differential privacy guarantee against the cloud server. Just like the exemplary question about party membership, the sensitive question here is "Do you have a certain index?", asked by the cloud server. The clients participating in one round of federated submodel learning make up the population to be surveyed. Thus, the clients can use randomized response to answer "Yes" or "No", which provides measurable deniabilities of their true answers. However, as sketched in Section IV-A, several practical issues need to be resolved so that randomized response can truly apply here. In what follows, we elaborate on our solution details on these issues.

As shown in Algorithm 2, we view the union of the chosen clients' real index sets, rather than the full index set, as the scope of the cloud server's questionnaire (Input). We reason about necessity and sufficiency as follows. Without loss of generality, we consider client i and the other chosen clients $\mathcal{C}\setminus\{i\}$. If any client in $\mathcal{C}\setminus\{i\}$ wants to obtain deniability of her real indices, she requires client i to join as an audience to answer the questions about her real indices, which implies that client i should know the union of the other chosen clients' real index sets. By incorporating client i's own real index set, the questionnaire to client i should contain the union of all chosen clients' real index sets. We further apply the above reasoning to all clients $\forall i \in \mathcal{C}$ and can derive that the global questionnaire should contain the union of all chosen clients' real index sets. Next, we illustrate whether the union suffices. We consider any index outside the union. Under the conventional federated learning framework, each client will upload a zero vector to the cloud server for this index. When the cloud server learns that the sum is a zero vector, she can infer that all chosen clients do not have this index. Please see index 5 in Fig. 2(a) for an intuition. From this perspective, any index outside the union does not need to be preserved in federated submodel learning as well. Nevertheless, suppose an index outside the union is introduced by chance, e.g., due to a false positive of Bloom filter when reconstructing union in Algorithm 3. A nice phenomenon occurs. On one hand, the privacy of federated submodel learning can be enhanced in the sense that the cloud server can only ascertain that those clients who return "Yes" answers do not really have this index, but cannot ascertain the states of the other clients due to plausible deniability. On the other hand, those clients with "Yes" answers need to download the row with respect to this index and further to upload zero vectors through secure aggregation, which are useless and increase their overheads.

Given the questionnaire, client i basically uses two probability parameters $p_1^{(i)}, p_2^{(i)}$ in randomized response to finetune the tension among effectiveness, efficiency, and privacy (Lines 3–6). In particular, $p_1^{(i)}$ denotes the probability that an index in client i's real index set will return a "Yes" answer and controls the factual size of a client's user data contributed to federated submodel learning. Thus, a larger $p_1^{(i)}$ implies

better effectiveness in terms of convergency rate. In addition, $p_2^{(i)}$ denotes the probability that an index outside client i's real index set will return a "Yes" answer and determines the number of redundant rows to be downloaded and the number of padded zero vectors to be uploaded through the secure aggregation protocol. Hence, given a fixed $p_1^{(i)}$, a smaller $p_2^{(i)}$ indicates higher efficiency. Furthermore, $p_1^{(i)}$, $p_2^{(i)}$ jointly adjust the level of local differential privacy, where a pair of closer $p_1^{(i)}$, $p_2^{(i)}$ provides a better privacy guarantee. We examine three typical examples: (1) The randomized response in the party membership survey takes $p_1^{(i)} = 75\%$ and $p_2^{(i)} = 25\%$ for each respondent; (2) conventional federated learning essentially leverages the full index, takes $p_1^{(i)} = 1$ and $p_2^{(i)} = 1$ for each client, and offers the best privacy and effectiveness guarantees but the worst efficiency guarantee; and (3) federated submodel learning adopts $p_1^{(i)} = 1$ and $p_2^{(i)} = 0$ for each client and provides the best effectiveness and efficiency guarantees but the worst privacy guarantee.

Considering that client i can be chosen to participate in multiple communication rounds and needs to repeatedly respond to some answered indices, we extend the basic randomized response mechanism to offer a rigorous privacy guarantee, also called longitudinal privacy in the literature [47], [48], [59]. We adopt a memoization technique from RAPPOR. The core idea of RAPPOR is to play the randomized response game twice with a memoization step between. The first perturbation step, called permanent randomized response, is used to create a noisy answer, which is memoized by the client and permanently reused in place of the real answer. The second perturbation step, called instantaneous randomized response, reports on the memoized answer over time, eventually completely revealing it. In other words, the privacy level, guaranteed by the underlying memoized answer in the permanent randomized response, imposes a lower bound on the privacy level, ensured by each instantaneous/revealed response. When the memoization technique is applied to federated submodel learning, we let client i maintain two index sets with "Yes" and "No" answers in the permanent randomized response, respectively (Input). Here, the permanent randomized response mechanism is parameterized by two probabilities $p_1^{(i)}, p_2^{(i)}$ to tune privacy and utility (Lines 3-6), as illustrated in the preceding paragraph. In addition, given that one client can be grouped with distinct clients in different communication rounds while the union of real index sets varies from one round to another, the client needs to handle new indices. As a new index comes (Line 2), client i generates a permanent noisy answer for it and further updates her two memoized sets (Lines 7–10). Moreover, client i obtains her final perturbed index set by performing an instantaneous randomized response over the memoized answers to the union of real index sets in the current communication round (Lines 11-16). In particular, the instantaneous randomized response is parameterized with another two probabilities $p_3^{(i)}, p_4^{(i)}$ (Lines 13 and 15), similar to $p_1^{(i)}, p_2^{(i)}$ in the permanent randomized response. Now, these four probability parameters jointly support tuning the tension among privacy, effectiveness, and efficiency. More specifically, $p_5^{(i)}=p_1^{(i)}(p_3^{(i)}-p_4^{(i)})+p_4^{(i)}$, denoting the probability that an index in client i's real index set finally returns a "Yes" answer, and $p_6^{(i)}=p_2^{(i)}(p_3^{(i)}-p_4^{(i)})+p_4^{(i)}$, denoting the probability that an index not in client i's real index set finally returns a "Yes" answer, now play the same roles as $p_1^{(i)}$ and $p_2^{(i)}$, respectively. Detailed derivations of $p_5^{(i)},p_6^{(i)}$ are deferred to Section V-A.

Finally, we provide some comments on the above design. First, our design is different from conventional locally differentially private schemes (e.g., randomized response and RAPPOR), which require each participating user to choose the same probability parameters (i.e., $\forall i \in \mathcal{C}, p_1^{(i)} = p_1, p_2^{(i)} = p_2, p_3^{(i)} = p_3, p_4^{(i)} = p_4$), so that true statistics (e.g., heavy hitter, histogram, and frequency) can be well estimated using collected noisy answers, particularly after additional corrections. Such a requirement/assumption is no longer needed in our design because the cloud server, as the aggregator, performs aggregate statistics based on secure aggregation rather than over the noisy answers, e.g., counting how many samples from the chosen clients involve a certain index in total (Algorithm 1, Line 8). Therefore, as mentioned above, different clients can customize probability parameters to tune privacy and utility. Second, our index perturbation mechanism in Algorithm 2 needs a prerequisite that the real index set of a client does not change when she participates in different communication rounds. Considering that the real index set corresponds to the client's local data, this prerequisite can be further converted to the invariance of the client's local data. One feasible way to meet this prerequisite is to introduce the concept of "period" into federated submodel learning, e.g., one period can be set to one month. In a concrete period, a client uses the historical data in the previous one period to participate in federated submodel learning for several communication rounds. In addition, when entering a new period, the client just restarts Algorithm 2. The other feasible way is to allow changes in a client's real index set from one communication round to another. This implies that the underlying binary states of some indices may change. For example, if a client's local data and thus her real index set expand incrementally, some indices, which were not in, can fall into the real index set in later rounds. Recently, Erlingsson et al. [59] considered a similar setting, in particular the collection of user statistics (e.g., software adoption) for multiple times with each user changing her underlying boolean value for a limited number of times. Therefore, their design, based on binary tree and Bernoulli distribution, can be leveraged to extend Algorithm 2, allowing a client to change her local data and thus her real index set in different communication rounds.

3) Private Set Union: We introduce the last module of our design: private set union. We first briefly review related work about private set operations, with a focus on the often overlooked but significantly important private set union. Then, we outline the practical infeasibility of existing protocols when adapted to federated submodel learning. We finally present our efficient and scalable scheme.

Algorithm 3: Private Set Union

Input: Client *i*'s real index set $\mathcal{S}^{(i)}$ for all $i \in \mathcal{C}$ **Output:** Union of real index sets $\bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)}$

- 1 Cloud server determines the partitions of the full index set S;
- 2 foreach Client $i \in C$ do
- 3 Represents $S^{(i)}$ as a Bloom filter $\mathbf{b}^{(i)}$;
- Perturbs $\mathbf{b}^{(i)}$ to an integer vector $\mathbf{b}'^{(i)}$ by replacing each bit 1 in $\mathbf{b}^{(i)}$ with a random integer from \mathbb{Z}_R ;
- 5 Uses a bit vector $\mathbf{a}^{(i)}$ to indicate whether there exists an element in $\mathcal{S}^{(i)}$ falling into the partitions of \mathcal{S} ;
- Perturbs $\mathbf{a}^{(i)}$ to an integer vector $\mathbf{a}'^{(i)}$ by replacing each bit 1 in $\mathbf{a}^{(i)}$ with a random integer from \mathbb{Z}_R ;
- 7 Submits materials for securely aggregating $\mathbf{b}'^{(i)}, \mathbf{a}'^{(i)};$
- 8 Cloud server obtains $\sum_{i \in \mathcal{C}} \mathbf{b}'^{(i)}$ and $\sum_{i \in \mathcal{C}} \mathbf{a}'^{(i)}$ with the secure aggregation protocol, reconstructs the union $\bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)}$, and delivers $\bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)}$ to each live client $i \in \hat{\mathcal{C}}$.

The goal of a private set operation protocol is to allow multiple parties, where each party holds a private set, to obtain the result of an operation over all the sets, without revealing each individual private set and without introducing a trusted third party. Compared with Private Set Intersection (PSI) [72]-[77] and Private Set union Cardinality (PSC) [78]-[80], which have received tremendous attention and also have seen several practical applications, such as in social networks [81], [82]; human genome testing [83]; locationbased services [84]; security incident information sharing [85]; online advertising [85]; private contact discovery [75]; and the Tor anonymity network [86], there is little work and negligible focus on Private Set Union (PSU). Nevertheless, private set union promises a wide range of applications in practice, e.g., union queries over several databases, and, more generally, integration/sharing of datasets from multiple private sources. Thus, independent of federated submodel learning, the task of designing a practical private set union protocol itself is highly desired and urgent. Existing protocols mainly come from the fields of data mining and cryptography. In the data mining field, the representative design of private set union [87] is based on commutative encryption and requires direct communication between any pair of two parties. Unfortunately, the design leaks the cardinality of any two-party set intersection, and the underlying commutative encryption is fragile as well. In the cryptography field, according to the representation format of a set, the protocols can be generally divided into two categories: polynomial based [88]-[92] and Bloom filter based [93]-[95]. For the polynomial-based protocols, elements of a set are represented as the roots of a polynomial, and the union of two sets is converted to the multiplication of two polynomials. For the protocols based on Bloom filter, the union operation over sets is normally transformed to the element-wise OR operation over Bloom filters, as demonstrated in Section III-B3, whereas the logic OR operation can be further converted to bit addition and bit multiplication. To obliviously perform addition and multiplication operations, the above two kinds of protocols mainly turn to generic secure two-party/multiparty computation (e.g., garbled circuit, homomorphic encryption, secret sharing, and oblivious transfer), or outsource secure computation to multiple noncolluding servers. Due to unaffordable computation and communication overheads, none of the existing private set union protocols have been deployed in practice. In addition to inefficiency, the basic setting of these protocols significantly differs from that of federated submodel learning, where clients cannot directly communicate with each other and should mediate through an untrusted cloud server. Additionally, the set elements here can come from a billion-scale domain, which has not been touched in previous work as of yet.

Given the infeasibility of existing protocols and the atypical setting of federated submodel learning, we present our new private set union scheme in Algorithm 3. First, each client represents her real index set as a Bloom filter (Line 3). The details about how to set the parameters of the Bloom filter can be found in Section III-B3. Second, different from the common practice to derive the union of sets by performing bitwise OR operations over their Bloom filters, which requires both addition and multiplication operations, we let the cloud server directly sum the Bloom filters. Here, the sum operation can be conducted obliviously and efficiently under the coordination of the untrusted cloud server with secure aggregation (Line 8). In addition, the resulting Bloom filter is actually a counting Bloom filter, equivalent to constructing it from scratch by sequentially inserting each real index set. In addition to membership information, the counting Bloom filter also contains the count numbers of elements in the union of real index sets. In other words, the cloud server may learn how many clients have a certain index, which is undesired in our context. To prevent the leakage of count numbers, we let each client generate a perturbed integer vector, which replaces each bit 1 in her Bloom filter with a random integer and keeps each bit 0 unchanged (Line 4). Such a perturbation process can obscure count numbers while retaining membership information. Third, after obtaining the sum of perturbed Bloom filters, the cloud server can recover the union of real index sets by doing membership tests for the full index set. For example, to judge whether an index belongs to the union, we check the resulting integer vector at the positions of its hash values. The index is considered to be in the union only if all the values are nonzero. However, one practical issue arises: the domain of index can be huge, e.g., the full size of the goods IDs in Taobao is in the magnitude of billions. Thus, it can be prohibitively time consuming to enumerate all indices. Even worse, the direct enumeration method can also introduce a large number of false positives in the union, i.e., those indices not in the union are falsely judged to be in, which can further lead to unnecessary redundancy in the download and upload phases. To handle these problems, we further incorporate a private "partition" union to narrow down the scope of index for union reconstruction above. We let the cloud server divide the full domain of the index into a certain number of partitions ahead of time (Line 1). A good partition scheme needs to well balance the pros in the union

reconstruction phase and the cons of additional cost. Given the partitions, each client first uses a bit vector to record whether there exists an index in her real index set falling into the partitions (Line 5). Just the same as the Bloom filter to hide the concrete count numbers, the client further replaces each bit 1 with a random integer (Line 6). Then, the cloud server obtains the sum of the integer vectors using the secure aggregation protocol and reveals those partitions with nonzero integers in the corresponding positions. By simply doing membership tests for the indices falling into these partitions, the cloud server can efficiently construct the union. Last, the union is delivered to all live clients (Line 8).

V. SECURITY AND PERFORMANCE ANALYSES

In this section, we first analyze the privacy guarantees of our secure federated submodel learning scheme according to Definition 1 and Definition 2, i.e., Theorem 1 and Theorem 2. We also provide an instantiation of our scheme, where each client consistently uses the union of the chosen clients' real index sets when interacting with the cloud server, and prove that its security and privacy guarantees are as strong as conventional federated learning with secure aggregation (hereinafter also called "Secure Federated Learning" and abbreviated as "SFL"), i.e., Theorem 3. We then show the proven security of our proposed private set union protocol, i.e., Theorem 4. We finally analyze the performance of our scheme by comparing with that of secure federated learning.

A. Security and Privacy Analyses

By Definition 1, we analyze the local differential privacy guarantee of index set perturbation in Algorithm 2. As stepping stones, we first analyze permanent randomized response and a one-time instantaneous randomized response in Lemma 1 and Lemma 2, which impose an upper bound and a lower bound on the privacy level of Algorithm 2, namely Theorem 1.

Lemma 1. Permanent randomized response in Algorithm 2 for client i achieves local differential privacy at the level $\epsilon_{\infty}^{(i)} = \ln\left(\max\left(\frac{p_1^{(i)}}{p_2^{(i)}},\frac{p_2^{(i)}}{p_1^{(i)}},\frac{1-p_1^{(i)}}{1-p_2^{(i)}},\frac{1-p_2^{(i)}}{1-p_1^{(i)}}\right)\right)$.

Proof. We focus on a certain index $j \in \bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)}$. According to Definition 1, we need to consider all possible pairs of inputs from client i and all possible outputs of the permanent randomized response in Algorithm 2. Here, the input pair is j in and not in client i's real index set, namely $j \in \mathcal{S}^{(i)}$ and $j \notin \mathcal{S}^{(i)}$. In addition, the possible outputs are j obtaining "Yes" and "No" noisy answers for memoization, namely $j \in \mathcal{Y}^{(i)}$ and $j \in \mathcal{N}^{(i)}$. We thus can compute four ratios between the conditional probabilities of a permanent output with a pair of distinct inputs: $\frac{\Pr(j \in \mathcal{Y}^{(i)} | j \in \mathcal{S}^{(i)})}{\Pr(j \in \mathcal{Y}^{(i)} | j \notin \mathcal{S}^{(i)})} = \frac{p_1^{(i)}}{p_2^{(i)}}, \quad \frac{\Pr(j \in \mathcal{Y}^{(i)} | j \notin \mathcal{S}^{(i)})}{\Pr(j \in \mathcal{Y}^{(i)} | j \notin \mathcal{S}^{(i)})} = \frac{p_2^{(i)}}{p_1^{(i)}}, \quad \frac{\Pr(j \in \mathcal{N}^{(i)} | j \notin \mathcal{S}^{(i)})}{\Pr(j \in \mathcal{N}^{(i)} | j \notin \mathcal{S}^{(i)})} = \frac{1-p_1^{(i)}}{1-p_1^{(i)}}.$ By Definition 1, we can derive the level of local differential privacy $\epsilon_{\infty}^{(i)}$: $\exp\left(\epsilon_{\infty}^{(i)}\right) = \max\left(\frac{p_1^{(i)}}{p_2^{(i)}}, \frac{p_2^{(i)}}{p_1^{(i)}}, \frac{1-p_1^{(i)}}{1-p_2^{(i)}}, \frac{1-p_2^{(i)}}{1-p_1^{(i)}}\right) \Rightarrow \epsilon_{\infty}^{(i)} = \ln\left(\max\left(\frac{p_1^{(i)}}{p_2^{(i)}}, \frac{p_2^{(i)}}{p_1^{(i)}}, \frac{1-p_2^{(i)}}{1-p_1^{(i)}}, \frac{1-p_2^{(i)}}{1-p_1^{(i)}}\right)\right).$

Lemma 2. A one-time instantaneous randomized response in Algorithm 2 for client i satisfies local differential privacy at the level $\epsilon_1^{(i)}$, where $\epsilon_1^{(i)} = \ln\left(\max\left(\frac{p_5^{(i)}}{p_6^{(i)}}, \frac{p_6^{(i)}}{p_5^{(i)}}, \frac{1-p_5^{(i)}}{1-p_6^{(i)}}, \frac{1-p_6^{(i)}}{1-p_5^{(i)}}\right)\right)$, $p_5^{(i)} = \Pr\left(j \in \mathcal{S}''^{(i)}|j \in \mathcal{S}^{(i)}\right) = p_1^{(i)}\left(p_3^{(i)} - p_4^{(i)}\right) + p_4^{(i)}$, and $p_6^{(i)} = \Pr\left(j \in \mathcal{S}''^{(i)}|j \notin \mathcal{S}^{(i)}\right) = p_2^{(i)}\left(p_3^{(i)} - p_4^{(i)}\right) + p_4^{(i)}$.

Proof. The proof is similar to that of Lemma 1. The difference is that the possible outputs are index j being in and not in the final perturbed index set, namely $j \in \mathcal{S}''^{(i)}$ and $j \notin \mathcal{S}''^{(i)}$. We first compute two conditional probabilities $p_5^{(i)}$ and $p_6^{(i)}$, denoting the probabilities of j in the final perturbed index set given an index j is and is not in client i's real index set, respectively. In particular, we can derive $p_5^{(i)}$ through

$$\begin{split} p_5^{(i)} &= \operatorname{Pr} \left(j \in \mathcal{S}^{\prime\prime(i)} \middle| j \in \mathcal{S}^{(i)} \right) \\ &= \operatorname{Pr} \left(j \in \mathcal{S}^{\prime\prime(i)} \middle| j \in \mathcal{S}^{(i)}, j \in \mathcal{Y}^{(i)} \right) \operatorname{Pr} \left(j \in \mathcal{Y}^{(i)} \middle| j \in \mathcal{S}^{(i)} \right) \\ &+ \operatorname{Pr} \left(j \in \mathcal{S}^{\prime\prime(i)} \middle| j \in \mathcal{S}^{(i)}, j \in \mathcal{N}^{(i)} \right) \operatorname{Pr} \left(j \in \mathcal{N}^{(i)} \middle| j \in \mathcal{S}^{(i)} \right) \\ &= \operatorname{Pr} \left(j \in \mathcal{S}^{\prime\prime(i)} \middle| j \in \mathcal{Y}^{(i)} \right) \operatorname{Pr} \left(j \in \mathcal{Y}^{(i)} \middle| j \in \mathcal{S}^{(i)} \right) \\ &+ \operatorname{Pr} \left(j \in \mathcal{S}^{\prime\prime(i)} \middle| j \in \mathcal{N}^{(i)} \right) \operatorname{Pr} \left(j \in \mathcal{N}^{(i)} \middle| j \in \mathcal{S}^{(i)} \right) \\ &= p_3^{(i)} p_1^{(i)} + p_4^{(i)} \left(1 - p_1^{(i)} \right) \\ &= p_1^{(i)} \left(p_3^{(i)} - p_4^{(i)} \right) + p_4^{(i)}, \end{split}$$

where Equation (2) follows from the law of total probability, and Equation (3) follows that $j \in \mathcal{S}''^{(i)}$ is independent of $j \in \mathcal{S}^{(i)}$ conditioned on $j \in \mathcal{Y}^{(i)}$ or $j \in \mathcal{N}^{(i)}$. In a similar way, we can get $p_6^{(i)} = \Pr\left(j \in \mathcal{S}''^{(i)} \middle| j \notin \mathcal{S}^{(i)}\right) = p_2^{(i)}\left(p_3^{(i)} - p_4^{(i)}\right) + p_4^{(i)}$. Based on $p_5^{(i)}$ and $p_6^{(i)}$, we can still compute four ratios between the conditional probabilities of an instantaneous output given a pair of different inputs and draw the level of local differential privacy $\epsilon_1^{(i)} = \ln\left(\max\left(\frac{p_5^{(i)}}{p_6^{(i)}},\frac{p_6^{(i)}}{p_5^{(i)}},\frac{1-p_6^{(i)}}{1-p_6^{(i)}},\frac{1-p_6^{(i)}}{1-p_5^{(i)}}\right)\right)$. \square

From the above deduction, we can draw that $p_5^{(i)}$ and $p_6^{(i)}$ in the instantaneous randomized response play the same roles as $p_1^{(i)}$ and $p_2^{(i)}$ in the permanent randomized response. This intuition has been given in Section IV-B2, and is now formally verified here.

By combining the above two lemmas, we show the level of local differential privacy ensured by Algorithm 2.

Theorem 1. When client i is chosen to participate in an arbitrary number of communication rounds, Algorithm 2 satisfies $\epsilon^{(i)}$ -local differential privacy, where $\epsilon_1^{(i)} \leq \epsilon^{(i)} \leq \epsilon_{\infty}^{(i)}$.

Proof. We consider that client i participates in k communication rounds of federated submodel learning, and Algorithm 2 guarantees $\epsilon_k^{(i)}$ -local differential privacy. Thus, client i should generate k instantaneous randomized responses. On one hand, suppose that an attacker only leverages the k-th

instantaneous randomized response while ignoring all previous k-1 instantaneous randomized responses. This corresponds to the strongest possible local differential privacy guarantee, namely the lower bound on $\epsilon_k^{(i)}$. According to Lemma 2, a one-time instantaneous randomized response guarantees $\epsilon_1^{(i)}$ -local differential privacy. Therefore, $\epsilon_1^{(i)} \leq \epsilon_k^{(i)}$. On the other hand, if the attacker leverages all k instantaneous randomized responses, and as k approaches positive infinity, the worst case is that the attacker reveals the permanent randomized response. This corresponds to the weakest possible local differential privacy guarantee, namely the upper bound on $\epsilon_k^{(i)}$. By Lemma 1, the permanent randomized response can guarantee $\epsilon_k^{(i)}$ -local differential privacy. Hence, $\epsilon_k^{(i)} \leq \epsilon_k^{(i)}$. We complete the proof.

In fact, to derive the detailed local differential privacy guarantee when client i participates in k communication rounds, namely $\epsilon_k^{(i)}$, we need to make an additional assumption on how effectively the attacker infers client i's permanent randomized response from all k instantaneous randomized responses. We defer this explorative problem as our future work. Additionally, if we set $p_1^{(i)} = p_2^{(i)} = p_3^{(i)} = p_4^{(i)} = 1$, implying $p_5^{(i)} = p_6^{(i)} = 1$, $\epsilon_1^{(i)} = 0$, $\epsilon_\infty^{(i)} = 0$, and $\epsilon^{(i)} = 0$, this corresponds to that any index, no matter whether it is or is not in client i's real index set, will receive a permanent "Yes" answer and an instantaneous "Yes" answer. In other words, if client i takes the union of the chosen clients' real index sets as her perturbed index set, the local differential privacy guarantee of Algorithm 2 is 0, which is the strongest case, as in conventional federated learning.

We next analyze our secure submodel updates aggregation in Algorithm 1, according to Definition 2.

Theorem 2. Algorithm 1 is a client-controllable privacy protection mechanism for submodel updates aggregation. In particular, for any index $j \in \bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)}$, we let $n_{j,0}$ and $n_{j,1}$ denote the numbers of live clients who do not have and have j in reality, respectively. If each live client chooses the same probability parameters (i.e., $\forall i \in \hat{\mathcal{C}}, p_1^{(i)} = p_1, p_2^{(i)} = p_2, p_3^{(i)} = p_3, p_4^{(i)} = p_4, p_5 = p_1(p_3 - p_4) + p_4$, and $p_6 = p_2(p_3 - p_4) + p_4$, then Algorithm 1 can guarantee: Event 1 happens with probability $p_7 = p_5(1 - p_5)^{n_{j,1}-1}(1 - p_6)^{n_{j,0}}$, and Event 2 happens with probability $p_8 = (1 - p_5)^{n_{j,1}}(1 - (1 - p_6)^{n_{j,0}})$.

Proof. Event 1 happens when only one of the $n_{j,1}$ clients who have j in reality submits a nonzero update while all $n_{j,0}$ clients who do not have j in reality also do not submit zero updates. By the product rule, we can compute the joint probability of Event 1 as $p_7 = p_5(1-p_5)^{n_{j,1}-1}(1-p_6)^{n_{j,0}}$.

Event 2 happens when all $n_{j,1}$ clients who have j in reality do not submit nonzero updates. Under such a circumstance, if part (at least one) of $n_{j,0}$ clients who do not have j in reality submit zero updates, from the aggregate zero update, the cloud server almost ascertains that these clients do not have j in reality. According to the product rule, we can compute

the probability of Event 2 as $p_8 = (1 - p_5)^{n_{j,1}} (1 - (1 - p_6)^{n_{j,0}})$.

Theorem 2 enables participating clients to jointly adjust the privacy level of secure submodel updates aggregation by choosing different probability parameters. Details on finetuning are deferred to Appendix A. Additionally, we still examine the case that each client uses the union of the chosen clients' real index sets to upload the submodel update by setting $p_1 = p_2 = p_3 = p_4 = 1$, implying $p_5 = p_6 = 1$ and $p_7 = p_8 = 0$. This is the strongest possible clientcontrollable privacy for secure submodel updates aggregation, as in secure federated learning. Combining with the local differential privacy guarantee, we can see that secure federated submodel learning with the setting $p_1 = p_2 = p_3 = p_4 = 1$ holds the same security and privacy levels as secure federated learning under both Definition 1 and Definition 2. We further generalize this observation in Theorem 3, which is free of any security or privacy definition.

Theorem 3. If the probability parameters $p_1^{(i)}, p_2^{(i)}, p_3^{(i)}, p_4^{(i)}$ all take 1s for each chosen client $i \in C$, the security and privacy of secure federated submodel learning scheme in Algorithm 1 are as strong as secure federated learning.

Proof. We consider any index j from the full index set S (i.e., $j \in S$), and prove in two complementary cases as follows.

Case 1 $(j \in \bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)})$: In both secure federated submodel learning and secure federated learning, each client $i \in \mathcal{C}$ will download j-th row of the full model and then upload the update of this row to the cloud server through secure aggregation. Specifically, if j is not in a client's real index set, then this client will upload a zero vector. The whole processes of two schemes are consistent, implying the same security and privacy guarantees.

Case 2 $(j \notin \bigcup_{i \in \mathcal{C}} \mathcal{S}^{(i)})$: In secure federated submodel learning, each client will not download j-th row and also will not upload the zero vector. Thus, the adversary knows that each client doesn't not have index j (i.e., $\forall i \in \mathcal{C}, j \notin \mathcal{S}^{(i)}$), and each client's update is a zero vector. In contrast, in secure federated learning, each client will download j-th row and update a zero vector as her update. From the result that the aggregate update is a zero vector, the adversary still ascertains that each client does not have j and her update is a zero vector, which are the same as in secure federated submodel learning.

By summarizing two cases, we complete the proof.

We finally analyze the security of private set union. As a springboard, we give Lemma 3, which states that the modular addition of one or more random integers from $\mathbb{Z}_R = \{0,1,\ldots,R-1\}$ is still uniformly random in \mathbb{Z}_R . We note that the elementary operation in secure aggregation [18] is modular addition with a large modulus rather than original addition, which is consistent with Lemma 3. In addition, in the field of number theory, the set of integers $\mathbb{Z}_R = \{0,1,2,\ldots,R-1\}$ is called the least residue system modulo R, or the ring of the integers modulo R. Moreover, the set \mathbb{Z}_R together with the operation of modular addition form a finite cyclic group.

Lemma 3. For any nonempty set $C_1 \neq \emptyset$ and for all $i \in C_1$, r_i is randomly and independently chosen from $\mathbb{Z}_R = \{0, 1, \dots, R-1\}$, denoted as $r_i \in_R \mathbb{Z}_R$, then $\sum_{i \in C_1} r_i \mod R$ is still uniformly random in \mathbb{Z}_R .

Proof. We prove by induction on the cardinality of C_1 , denoted as $|C_1|$, where $|C_1| \ge 1$.

The base case is to show that the statement holds for $|\mathcal{C}_1| = 1$. We let $\{r_a\}$ denote the element, where $r_a \in_R \mathbb{Z}_R$. Thus, it is trivial that $\sum_{i \in \mathcal{C}_1} r_i \mod R = r_a \in_R \mathbb{Z}_R$.

The inductive step is to prove that if the statement for any nonempty $\bar{\mathcal{C}}_1 \subset \mathcal{C}_1$ holds, then the statement for $\bar{\mathcal{C}}_1 \bigcup \{b\}$ where $b \notin \bar{\mathcal{C}}_1, b \in \mathcal{C}_1$ still holds. We let r_a denote $\sum_{i \in \bar{\mathcal{C}}_1} r_i \mod R$. Hence, it suffices to show that $r_a \in_R \mathbb{Z}_R, r_b \in_R \mathbb{Z}_R \Rightarrow (r_a + r_b) \mod R \in_R \mathbb{Z}_R$. We prove this statement by showing $\Pr((r_a + r_b) = r \mod R) = 1/R$ for any $r \in \mathbb{Z}_R$. The detailed deduction is shown as follows:

$$\Pr\left((r_a + r_b) = r \mod R\right)$$

$$= \sum_{k=0}^{R-1} \Pr\left(r_a = k, r_b = r - k \mod R\right)$$

$$= \sum_{k=0}^{R-1} \Pr\left(r_a = k\right) \Pr\left(r_b = r - k \mod R\right) \qquad (4)$$

$$= \sum_{k=0}^{R-1} \frac{1}{R} \frac{1}{R}$$

$$= \frac{1}{R},$$

where Equation (4) follows that r_a, r_b are independent. By mathematical induction, we complete the proof.

Theorem 4. The private set union protocol in Algorithm 3 is proven secure in the sense that only the union of the chosen clients' real index sets is revealed.

Proof. We recall that in Algorithm 3, a client first represents her real index set as a Bloom filter $\mathbf{b}^{(i)}$, then replaces bit 1s with random integers, denoted as $\mathbf{b}'^{(i)}$, and finally executes secure aggregation with other clients under the coordination of the cloud server. Additionally, just in the same manner as computing the union, the client joins in the survey of whether there exists an index in her real index set falling into predivided partitions. Hence, for the sake of conciseness, we here only elaborate on secure union computation.

First, according to the security analysis in [18], the secure aggregation protocol is proven secure in both honest-but-curious and active adversary settings, where the adversaries can be the cloud server and participating clients. In more detail, from security and robustness, the secure aggregation protocol can guarantee that nothing but the aggregate result is revealed to the cloud server and all participating clients even if part of clients may drop out during the aggregation process. When the secure aggregation guarantee applies to our context, only $\sum_{i \in \mathcal{C}} \mathbf{b}'^{(i)}$ is revealed, while any individual $\mathbf{b}'^{(i)}$ for $i \in \mathcal{C}$ is concealed from both the cloud server and the other chosen clients $\mathcal{C} \setminus \{i\}$. Given that $\mathbf{b}'^{(i)}$ is a postprocessing of

 $\mathbf{b}^{(i)}$, the underlying Bloom filter $\mathbf{b}^{(i)}$ and thus each client *i*'s real index set $\mathcal{S}^{(i)}$ are obscured.

Second, we prove that the revealed sum $\sum_{i\in\mathcal{C}}\mathbf{b}'^{(i)}$ contains no information other than the union $\bigvee_{i\in\mathcal{C}}\mathbf{b}^{(i)}$ from the view of adversary with any prior knowledge. Considering both $\bigvee_{i\in\mathcal{C}}\mathbf{b}^{(i)}$ and $\sum_{i\in\mathcal{C}}\mathbf{b}'^{(i)}$ are element-wise operation, we thus just need to focus on one-dimensional cases of $\mathbf{b}^{(i)}$ and $\mathbf{b}'^{(i)}$, i.e., $\mathbf{b}^{(i)}$ degenerates to one bit $b^{(i)} \in \{0,1\}$ and $\mathbf{b}'^{(i)}$ degenerates to one random integer $b'^{(i)} \in \mathbb{Z}_R$. We now consider two complementary cases of $\{b^{(i)}|i\in\mathcal{C}\}$ as follows.

Case 1 ($\forall i \in \mathcal{C}, b^{(i)} = 0$): The union $\bigvee_{i \in \mathcal{C}} b^{(i)} = 0$ is the same as the $\sum_{i \in \mathcal{C}} b'^{(i)} = 0$. This implies nothing but the union is leaked from the sum in Case 1.

Case 2 ($\exists i \in \mathcal{C}, b^{(i)} = 1$): The union is $\bigvee_{i \in \mathcal{C}} b^{(i)} = 1$. We compute the sum $\sum_{i \in \mathcal{C}} b'^{(i)}$ by splitting \mathcal{C} into two parts: $\mathcal{C}_0 = \{i | i \in \mathcal{C}, b^{(i)} = 0\}$ and $\mathcal{C}_1 = \{i | i \in \mathcal{C}, b^{(i)} = 1\}$. Then, we can derive

$$\sum_{i \in \mathcal{C}} b'^{(i)} = \sum_{i \in \mathcal{C}_0} b'^{(i)} + \sum_{i \in \mathcal{C}_1} b'^{(i)}$$

$$= 0 + \sum_{i \in \mathcal{C}_1} b'^{(i)}$$

$$\in_{R} \mathbb{Z}_{R}. \tag{5}$$

Here, Equation (5) holds as follows. According to the antecedent $\exists i \in \mathcal{C}, b^{(i)} = 1$, we have $\mathcal{C}_1 \neq \emptyset$. Additionally, in Algorithm 3 (Line 4), if $b^{(i)} = 1$, we have $b'^{(i)} \in_R \mathbb{Z}_R$. Now, by using Lemma 3, we have $\sum_{i \in \mathcal{C}_1} b'^{(i)} \in_R \mathbb{Z}_R$. This indicates that the sum is just a uniformly random integer to the adversary's view. Furthermore, in the context of union, (1) this random integer takes positive values with probability 1 - 1/R and is further decoded as element "1"; (2) the random integer takes value 0 with a negligible probability 1/R and further is wrongly² decoded as element "0", i.e., a false negative happens with a negligible probability 1/R. Hence, the sum reveals nothing but the union in Case 2.

By incorporating no leakage of any individual real index set and no leakage but the union from the aggregate result, we complete the proof.

We give more interpretations about the word "only" in Theorem 4. With no exception, the count number of each element in the union (i.e., the cardinality of \mathcal{C}_1 in the proof) is also hidden. The reason is that for any non-empty \mathcal{C}_1 , the sum is uniformly random (i.e., Equation (5) holds). This further indicates that from the sum, all possible \mathcal{C}_1 's are computationally indistinguishable to the adversary. In other words, the adversary learns the exact cardinality with probability $1/|\mathcal{C}|$, which is the same as the probability of a random guess among all possible cardinalities $\{1, 2, \ldots, |\mathcal{C}|\}$.

B. Performance Analysis and Comparison

We now analyze the performance of our proposed secure federated submodel learning scheme. We first analyze the

²Even if each bit "1" is replaced with a positive rather than non-negative integer, the sum under modular addition can still take value 0 with a negligible probability. However, the sum can be no longer uniformly random.

TABLE III

COMPLEXITIES OF SECURE FEDERATED LEARNING (SFL) AND SECURE FEDERATED SUBMODEL LEARNING (SFSL) AT THE SAME SECURITY AND PRIVACY LEVELS AS WELL AS PRIVATE SET UNION (PSU) IN SFSL.

		Communication	Computation	Storage
Client	SFL SFSL PSU	$O(n+md) \ O(nsd) \ O(ns)$	$O(n^2 + nmd)$ $O(n^2sd)$ $O(n^2s)$	O(n+md) $O(nsd)$ $O(ns)$
Server	SFL SFSL PSU	$O(n^2 + nmd)$ $O(n^2sd)$ $O(n^2s)$	$O(n^2md)$ $O(n^3sd)$ $O(n^3s)$	$O(n^2 + md)$ $O(n^2 + nsd)$ $O(n^2 + ns)$

 $^{*|\}bigcup_{i\in\mathcal{C}}\mathcal{S}^{(i)}|\ll |\mathcal{S}|\Rightarrow ns\ll m.$

communication, computation, and storage (including both memory and disk loads) complexities of the client and the cloud server. Then, we introduce secure federated learning as a benchmark for comparison. For clarity, Table III summarizes the complexities of two secure schemes at the same levels of security and privacy along with the complexities of our proposed private set union protocol.

1) Performance of Secure Federated Submodel Learning: We focus on a certain communication round. For a concrete phase within the round, e.g., private set union or secure submodel updates and count vectors aggregations, we consider that all n chosen clients are alive at the beginning, but may drop out during the process, which imposes upper bounds on the overheads of the phase. In addition, for feasibility and clarity in analysis, we let each client choose the same probability parameters p_1, p_2, p_3, p_4 , resulting in the same p_5, p_6 . Moreover, we assume that the expected cardinality of each client's real index set is s, which indicates that the expected cardinality of the union of n chosen clients' real index sets $\bigcup_{i \in C} S^{(i)}$ is upper bounded by ns. Here, ns is much less than the cardinality of the full index set m. Furthermore, the expected cardinality of each client's perturbed index set $S''^{(i)}$ is upper bounded by $sp_5 + (n-1)sp_6$, and the expected cardinality of each client's succinct index set $S^{(i)} \cap S''^{(i)}$ is sp_5 . In what follows, we first present the detailed complexity formulas containing the probability parameters p_5, p_6 and then instantiate with $p_5 = p_6 = 1$. This corresponds to that each client uses the union of all chosen clients' real index sets to interact with the cloud server. As demonstrated in Theorem 3, this case can provide the same levels of privacy and security as secure federated learning. Further, given $0 \le p_5, p_6 \le 1$ while the complexities are non-decreasing with p_5 , p_6 , this case still upper bounds the complexities of our scheme.

First regards the communication complexities of the client and the cloud server. Each client's communication overhead can be broken up as: (1) Participating in the private set union protocol, where two vectors need to be securely aggregated, and the final union result needs to be downloaded with size O(ns). The first vector is the perturbed Bloom filter with size β , and the second vector is the perturbed indicator vector with a preset constant size. According to Equation (1), the optimal

length of Bloom filter β is proportional to the cardinality of the set to be filtered ϕ , here the cardinality of the union of n chosen clients' real index sets, which implies that $\beta \propto$ $\phi = O(ns)$. Therefore, the communication complexity of each client in private set union is O(n + ns + ns) = O(ns), where O(n+ns) corresponds to the cost of securely aggregating the perturbed Bloom filters and is obtained by letting the size of vector l in Table II take O(ns); (2) Downloading the training hyperparameters as well as the submodel with the perturbed index set $S''^{(i)}$, namely a $(sp_5 + (n-1)sp_6) \times d$ matrix; (3) Uploading the weighted submodel update and corresponding count vector with respect the perturbed index set through secure aggregation. In particular, the total size of vector to be securely aggregated with at most n-1 other clients is $(sp_5 +$ $(n-1)sp_6(d+1)$. According to Table II, the communication complexity of this part is $O(n+(sp_5+(n-1)sp_6)(d+1))$. In summary, the total communication complexity of each client is $O(ns)+O((sp_5+(n-1)sp_6)d)+O(n+(sp_5+(n-1)sp_6)(d+$ 1)) = $O(ns + (sp_5 + (n-1)sp_6)(2d+1))$. If the probability parameters p_5, p_6 both take the value 1, the communication complexity of each client is O(nsd). Correspondingly, the cloud server's communication cost can be broken up into: (1) Working as a mediation in the private set union protocol, and delivering the union result to clients with $O(n^2s)$ communication cost. Hence, the communication complexity of the cloud server in private set union is $O(n^2 + n^2s + n^2s) = O(n^2s)$, where $O(n^2 + n^2 s)$ is related to secure aggregation; (2) Returning training hyperparameters and requested submodels to clients with $O(n + n(sp_5 + (n-1)sp_6)d)$ communication cost; (3) Working as a mediation in the secure aggregations of weighted submodel updates and count vectors. Its communication complexity is $O(n^2 + n(sp_5 + (n-1)sp_6)(d +$ 1)). Overall, the cloud server's communication complexity is $O(n^2s) + O(n + n(sp_5 + (n-1)sp_6)d) + O(n^2 + n(sp_5 + n(sp_5$ $(n-1)sp_6(d+1) = O(n^2s + n(sp_5 + (n-1)sp_6)(2d+1)).$ We still examine that when $p_5 = p_6 = 1$, the communication complexity of the cloud server is $O(n^2sd)$.

Second regards time complexities. Despite of the local training phase, the computation cost of each client is dominated by: (1) Perturbing the real index set, which costs O(ns) time. Here, one set lookup operation can be implemented with O(1)complexity; (2) Participating in the secure aggregation based stages, including private set union and secure weighted submodel updates and count vectors aggregations, which consume $O(n^2+n^2s) = O(n^2s)$ and $O(n^2+n(sp_5+(n-1)sp_6)(d+1))$ time, respectively. Thus, the overall time complexity of each client is $O(ns) + O(n^2s) + O(n^2 + n(sp_5 + (n-1)sp_6)(d +$ 1) = $O(n^2s + n(sp_5 + (n-1)sp_6)(d+1))$. When $p_5 = p_6 = 1$, it turns to $O(n^2sd)$. Correspondingly, the computation overhead of the cloud server is dominated by mediating private set union and secure aggregations of weighted submodel updates and count vectors, which consume $O(n^2ns) = O(n^3s)$ and $O(n^2(sp_5 + (n-1)sp_6)(d+1))$ time, respectively, and $O(n^3s + n^2(sp_5 + (n-1)sp_6)(d+1))$ time in total. When $p_5 = p_6 = 1$, the cloud server's time complexity is $O(n^3 s d)$.

Third regards storage complexities. Each client's storage

overhead comes from: (1) Storing the materials in private set union and in secure aggregations of weighted submodel updates and count vectors, which occupy O(n+ns) = O(ns)and $O(n + (sp_5 + (n-1)sp_6)(d+1))$ space, respectively; (2) Storing her permanent answers, which occupies O(ns)space. Thus, the storage overhead of each client is O(ns) + $O(n + (sp_5 + (n-1)sp_6)(d+1)) + O(ns) = O(ns + (sp_5 + (n-1)sp_6)(d+1))$ $(sp_5 + (n-1)sp_6)(d+1)$ in total. When $p_5 = p_6 = 1$, the client does not need to store the permanent answers any more, and her storage overhead is O(nsd). Correspondingly, the cloud server's storage overhead can be broken up into: (1) Storing the materials in private set union and in secure weighted submodel updates and count vectors aggregations, which occupy $O(n^2+ns)$ and $O(n^2+(sp_5+(n-1)sp_6)(d+1))$ space, respectively; (2) Storing n chosen clients' perturbed index sets, which occupies $O(n(sp_5 + (n-1)sp_6))$ space. In summary, the overall storage overhead of the cloud server is $O(n^2+ns)+O(n^2+(sp_5+(n-1)sp_6)(d+1))+O(n(sp_5+(n-1)sp_6)(d+1)+O(n(sp_5+(n-1)sp_6)(d+1)+O(n(sp_5+(n-1)sp_6)(d+1)+O(n(sp_5+(n-1)sp_6)(d+1)+O(n(sp_5+(n-1)sp_6)(d+1)+O(n(sp_5+(n-1)sp_6)(d+1)+O(n(sp_5+$ $(1)sp_6) = O(n^2 + ns + (sp_5 + (n-1)sp_6)(n+d+1))$. When $p_5 = p_6 = 1$, the cloud server does not need to store clients' perturbed index sets, and her storage overhead is $O(n^2 + nsd)$.

2) Comparison with Secure Federated Learning: We first analyze the complexities of secure federated learning. Each client's communication overhead mainly comes from: (1) Downloading the full model, namely a $m \times d$ matrix; (2) Uploading the update of the full model through the secure aggregation protocol with O(n+md) complexity, which is obtained by letting the size of vector l in Table II take md. Thus, the overall communication complexity of each client is O(n+md). In correspondence, the cloud server's communication cost comes from: (1) Sending the full model to all nclients, with complexity O(nmd); (2) Working as a mediation in the secure aggregation of the full model updates, with complexity $O(n^2+nmd)$. Thus, the cloud server's communication complexity is $O(n^2 + nmd)$. Regardless of the local training phase, the computation and storage overheads of each client and the cloud server are dominated by secure aggregation. In particular, the time complexities of each client and the cloud server are $O(n^2 + nmd)$ and $O(n^2md)$, respectively. Additionally, the storage complexities of each client and the cloud server are O(n+md) and $O(n^2+md)$, respectively.

We next compare our secure scheme with secure federated learning at the same security and privacy levels, i.e., all the probability parameters in our scheme are set to 1. Given that the cardinality of the union of n chosen clients' real index sets is much smaller than the cardinality of the full index set (i.e., $|\bigcup_{i\in\mathcal{C}}\mathcal{S}^{(i)}|\ll|\mathcal{S}|\Rightarrow ns\ll m$), we can draw from Table III that the complexities of each client and the cloud server in our secure federated submodel learning are both much lower than those in secure federated learning. We can further derive from Table III that our secure scheme is quite scalable because its complexities are independent of the size of the full model, which is controlled by the number of total rows m.

We finally compare the overheads, typically the computation and memory overheads, of each client in the local training phase under two secure frameworks. We mainly focus on the

TABLE IV Statistics of Taobao dataset.

Туре	#User(s)	#Goods	#Categories	#Samples
Test (Full)	24,790	138,829	4,758	1,010,284
Train (Full)	49,023	143,534	4,815	15,854,357
Train (Per Client)	1	301	117	323

size of local model/submodel. For secure federated learning, each client trains the full model \mathbf{W} with size md. In contrast, for our secure federated submodel learning, each client trains the succinct submodel $\mathbf{W}_{\mathcal{S}^{(i)} \bigcap \mathcal{S}''^{(i)}}$ with size sp_5d , which is sd when $p_5 = 1$. Given $md \gg nsd > sd \geq sp_5d$, we can find that our secure scheme is still far more efficient than secure federated learning in the local training phase.

VI. EVALUATION

In this section, we introduce our evaluation results from model accuracy and convergency, practical computation, communication, and storage overheads.

Dataset: We use an industrial dataset built from 30-day impression and click logs of Taobao users from June 15, 2019 to July 15, 2019. For a certain Taobao user, we leverage her click behaviors in previous 14 days as historical data to predict her click and non-click behaviors in the following 1 day. We leave out the samples within the last 1 day as the target items of the test set while putting the others into the training set. Specifically, the test set is located at the cloud server to judge the accuracy and convergency of the global model. In contrast, for the full training set, we further cluster each Taobao user's data as a training set located at a client. Some statistical information about the numbers of Taobao users, goods IDs, category IDs, and samples, in the test set, the full training set, and one client's training set is shown in Table IV.

Model, Hyperparameters, and Metrics: We take the Deep Interest Network (DIN) [4] as the e-commerce recommendation model for federated submodel learning, where the number of columns in the embedding matrix is set to 18. Except the embedding layer for user IDs, goods IDs, and category IDs, the parameters of the other network layers in DIN, including the attention layer and the fully connected layer, are of size 64,327. Hence, the global model at the cloud server is of size 3,617,023, whereas a desired submodel at the client is of size 71,869 in average, which is only 1.99% of the global model's size and roughly requires 0.27MB space using 32bit representation. For each client's local training, we choose mini-batch SGD as the optimization algorithm, set the batch size to 2, and set the local epoch number to 1. In addition, we initially set the learning rate to 1 and further apply exponential decay with the decay rate of 0.999 per communication round. For the cloud server's global testing, we adopt a golden metric in the task of Click-Through Rate (CTR) prediction, called Area Under the Curve (AUC), and set the batch size to 1,024.

Prototypes and Configurations: We implemented prototypes of our secure federated submodel learning and secure

federated learning in Python 2.7.16. Due to the synchronization requirement of secure aggregation, we adopted a synchronous architecture on top as suggested by Google's deployment practice [40], where the cloud server delivers instructions to chosen clients, waits for them to finish their tasks, and then moves on to the next step. In particular, we implemented a communication module between the cloud server and each client with standard socket programming. In addition, we used TensorFlow 1.12.0 to implement DIN. Moreover, we mainly used PyCryptodome 3.7.3 to implement the secure aggregation protocol: Diffie-Hellman key exchange was implemented over RFC 3526 Group 14, which is a 2048bit modular exponential (MODP) group [96]; the secret sharing adopted the standard Shamir's version; the authenticated encryption used Advanced Encryption Standard (AES) in the Cipher Block Chaining (CBC) mode, where the secret key is set to be 128-bit long; Pseudo-Random Number Generator (PRNG) was implemented using Hash-based Message Authentication Code (HMAC) suggested by NIST SP 800-90A [97], and the security strength takes 128 bits.

Our running environment is a Linux workstation with 64-bit Ubuntu 18.04.2 OS. The processor is Intel(R) Core(TM) i9-9900K with 8 cores, the base frequency is 3.60GHz, the memory size is 64GB, and the cache size is 16MB. The workstation is also equipped with 2 NVIDIA's GeForce RTX 2080 Ti graphics cards. In our evaluation, to manifest the resource differentiation between clients and the cloud server, from hardware, we ran all the clients only on CPU, but allowed the cloud server to accelerate her operations using GPU. Additionally, from parallelism and concurrency, we optimized some of the cloud server's hotspot functions with the multiprocessing library in Python.

Implementation Overview: We revisit some key points and also present more implementation details when our secure federated submodel learning is instantiated with Taobao's e-commerce recommendation. The cloud server chooses nclients, namely n Taobao users, in every communication round, where n can range from 20 to 100 with a step of 20. Each chosen client extracts goods IDs from her data as her real index set. Then, she joins in private set union to obtain the union of n chosen clients' goods IDs. In our Taobao dataset, the maximum size of the union of 100 clients' goods IDs reaches 32,904, which corresponds to the cardinality of set to be filtered ϕ . From Equation (1), we can derive that the optimal length of Bloom filter β should be 630,774 at a desired false positive rate of 0.01%. This is larger than the total number of goods IDs 143,534 involved in this Taobao dataset. Thus, we set the length of Bloom filter to be 143,534, take an identity map as the only hash function, and omit the partition steps. The false positive rate now is 0. In addition, we set the modulus R to 2^{32} in private set union, and any false negative due to modular addition in recovering union never happened throughout our evaluation. It is worth noting that our evaluation results of private set union presented here can sufficiently embody the practical case, where the full size of the goods IDs in Taobao scales to billions. The reason is

TABLE V
CHOICES OF PROBABILITY PARAMETERS (CPPS) AND RESULTING
PRIVACY LEVELS.

	p_1, p_3	p_2, p_4	p_5	p_6	ϵ_1	ϵ_{∞}	p_7	p_8
CPP1	1	0	1	0		∞	86.7%	0
CPP2 CPP3	15/16 7/8	$\frac{1/16}{1/8}$	88.3% 78.1%	11.7% $21.9%$	$\begin{array}{ c c c } 2.02 \\ 1.27 \end{array}$	$\frac{2.71}{1.95}$	0	10.3% $19.5%$
CPP4	3/4	1/4	62.5%	37.5%	0.51	1.10	0	34.2%
CPP5	1	1	1	1	0	0	0	0

*Smaller ϵ_1 , ϵ_∞ , p_7 , and p_8 indicate better privacy.

that as shown in Table III, the overheads of our private set union scheme only hinge on the number of chosen clients n and the optimal length of Bloom filter β , where the latter is independent of the domain of the represented set's elements and is at the same level of 143,534 evaluated here.

After obtaining the union of goods IDs, each client runs Algorithm 2 to generate a set of perturbed goods IDs. For clarity in presenting results, we let each client use the same Choice of Probability Parameters (CPP). Table V lists 5 CPPs in our evaluation and their resulting privacy levels of index set perturbation and secure submodel updates aggregation, where the number of chosen clients is 100 for the secure submodel updates aggregation. We provide some insights from Table V as follows: (1) CPP1 corresponds to the worst local differential privacy guarantee where each client reveals her real goods IDs, while CPP5 corresponds to the best privacy guarantee as strong as secure federated learning where each client uses the union of chosen clients' goods IDs as her perturbed index set. Additionally, as the serial number of CPP becomes larger, the local differential privacy guarantee is stronger; (2) the resulting $p_7 = 86.7\%$ at CPP1 indicates that if each chosen client submits the submodel update using her real goods IDs, then Event 1 (i.e., from the securely aggregated submodel update, the cloud server ascertains that some goods ID belongs to a concrete client and also learns the detailed update with respect to this goods ID.) happens with probability 86.7%. This further implies that 86.7% of the goods IDs in the union of 100 chosen clients' real goods IDs involve single client. Therefore, we can draw that user data in our e-commerce context is highly heterogeneous, and the truly required submodels of different clients are highly differentiated; (3) with CPP from CPP2 to CPP5, the cloud server cannot ensure that any goods ID belongs to some client from the aggregate submodel update and also cannot learn any client's detailed update, namely Event 1 happens with probability $p_7 = 0$; and (4) with CPP from CPP2 to CPP4, the proportion of each client's real goods IDs falling into her perturbed goods IDs controlled by p_5 decreases, while the proportion of redundant goods IDs controlled by p_6 increases. Hence, the proportion of aggregate zero updates grows, implying a higher probability p_8 of the cloud server in ascertaining that a certain goods ID does not belong to some client. Further, given an observation that p_8 is approaching $1 - p_5$, we can still draw that the real goods IDs of different Taobao users are highly differentiated.

Based on the perturbed goods IDs, each client can further generate the perturbed category IDs by leveraging the global mapping between goods IDs and category IDs, which is publicly shared among all clients and the cloud server. Upon receiving the perturbed goods IDs from a client, the cloud server can still generate her perturbed category IDs, thus returning a submodel to the client. Here, the submodel returned to the client comprises the embedding vectors for her perturbed goods and category IDs as well as the other network parameters of DIN. In addition, the client should not generate perturbed goods and category IDs independently by running Algorithm 2 twice. The reason is that there exist strong correlations, particularly a surjective but not injective mapping, between goods IDs and category IDs, whereas independent randomized responses to dependent questions can lead to inconsistency and may leak information about true answers to the adversary. Therefore, we should apply randomized response to major questions and then generate answers to secondary/correlated questions via their dependence relationships, which can keep consistency. For example, suppose you are asked two types of questions: one type is about whether you have a certain fruit by enumerating all possible fruits, such as "Do you have an apple?" or "Do you have a banana?"; the other type is about whether you have the category fruit. To preserve plausible deniability and consistency of your answers, you should simply apply randomized response to the first type of questions, and then directly respond to the second question based on the noisy answers to the first type of questions. In other words, if you respond any "Yes" to the first type of questions, then you should respond "Yes" to the second question; otherwise, you should respond "No".

By performing set intersection between real and perturbed goods IDs, each client obtains her succinct set of goods IDs and also gets her succinct set of category IDs based on the global goods-category mapping. Then, the client extracts her succinct submodel from the downloaded submodel and also extracts her succinct training set from the original training set by following two rules: (1) For the goods ID to be predicted in a sample, if it does not belong to the succinct set of goods IDs, this sample will be filtered out; and (2) for the sequence of historical goods IDs in a sample, we only keep those goods IDs in the succinct set of goods IDs as well as their corresponding category IDs. Of course, if none of goods IDs are left in the historical sequence, this sample will be filtered out. Afterwards, each client trains her succinct model, obtains the update of the submodel, and prepares the weighted submodel update and the count vector to be uploaded.

To facilitate the cloud server in obliviously adding up the weighted submodel updates based on secure aggregation, the float-type parameters need to be converted to integers so that modular addition is supported. One common practice is that each client first scales up the parameters through multiplying by a large constant (e.g., a power of 2), and only keeps the integral parts. Later, the cloud server scales down the aggregate result through dividing by the same constant. Different from the common scaling technique, in this work, we perform float-

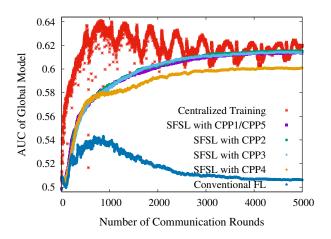


Fig. 4. Accuracies and convergencies of global model in centralized training, Secure Federated Submodel Learning (SFSL) with different Choices of Probability Parameters (CPPs), and conventional Federated Learning (FL).

to-integer conversion in a new manner, particularly by means of a celebrated model compression algorithm, called stochastic γ -level quantization [29], which maps each parameter into $\{0,1,\ldots,2^{\gamma}-1\}$. The detailed procedure is shown as follows: Each client first compresses her submodel update and then weights the parameters in the compressed submodel update by multiplying with corresponding count numbers. After the secure aggregation, the cloud server gets the sum of the weighted compressed submodel updates, and then divides by corresponding total count numbers. Finally, the cloud server applies the decompression algorithm and updates the global model by adding the aggregate submodel update. We note that the postprocessing with a weighted average/mean in the compressed space does not introduce any bias/error to the final decompressed result. Detailed proof is deferred to Appendix B. Furthermore, in our following evaluation, we set γ to 2^{15} and set the modulus in the secure aggregation of weighted compressed submodel updates to 2^{32} .

A. Model Accuracy and Convergency

We bring in centralized training and conventional federated learning as two baselines and plot their AUCs as well as the AUCs of our secure federated submodel learning with different CPPs in Fig. 4. Here, the number of chosen clients in each communication round n is set to 100, and the total number of communication rounds is set to 5,000. In addition, centralized training refers to the traditional case that the cloud server first collects data from all clients, then trains the DIN model, and tests the model once training over the samples with a similar size to the total size of n chosen clients' datasets.

From Fig. 4, we can see that compared with centralized training, which reaches the highest AUC of 0.641 in 803 communication rounds, our secure federated submodel learning with CPP2 achieves the highest AUC of 0.615 in 4,908 communication rounds, slightly decreasing by 0.026. In contrast, conventional federated learning performs worst among all schemes, only achieving the highest AUC of 0.543 in 867

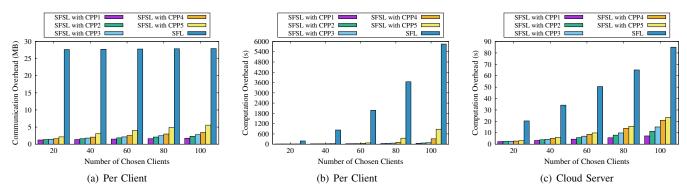


Fig. 5. Total communication and computation overheads of the client and the cloud server per communication round in Secure Federated Submodel Learning (SFSL) with different Choices of Probability Parameters (CPPs) and Secure Federated Learning (SFL).

communication rounds. Even worse, it does not converge to a good model at the end of many communication rounds³. The major reason for conventional federated learning not working well in the e-commerce recommendation context is that it coarsely computes the weighted average of the clients' updates of the full model proportional to their training set sizes, no matter whether one client's whole training set actually involves some network parameters (the full model excluding her submodel, e.g., some embedding vectors here), thus inaccurately counting in the weights (i.e., the training set sizes) of those clients who contribute zero/no updates for these network parameters. In addition, with a higher heterogeneity of user data and thus a higher differentiation of submodels, the roughness and inaccuracy of conventional federated learning will be exposed more completely, which clarifies why it can work in the natural language context with a 10,000 word vocabulary considered by Google, but does not work well in our e-commerce context with billion-scale goods IDs. More thorough demonstrations are deferred to Appendix C.

From Fig. 4, we can also observe some differences among the AUCs of our secure federated submodel learning with different CPPs. In particular, secure federated submodel learning with CPP4 is the worst among all CPPs, and it achieves the highest AUC of 0.601, decreasing by 0.014 compared with CPP2. We clarify the reason through the resulting probabilities p_5 's of different CPPs listed in Table V, where p_5 denotes the probability that an index in a client's real index set finally is put into her perturbed index set, and it dominates the size of the client's succinct training set. Further considering the process of generating the succinct training set in our evaluation, in addition to the size of each client's factual local training samples, p_5 also controls the length of historical goods and category IDs in every sample. Thus, a smaller p_5 tends to imply worse model performance in general. This accounts for different model performances under CPP2 and CPP4, and also explains the observation that CPP1 and CPP5, sharing the same $p_5 = 1$, have identical model performance.

B. Communication Overhead

We show the total communication overhead of secure federated submodel learning and first introduce secure federated learning as a baseline. Fig. 5(a) plots the communication overhead per client per communication round. Here, we do not plot the communication overhead of the cloud server, since it is equal to the communication overhead per client multiplying by the number of chosen clients. In more detail, the incoming data of the cloud server is exactly the total outgoing data of all chosen clients, and vice verse. Additionally, we also do not plot for different dropout ratios because this factor has little impact on the bandwidth cost.

One key observation from Fig. 5(a) is that our secure federated submodel learning can significantly reduce the total communication overhead, compared with secure federated learning. In particular, when the number of chosen clients is 100, the communication overheads per client per communication round are 1.76MB, 2.33MB, 2.78MB, 3.40MB, and 5.57MB in secure federated submodel learning with CPP1, CPP2, CPP3, CPP4, and CPP5, respectively, reducing 93.72%, 91.65%, 90.06%, 87.81%, and 80.05% than secure federated learning, which incurs 27.94MB per client per round. Considering secure federated submodel learning with CPP5 share the same levels of security and privacy with secure federated learning (i.e., Theorem 3), we can draw that our secure scheme can reduce communication overhead even with not scarifying any security or privacy. These results coincide with our complexity analysis in Section V-B and Table III.

The second key observation from Fig. 5(a) is that for secure federated submodel learning with a certain CPP, the communication overhead per client increases with the number of chosen clients n. In addition, for a certain number of chosen clients, the communication overhead per client increases with the serial number of CPP. We clarify the reasons by adopting the detailed communication complexity formula of each client from Section V-B: $O(ns+(sp_5+(n-1)sp_6)(2d+1))$. On one hand, the communication complexity grows linearly with n. On other other hand, it is increasing with p_5 and p_6 , and thus

³We have tried several different pairs of an initial learning rate and a decay rate and all observed divergences in conventional federated learning. For example, when the initial learning rate is 4, and the exponential decay rate is 0.996, conventional federated learning achieves the best AUC of 0.554 in 230 rounds but diverges to 0.503 at the end of 5,000 rounds.

TABLE VI
COMMUNICATION AND COMPUTATION OVERHEADS OF THE CLIENT AND
THE CLOUD SERVER PER ROUND IN PRIVATE SET UNION.

#Chosen Clients n	20	40	60	80	100
Client's Comm. Overhead (MB)	0.63	0.70	0.77	0.84	0.91
Client's Comp. Overhead (s)	5.04	10.09	15.53	26.49	33.69
Server's Comp. Overhead (s)	1.12	1.68	2.40	3.13	4.19

CPP5 is most communication expensive. Additionally, given $p_5+p_6=1$ for CPPs from CPP1 to CPP4 in Table V, we can simplify the formula to $O(ns+(s+(n-2)sp_6)(2d+1))$, which increases with p_6 for n>2. From Table V, we can see that p_6 increases with the serial number of CPP, implying a higher communication overhead as depicted in Fig. 5(a). Intuitively, p_6 denotes the probability that an index not in a client's real index set finally falls into the perturbed index set and controls the size of the redundant/zero parameters to be downloaded and securely uploaded. Thus, when the sum of two probabilities is fixed, particularly $p_5+p_6=1$, the increase of bandwidth cost due to introducing redundant parameters exceeds the decrease due to discarding original parameters. In other words, the introduced redundant parameters controlled by p_6 dominates the holistic trend of communication overhead.

We next introduce the pure versions of federated submodel learning and conventional federated learning shown in Fig. 1 as another type of baselines, and investigate the expansion factors due to the security and privacy guarantees. In particular, the pure federated submodel learning (resp., federated learning) means that each client directly downloads her required submodel (resp., the full model) from the cloud server, and then uploads the submodel update and the count vector (resp., the full model update and the training set size) to the cloud server. The communication overheads per client per round are 0.41MB and 20.70MB in the federated submodel learning and federated learning, respectively, and are irrelevant with the number of chosen clients. Compared with the pure version, when the number of chosen clients is 100, the communication overhead of secure federated submodel learning with CPP2 (resp., secure federated learning) expands $5.65 \times$ (resp., $1.35\times$). Three major reasons account for a larger expansion factor in secure federated submodel learning: (1) First is that the bandwidth cost of pure federated submodel learning, as the denominator, is much lower than, particularly 1.99\% of, the pure federated learning's bandwidth cost; (2) second is that the size of model parameters in secure federated learning is much larger than that in secure federated submodel learning, which can amortize the communication cost spent in transferring security and privacy related parameters; and (3) third is that secure federated submodel learning requires an extra process of private set union to facilitate later index set perturbation.

We finally present the communication overhead per client per round of our private set union protocol. Table VI lists the detailed bandwidth cost. We can see that the communication overhead per client in private set union increases linearly with the number of chosen clients, roughly with an increase of 0.07MB per 20 clients. In addition, we can also see that our private set union is communication efficient, and it only incurs 0.91MB when the number of chosen clients reaches 100. These evaluation results conform to our complexity analysis in Section V-B and Table III.

C. Computation Overhead

We now report the practical computation overhead, mainly by investigating the effects of the number of chosen clients, the choices of probability parameters, as well as the ratios of dropout. To be consistent with our time complexity analysis, the computation overhead here only includes the run time of the client or the cloud server in executing the protocol but ignores synchronization delay and the time overhead of testing global model. Of course, given mobile clients are highly parallel in practice, the total run time per communication round can be estimated by adding up the computation overheads of the client and the cloud server shown here. In addition, testing the global model per round costs the cloud server 32.12s.

We first show the computation overheads of the client and the cloud server in our secure federated submodel learning with different CPPs in Fig 5(b) and Fig 5(c), respectively. We still introduce secure federated learning as a baseline. First, we compare two figures and find that the computation overhead per client is higher than that of the cloud server. For example, in secure federated submodel learning with CPP2, when the number of chosen clients is 100, the computation overheads of the client and the cloud server are 71.80s and 11.31s, respectively. In addition to the superiorities of hardware and multiprocessing, under the synchronous architecture, the cloud server actually takes up the entire computing resources of the physical workstation alone, which are instead shared by all n chosen clients simultaneously. These factors jointly account for the differences between Fig 5(b) and Fig 5(c). Second, we focus on a certain side, either the client or the cloud server, and we can observe that her computation overhead grows with the number of chosen clients or the serial number of CPP. Here, we explain the reason by means of the detailed time complexities in Section V-B. For example, the time complexity of the client is $O(n^2s + n(sp_5 + (n-1)sp_6)(d+1))$, which is increasing with the number of chosen clients n as well as p_5 and p_6 , where the latter explains why CPP5 is the most timeconsuming one among all CPPs. Regarding CPP1 to CPP4 where $p_5 + p_6 = 1$, we can simplify the complexity formula to $O(n^2s+n(s+(n-2)sp_6)(d+1))$, which is increasing with p_6 for n > 2. The intuition behind the above analysis is analogous to that behind communication overhead, i.e., the size of redundant/zero parameters dominates the holistic computation overhead. Third, we compare our secure federated submodel learning with secure federated learning, and can find that our secure scheme significantly outperforms the baseline on both the sides of the client and the cloud server. In particular, when the number of chosen clients is 100, at the same security and privacy levels, secure federated submodel learning with CPP5 reduces 85.02% and 72.51% of computation overheads than

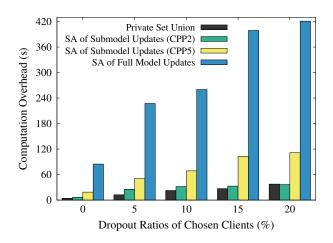


Fig. 6. Computation overheads of the cloud server with varying dropout ratios of chosen clients, in private set union, Secure Aggregations (SAs) of the weighted submodel updates and count vectors in SFSL with different CPPs, and Secure Aggregation (SA) of the full model updates in SFL.

secure federated learning on the sides of the client and the cloud server, respectively. When security and privacy become weaker, the advantages of our scheme are more evident under CPP1 to CPP4, e.g., CPP2 reduces 98.77% and 86.70% on the sides of the client and the cloud server, respectively.

We next investigate the effect of the dropout ratio of chosen clients and depict the evaluation results in Fig. 6. Here, we mainly focus on the secure aggregation based stages while ignoring the other stages which are irrelevant with dropout. In addition, we fix the number of chosen clients in each communication round at 100. Moreover, each client can randomly go offline after sending the encrypted shares of her private PRNG seed for the self mask and her private key for the mutual mask to the other clients because only this type of dropout behaviors introduce additional overhead. In particular, live clients still incorporate the public keys of dropped clients for mutual masks, and the cloud server must perform an expensive recovery computation to remove these mutual masks. The higher the dropout ratio is, the more costly the recovery computation will be. This trend is clear in Fig. 6. Last, we only report the computation overhead of the cloud sever because dropped clients do not introduce additional operation cost to live clients, like for communication overhead. Of course, the case that more clients are dropped can mitigate the competition for system resources, especially when the number of chosen clients is large and the size of vector for secure aggregation is large. For example, when the number of chosen clients is 100 and the dropout ratio is 20\%, for the secure aggregations of weighted submodel updates and count vectors in our secure scheme with CPP5, the computation overhead per client reduces by 16.75\% than the overhead in the case of no dropout.

We further observe Fig. 6 more carefully. First, we compare our secure scheme with secure federated learning in the stages of obliviously summing submodel or full model updates. We can find that our scheme significantly outperforms secure federated learning at any dropout ratio. Specifically, when the dropout ratio is 20%, our scheme with CPP2 and CPP5 reduce 91.33% and 73.55% of computation overhead, respectively. Second, we examine our private set union protocol and can see that it is quite efficient, even when the dropout ratio is high. In particular, when the dropout ratio reaches 20%, the computation overhead of the cloud server is 37.66s.

D. Memory and Disk Loads

We finally present the practical storage overheads of secure federated submodel learning and secure federated learning, including memory and disk loads. In particular, the materials for private set union and secure aggregations of submodel/full model updates and count vectors are stored in memory for immediate use, whereas permanent answers of each client are written into disk if necessary.

First is about memory overhead. The cloud server requires the video memory of 551MB, mainly for testing the global model at the end of each communication round, which is the same for all schemes. In addition, when the number of chosen clients is 100, and there is no dropout, the memory overheads per client are 209MB and 281MB in our secure federated submodel learning with CPP2 and CPP5, respectively, reducing 59.40% and 45.43% than secure federated learning. Correspondingly, the memory overheads of the cloud server are 1.58GB and 3.15GB, respectively, reducing 81.88% and 63.77% than secure federated learning. Furthermore, compared with the pure version, the memory overheads of our secure federated submodel learning with CPP2 and CPP5 only slightly expand on the client's side, and expand $1.19\times$ and $2.38\times$ on the cloud server's side, respectively. In contrast, the expansion factor is $4.60\times$ on the cloud server's side in secure federated learning.

Regarding the disk load, only the client in our scheme with CPP2, CPP3, and CPP4 needs to store her permanent answers in index set perturbation for multiple communication rounds, which roughly occupies the disk space of 280KB within the total 5,000 rounds.

E. Summary and Discussion

The above evaluation results adequately demonstrate the effectiveness and efficiency of our proposed secure federated submodel learning as well as its superiority over the baselines of conventional federated learning. Additionally, when the size of the full model, depending on the full size of the goods IDs in the e-commerce scenario, scales further to billions in practice, it has no effect on the performance and overhead of our scheme. This is because as analyzed in Section V-B and summarized in Table III, the complexities of our scheme are independent of the size of the full model. However, the conventional federated learning framework, hinging on the full model, will be too prohibitively inefficiently to be applicable.

VII. CONCLUSION

In this paper, we have proposed a new framework, called secure federated submodel learning, for numerous clients to effectively and efficiently train large-scale deep learning models under the coordination of an untrusted cloud server while keeping their user data private. We further have applied our framework to the e-commerce recommendation scenario of Alibaba, implemented a prototype system, and extensively evaluated its performance over a Taobao dataset. Evaluation results have validated practical feasibility.

APPENDIX

A. Fine-Tuning the Privacy Level of Secure Submodel Updates Aggregation

We introduce how to enable each client to fine-tune the privacy level of secure submodel updates aggregation by analyzing the impacts of the probability parameters p_5 and p_6 on the probabilities of Event 1 and Event 2, i.e., $p_7 = p_5(1-p_5)^{n_{j,1}-1}(1-p_6)^{n_{j,0}}$ and $p_8 = (1-p_5)^{n_{j,1}}(1-(1-p_6)^{n_{j,0}})$, respectively. Without loss of generality, we consider three different policies: (1) Fixing p_5 and adjusting p_6 ; (2) fixing p_6 and adjusting p_5 ; and (3) fixing $p_5 + p_6 = 1$. In particular, as shown in Table V, we mainly adopted the first and the third policies in our evaluation.

We first analyze p_7 as follows: (1) If p_5 is fixed, p_7 decreases as p_6 increases; (2) If p_6 is fixed, the monotonicity is nontrivial. We need to compute the derivative of p_7 with respect to p_5 , i.e.,

$$\frac{\mathrm{d}p_7}{\mathrm{d}p_5} = (1 - p_5 n_{j,1}) (1 - p_5)^{n_{j,1} - 2} (1 - p_6)^{n_{j,0}}.$$

Thus, if $p_5 < 1/n_{j,1}$, p_7 increases as p_5 increases; otherwise, p_7 decreases as p_5 increases; (3) If $p_5 + p_6 = 1$, we first simplify p_7 into $p_5^{n_{j,0}+1}(1-p_5)^{n_{j,1}-1}$. Then, we compute the derivative of p_7 with respect to p_5 as

$$\frac{\mathrm{d}p_{7}}{\mathrm{d}p_{5}} = \left(n_{j,0} + 1 - \left(n_{j,0} + n_{j,1}\right)p_{5}\right)p_{5}^{n_{j,0}}\left(1 - p_{5}\right)^{n_{j,1} - 2}.$$

Hence, if $p_5 < (n_{j,0} + 1)/(n_{j,0} + n_{j,1})$, p_7 increases as p_5 increases; otherwise, p_7 decreases as p_5 increases.

We next analyze p_8 as follows: (1) If p_5 is fixed, p_8 increases as p_6 increases; (2) If p_6 is fixed, p_8 decreases as p_5 increases; (3) If $p_5 + p_6 = 1$, we simplify p_8 into $(1 - p_5)^{n_{j,1}} (1 - p_5^{n_{j,0}})$. Thus, p_8 decreases as p_5 increases.

Finally, we can verify the above deductions by checking the differences among CPPs and the changes in the resulting privacy levels of secure submodel updates aggregation in Table V, where the number of chosen clients is 100 (i.e., $n_{j,0}+n_{j,1}=100$), and the number of chosen clients who have an arbitrary goods ID from the union is 1.17 in average (i.e., $n_{j,0}=98.83, n_{j,1}=1.17$).

B. \(\gamma\)-Level Stochastic Quantization and Weighted Average

We first briefly review the application of γ -level stochastic quantization in secure federated submodel learning. Considering that the compression, federated averaging (or mathematically, weighted average/mean), and decompression of submodel updates are element-wise, we thus just focus on the update of one parameter. Here, we consider that the parameter

update and the count number from client $i \in \mathcal{C}$ are $\Delta w^{(i)} \in \mathbb{R}$ and $v^{(i)} \geq 0 \in \mathbb{Z}$, respectively, and the cloud server wants to compute $\sum_{i \in \mathcal{C}} v^{(i)} \Delta w^{(i)} / \sum_{i \in \mathcal{C}} v^{(i)}$. In addition, we assume that $\forall i \in \mathcal{C}, \Delta w_{\min} \leq \Delta w^{(i)} \leq \Delta w_{\max}$. Moreover, for γ -level quantization, the interval from Δw_{\min} to Δw_{\max} should be equally divided into $\gamma - 1$ segments, where the length of each segment is $\Delta w_{\min} = (\Delta w_{\max} - \Delta w_{\min})/(\gamma - 1)$. In fact, to compress $\Delta w^{(i)}$, client i needs to find the segment that contains $\Delta w^{(i)}$, and then takes either the starting index or the ending index of the segment with probability (w.p.) inversely proportional to the distance between $\Delta w^{(i)}$ and the starting or ending point. More specifically, $\Delta w^{(i)}$ is mapped into $z^{(i)} \in \{0,1,\ldots,\gamma-1\}$, where

$$z^{(i)} = \begin{cases} \left\lfloor \frac{\Delta w^{(i)} - \Delta w_{\min}}{\Delta w_{\text{unit}}} \right\rfloor \\ \text{w.p. } \left\lceil \frac{\Delta w^{(i)} - \Delta w_{\min}}{\Delta w_{\text{unit}}} \right\rceil - \frac{\Delta w^{(i)} - \Delta w_{\min}}{\Delta w_{\text{unit}}}, \\ \left\lceil \frac{\Delta w^{(i)} - \Delta w_{\min}}{\Delta w_{\text{unit}}} \right\rceil \text{ otherwise.} \end{cases}$$

On one hand, we can compute that the expectation of $z^{(i)}$ is $\frac{\Delta w^{(i)} - \Delta w_{\min}}{\Delta w_{\text{unit}}}$, which is denoted as $z^{(i)*}$ and will be consistently used in our following derivations. On the other hand, we can decompress $z^{(i)}$ and get the recovered parameter as $z^{(i)} \Delta w_{\text{unit}} + \Delta w_{\min}$, the expectation of which is $\Delta w^{(i)}$.

We next demonstrate that the weighted averaging operation in the compression space will not introduce any bias/error. On each client's side, she compresses her parameter update into:

$$\forall i \in \mathcal{C}, z^{(i)*} = \frac{\Delta w^{(i)} - \Delta w_{\min}}{\Delta w_{\text{unit}}}.$$

Then, she further weights $z^{(i)*}$ through multiplying by $v^{(i)}$ (i.e., $v^{(i)}z^{(i)*}$), and uploads the materials for secure aggregation. On the cloud server's side, she divides the aggregate result $\sum_{i\in\mathcal{C}}v^{(i)}z^{(i)*}$ by $\sum_{i\in\mathcal{C}}v^{(i)}$, gets $\sum_{i\in\mathcal{C}}v^{(i)}z^{(i)*}/\sum_{i\in\mathcal{C}}v^{(i)}$, and finally performs decompression as

$$\begin{split} &\frac{\sum_{i \in \mathcal{C}} v^{(i)} z^{(i)*}}{\sum_{i \in \mathcal{C}} v^{(i)}} \Delta w_{\text{unit}} + \Delta w_{\text{min}} \\ &= \frac{\sum_{i \in \mathcal{C}} v^{(i)} \frac{\Delta w^{(i)} - \Delta w_{\text{min}}}{\Delta w_{\text{unit}}}}{\sum_{i \in \mathcal{C}} v^{(i)}} \Delta w_{\text{unit}} + \Delta w_{\text{min}} \\ &= \frac{\sum_{i \in \mathcal{C}} v^{(i)} \left(\Delta w^{(i)} - \Delta w_{\text{min}}\right)}{\sum_{i \in \mathcal{C}} v^{(i)}} + \Delta w_{\text{min}} \\ &= \frac{\sum_{i \in \mathcal{C}} v^{(i)} \Delta w^{(i)}}{\sum_{i \in \mathcal{C}} v^{(i)}} - \frac{\sum_{i \in \mathcal{C}} v^{(i)} \Delta w_{\text{min}}}{\sum_{i \in \mathcal{C}} v^{(i)}} + \Delta w_{\text{min}} \\ &= \frac{\sum_{i \in \mathcal{C}} v^{(i)} \Delta w^{(i)}}{\sum_{i \in \mathcal{C}} v^{(i)}} - \Delta w_{\text{min}} + \Delta w_{\text{min}} \\ &= \frac{\sum_{i \in \mathcal{C}} v^{(i)} \Delta w^{(i)}}{\sum_{i \in \mathcal{C}} v^{(i)}}, \end{split}$$

which is the same as the desired outcome in expectation.

In a nutshell, we can conclude that secure federated sub-model learning does not introduce any error/bias under the γ -level stochastic quantization mechanism.

C. Divergence of Conventional Federated Learning

The behavior of divergence has ever been observed due to large local epoch numbers in [9], which initially proposed the federated averaging algorithm. However, the local epoch number in our evaluation is set to the minimum 1, which cannot account for divergence here. As illustrated in Section VI-A, the major reason is that conventional federated learning inaccurately counts in the weights (i.e., the training set sizes) of those clients who contribute zero/no updates for some network parameters (e.g., some embedding vectors in DIN) when computing the weighted average of the full model updates. We examine an embedding vector for the goods with ID 1 for example. We consider that 100 clients are chosen, and assume that the sizes of their training sets are all 300s. Additionally, only 10 samples in client 1's training set involves goods 1 while the samples of the other 99 clients do not, which implies that only client 1 updates the embedding vector for goods 1 while the others do not. After each client trains locally for one epoch and uploads update, under the conventional federated learning framework, the cloud server will update the global model by adding $300/(300 \times 100) = 1\%$ of the client 1's update of the embedding vector, where the total weights count in the training set sizes of the other 99 clients who contribute zero updates. In contrast, under our federated submodel learning framework, the cloud server will add 10/10 = 100% of the client 1's update of the embedding vector to the global model, where the weight of client 1 along with the total weights only count in the size of involved samples. By comparing with centralized training using the same hyperparameters over all 100 clients' data in sequence for one epoch, which will update the global model by adding 100% rather than 1% of the client 1's update of the embedding vector, we can find that federated averaging of the full model updates may cause some network layers, which do not involve each client's whole training set, to be trained in an inaccurate way. This further indicates that we should leverage the fine-grained involved training set sizes (e.g., at the level of individual embedding vector here) as weights like in our federated submodel learning, rather than using the whole training set size at the level of the full model as a unified weight like in conventional federated learning. Of course, the roughness and inaccuracy of conventional federated learning are completely exposed in our e-commerce context, mainly due to the high heterogeneity of user data. In particular, the full size of goods and category IDs are huge, and different Taobao users tend to have highly differentiated or even mutually exclusive sets of goods and category IDs, thus involving and updating different rows of the embedding matrix. Nevertheless, in some other contexts, where user data and their truly required submodels are not heterogenous enough, these shortcomings may not appear. For example, in the natural language context considered by Google, clients use a small vocabulary of size 10,000 for local training in Gboard, which is similar to the full set of goods and category IDs but with a much smaller scale. Therefore, for the embedding vector of a certain word, federated averaging results of its updates using the size of a client's whole text samples and the size of the samples involving this word may not differ too much.

D. Leakage of A Client's Real Index Set in Multiple Communication Rounds and Countermeasures

We now introduce a privacy leakage that if a client participates in multiple communication rounds of federated (submodel) learning even with the strongest security and privacy guarantees, the cloud server, as an adversary, may reveal a client's real indices. We focus on client i with her real index set $S^{(i)}$ and assume that client i participates in two communication rounds, where the unions of chosen clients' real index sets are denoted as \mathcal{U}_1 and \mathcal{U}_2 , respectively. We note that \mathcal{U}_1 and \mathcal{U}_2 are both revealed to the cloud server, no matter whether in our secure federated submodel learning or in secure federated learning. First, the leakages of \mathcal{U}_1 and \mathcal{U}_2 in secure federated submodel learning are trivial. Regarding secure federated learning, the cloud server can still learn \mathcal{U}_1 and \mathcal{U}_2 from two aggregate full model updates by filtering out those indices with zero updates. Please see Fig. 2(a) for an intuition, and refer to the proof of Theorem 3 for formal reasoning. Now, the cloud server computes the intersection of \mathcal{U}_1 and \mathcal{U}_2 , which must contain $\mathcal{S}^{(i)}$, i.e., $\mathcal{S}^{(i)} \subset \mathcal{U}_1 \cap \mathcal{U}_2$. Here, we can further derive that: (1) if client i participates in more communication rounds, due to the properties of the intersection operation, the cloud server can narrow down the scope of $S^{(i)}$; (2) if the real index sets of different clients are mutually exclusive, and the cloud server selects totally different sets of clients (excluding client i) in two communication rounds, then $S^{(i)} = U_1 \cap U_2$.

To mitigate the above leakage, we give the following four kinds of countermeasures. (1) First, we adopt the concept of "period" introduced in Section IV-B2, where a period comprises multiple communication rounds, and within a certain period, each client just uses his historical data in the previous one period to participate in federated (submodel) learning. To avoid the above leakage, we here enforce that each client is only chosen to join in one communication round in one period. Under such a circumstance, the cloud server only obtains \mathcal{U}_1 and knows $\mathcal{S}^{(i)} \subset \mathcal{U}_1$, but cannot further get \mathcal{U}_2 and more unions to narrow down the scope of $S^{(i)}$. (2) Second, we adopt the concept of "group", where the clients in a group participate in federated (submodel) learning together. For example, if client i joins in a group of 50 clients, then the cloud server may only learn the union of these 50 clients' real index sets, even if this group participates in an infinite number of communication rounds. (3) Third is resorting to anonymization. For example, a client can use pseudo identities to participate in different communication rounds, which does not affect the execution of secure federated (submodel) learning. However, such an anonymization process breaks the linkability and thus disables the intersection operation between \mathcal{U}_1 and \mathcal{U}_2 because the cloud server does not know whether a specific client participates in both communication rounds or not. (4) Last is to let each client replace her real index set with a perturbed index set to participate in federated (submodel) learning. More specifically, just as index set perturbation in Algorithm 2, each client only keeps part of her real indices, randomly adds some other indices from the full index set or other clients' indices in her previously involved rounds, and thus generates a perturbed index set locally. Even though a client is chosen for an infinite number of communication rounds, the cloud server may only reveal her perturbed index set.

ACKNOWLEDGMENT

We would like to sincerely thank some members of Advanced Network Laboratory (ANL) in Shanghai Jiao Tong University, including Renjie Gu, Hongtao Lv, Hejun Xiao, and Zhenzhe Zheng, as well as many colleagues in Alibaba, for meaningful discussions and great engineering support.

REFERENCES

- L. Melis, C. Song, E. D. Cristofaro, and V. Shmatikov, "Exploiting unintended feature leakage in collaborative learning," in *Proc. of S&P*, 2019, pp. 497–512.
- [2] A. Hard, K. Rao, R. Mathews, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, "Federated learning for mobile keyboard prediction," arXiv: 1811.03604, 2018, http://arxiv.org/abs/1811.03604.
- [3] H. Cheng, L. Koc, J. Harmsen, T. Shaked, T. Chandra, H. Aradhye, G. Anderson, G. Corrado, W. Chai, M. Ispir, R. Anil, Z. Haque, L. Hong, V. Jain, X. Liu, and H. Shah, "Wide & deep learning for recommender systems," in *Proc. of the 1st Workshop on Deep Learning for Recommender Systems (DLRS)*, 2016, pp. 7–10.
- [4] G. Zhou, X. Zhu, C. Song, Y. Fan, H. Zhu, X. Ma, Y. Yan, J. Jin, H. Li, and K. Gai, "Deep interest network for click-through rate prediction," in *Proc. of KDD*, 2018, pp. 1059–1068.
- [5] G. Zhou, N. Mou, Y. Fan, Q. Pi, W. Bian, C. Zhou, X. Zhu, and K. Gai, "Deep interest evolution network for click-through rate prediction," in *Proc. of AAAI*, 2019.
- [6] Q. Chen, H. Zhao, W. Li, P. Huang, and W. Ou, "Behavior sequence transformer for e-commerce recommendation in alibaba," arXiv: 1905.06874, 2019, http://arxiv.org/abs/1905.06874.
- [7] European Parliament and Council of the European Union, "The General Data Protection Regulation (EU) 2016/679 (GDPR)," https://eurlex.europa.eu/eli/reg/2016/679/oj, Apr. 2016, took effect from May 25, 2018.
- [8] A. A. Ginart, M. Guan, G. Valiant, and J. Zou, "Making ai forget you: Data deletion in machine learning," in *Proc. of NeurIPS*, 2019.
- [9] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Proc. of AISTATS*, 2017, pp. 1273–1282.
- [10] J. Wang, P. Huang, H. Zhao, Z. Zhang, B. Zhao, and D. L. Lee, "Billion-scale commodity embedding for e-commerce recommendation in alibaba," in *Proc. of KDD*, 2018, pp. 839–848.
- [11] M. Chen, R. Mathews, T. Ouyang, and F. Beaufays, "Federated learning of out-of-vocabulary words," arXiv: 1903.10635, 2019, http://arxiv.org/ abs/1903.10635.
- [12] S. Ramaswamy, R. Mathews, K. Rao, and F. Beaufays, "Federated learning for emoji prediction in a mobile keyboard," arXiv: 1906.04329, 2019, http://arxiv.org/abs/1906.04329.
- [13] S. Angel, H. Chen, K. Laine, and S. Setty, "PIR with compressed queries and amortized query processing," in *Proc. of S&P*, 2018, pp. 962–979.
- [14] H. Chen, Z. Huang, K. Laine, and P. Rindal, "Labeled PSI from fully homomorphic encryption with malicious security," in *Proc. of CCS*, 2018, pp. 1223–1237.
- [15] S. Patel, G. Persiano, and K. Yeo, "Private stateful information retrieval," in *Proc. of CCS*, 2018, pp. 1002–1019.
- [16] S. Angel and S. Setty, "Unobservable communication over fully untrusted infrastructure," in *Proc. of OSDI*, 2016, pp. 551–569.
- [17] M. O. Rabin, "How to exchange secrets with oblivious transfer," IACR Cryptology ePrint Archive, Report 2005/187, 1981, https://eprint.iacr.org/2005/187.

- [18] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for privacy-preserving machine learning," in *Proc. of CCS*, 2017, pp. 1175–1191.
- [19] P. Paillier, "Public-key cryptosystems based on composite degree residuosity classes," in *Proc. of EUROCRYPT*, 1999, pp. 223–238.
- [20] D. Boneh, E. Goh, and K. Nissim, "Evaluating 2-dnf formulas on ciphertexts," in *Proc. of TCC*, 2005, pp. 325–341.
- [21] D. M. Freeman, "Converting pairing-based cryptosystems from composite-order groups to prime-order groups," in *Proc. of EURO-CRYPT*, 2010, pp. 44–61.
- [22] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," arXiv: 1908.07873, 2019, http://arxiv.org/abs/1908.07873.
- [23] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," ACM Transactions on Intelligent Systems and Technology (TIST), vol. 10, no. 2, pp. 12:1–12:19, 2019.
- [24] H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang, "Learning differentially private recurrent language models," in *Proc. of ICLR*, 2018.
- [25] M. Abadi, A. Chu, I. J. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep learning with differential privacy," in *Proc. of CCS*, 2016, pp. 308–318.
- [26] V. Feldman, I. Mironov, K. Talwar, and A. Thakurta, "Privacy amplification by iteration," in *Proc. of FOCS*, 2018, pp. 521–532.
- [27] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to backdoor federated learning," arXiv: 1807.00459, 2018, http://arxiv.org/ abs/1807.00459.
- [28] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," arXiv: 1610.05492, 2016, http://arxiv.org/abs/1610.05492.
- [29] A. T. Suresh, F. X. Yu, S. Kumar, and H. B. McMahan, "Distributed mean estimation with limited communication," in *Proc. of ICML*, 2017, pp. 3329–3337.
- [30] N. Agarwal, A. T. Suresh, F. Yu, S. Kumar, and H. B. McMahan, "cpsgd: Communication-efficient and differentially-private distributed SGD," in *Proc. of NeurIPS*, 2018, pp. 7575–7586.
- [31] S. Caldas, J. Konečný, H. B. McMahan, and A. Talwalkar, "Expanding the reach of federated learning by reducing client resource requirements," arXiv: 1812.07210, 2018, http://arxiv.org/abs/1812.07210.
- [32] H. Yu, S. Yang, and S. Zhu, "Parallel restarted SGD with faster convergence and less communication: Demystifying why model averaging works for deep learning," in *Proc. of AAAI*, 2019, pp. 5693–5700.
- [33] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of FedAvg on non-iid data," arXiv: 1907.02189, 2019, http://arxiv.org/ abs/1907.02189.
- [34] H. Yu, R. Jin, and S. Yang, "On the linear speedup analysis of communication efficient momentum SGD for distributed non-convex optimization," in *Proc. of ICML*, 2019, pp. 7184–7193.
- [35] V. Smith, C.-K. Chiang, M. Sanjabi, and A. Talwalkar, "Federated multitask learning," in *Proc. of NeurIPS*, 2017, pp. 4424–4434.
- [36] F. Chen, Z. Dong, Z. Li, and X. He, "Federated meta-learning for recommendation," arXiv: 1802.07876, 2018, http://arxiv.org/abs/1802.07876.
- [37] M. Mohri, G. Sivek, and A. T. Suresh, "Agnostic federated learning," in *Proc. of ICML*, 2019, pp. 4615–4625.
- [38] H. Eichner, T. Koren, H. B. McMahan, N. Srebro, and K. Talwar, "Semi-cyclic stochastic gradient descent," in *Proc. of ICML*, 2019, pp. 1764–1773.
- [39] 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: 1812.02903, 2018, http://arxiv.org/abs/1812.02903.
- [40] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingerman, V. Ivanov, C. Kiddon, J. Konečný, S. Mazzocchi, H. B. McMahan, T. V. Overveldt, D. Petrou, D. Ramage, and J. Roselander, "Towards federated learning at scale: System design," in *Proc. of SysML*, 2019.
- [41] "Tensorflow federated: Machine learning on decentralized data," https://www.tensorflow.org/federated, 2019.
- [42] S. Caldas, P. Wu, T. Li, J. Konečný, H. B. McMahan, V. Smith, and A. Talwalkar, "LEAF: A benchmark for federated settings," arXiv: 1812.01097, 2018, http://arxiv.org/abs/1812.01097.
- [43] M. Fredrikson, E. Lantz, S. Jha, S. Lin, D. Page, and T. Ristenpart, "Privacy in pharmacogenetics: An end-to-end case study of personalized warfarin dosing," in *Proc. of USENIX Security*, 2014, pp. 17–32.

- [44] M. Fredrikson, S. Jha, and T. Ristenpart, "Model inversion attacks that exploit confidence information and basic countermeasures," in *Proc. of CCS*, 2015, pp. 1322–1333.
- [45] G. Ateniese, L. V. Mancini, A. Spognardi, A. Villani, D. Vitali, and G. Felici, "Hacking smart machines with smarter ones: How to extract meaningful data from machine learning classifiers," *IJSN*, vol. 10, no. 3, pp. 137–150, 2015.
- [46] L. Zhu, Z. Liu, and S. Han, "Deep leakage from gradients," in *Proc. of NeurIPS*, 2019.
- [47] Ú. Erlingsson, V. Pihur, and A. Korolova, "RAPPOR: randomized aggregatable privacy-preserving ordinal response," in *Proc. of CCS*, 2014, pp. 1054–1067.
- [48] G. Fanti, V. Pihur, and Ú. Erlingsson, "Building a RAPPOR with the unknown: Privacy-preserving learning of associations and data dictionaries," *Proceedings on Privacy Enhancing Technologies (PoPETs)*, vol. 2016, no. 3, pp. 41–61, 2016.
- [49] Apple's Differential Privacy Team, "Learning with privacy at scale," Apple Machine Learning Journal, vol. 1, no. 8, 2017.
- [50] B. Ding, J. Kulkarni, and S. Yekhanin, "Collecting telemetry data privately," in *Proc. of NeurIPS*, 2017, pp. 3574–3583.
- [51] S. P. Kasiviswanathan, H. K. Lee, K. Nissim, S. Raskhodnikova, and A. Smith, "What can we learn privately?" in *Proc. of FOCS*, 2008, pp. 531–540.
- [52] J. C. Duchi, M. I. Jordan, and M. J. Wainwright, "Local privacy and statistical minimax rates," in *Proc. of FOCS*, 2013, pp. 429–438.
- [53] R. Bassily and A. Smith, "Local, private, efficient protocols for succinct histograms," in *Proc. of STOC*, 2015, pp. 127–135.
- [54] Z. Qin, Y. Yang, T. Yu, I. Khalil, X. Xiao, and K. Ren, "Heavy hitter estimation over set-valued data with local differential privacy," in *Proc.* of CCS, 2016, pp. 192–203.
- [55] R. Bassily, K. Nissim, U. Stemmer, and A. Thakurta, "Practical locally private heavy hitters," in *Proc. of NeurIPS*, 2017, pp. 2288–2296.
- [56] T. Wang, J. Blocki, N. Li, and S. Jha, "Locally differentially private protocols for frequency estimation," in *Proc. of USENIX Security*, 2017, pp. 729–745.
- [57] T. Wang, N. Li, and S. Jha, "Locally differentially private frequent itemset mining," in *Proc. of S&P*, 2018, pp. 127–143.
- [58] Z. Zhang, T. Wang, N. Li, S. He, and J. Chen, "CALM: consistent adaptive local marginal for marginal release under local differential privacy," in *Proc. of CCS*, 2018, pp. 212–229.
- [59] Ú. Erlingsson, V. Feldman, I. Mironov, A. Raghunathan, K. Talwar, and A. Thakurta, "Amplification by shuffling: From local to central differential privacy via anonymity," in *Proc. of SODA*, 2019, pp. 2468–2479.
- [60] T. Wang, B. Ding, J. Zhou, C. Hong, Z. Huang, N. Li, and S. Jha, "Answering multi-dimensional analytical queries under local differential privacy," in *Proc. of SIGMOD*, 2019, pp. 159–176.
- [61] T. Wang, M. Lopuhaä-Zwakenberg, Z. Li, B. Skoric, and N. Li, "Consistent and accurate frequency oracles under local differential privacy," arXiv: 1905.08320, 2019, http://arxiv.org/abs/1905.08320.
- [62] S. L. Warner, "Randomized response: A survey technique for eliminating evasive answer bias," *Journal of the American Statistical Association*, vol. 60, no. 309, pp. 63–69, 1965.
- [63] C. Dwork and A. Roth, "The algorithmic foundations of differential privacy," Foundations and Trends in Theoretical Computer Science, vol. 9, no. 3-4, pp. 211–407, 2014.
- [64] O. Goldreich, S. Micali, and A. Wigderson, "How to play any mental game or A completeness theorem for protocols with honest majority," in *Proc. of STOC*, 1987, pp. 218–229.
- [65] M. Ben-Or, S. Goldwasser, and A. Wigderson, "Completeness theorems for non-cryptographic fault-tolerant distributed computation (extended abstract)," in *Proc. of STOC*, 1988, pp. 1–10.
- [66] M. Burkhart, M. Strasser, D. Many, and X. Dimitropoulos, "SEPIA: privacy-preserving aggregation of multi-domain network events and statistics," in *Proc. of USENIX Security*, 2010, pp. 223–240.
- [67] H. Corrigan-Gibbs and D. Boneh, "Prio: Private, robust, and scalable computation of aggregate statistics," in *Proc. of NSDI*, 2017, pp. 259– 282.
- [68] B. H. Bloom, "Space/time trade-offs in hash coding with allowable errors," Communications of the ACM, vol. 13, no. 7, pp. 422–426, 1970.
- [69] D. Starobinski, A. Trachtenberg, and S. Agarwal, "Efficient PDA synchronization," *IEEE Transactions on Mobile Computing*, vol. 2, no. 1, pp. 40–51, 2003.

- [70] A. Broder and M. Mitzenmacher, "Network applications of bloom filters: A survey," *Internet Mathematics*, vol. 1, no. 4, pp. 485–509, 2004.
- [71] L. Fan, P. Cao, J. Almeida, and A. Z. Broder, "Summary cache: A scalable wide-area web cache sharing protocol," in *Proc. of SIGCOMM*, 1998, pp. 254–265.
- [72] M. J. Freedman, K. Nissim, and B. Pinkas, "Efficient private matching and set intersection," in *Proc. of EUROCRYPT*, 2004, pp. 1–19.
- [73] C. Dong, L. Chen, and Z. Wen, "When private set intersection meets big data: an efficient and scalable protocol," in *Proc. of CCS*, 2013, pp. 789–800.
- [74] P. Rindal and M. Rosulek, "Improved private set intersection against malicious adversaries," in *Proc. of EUROCRYPT*, 2017, pp. 235–259.
- [75] H. Chen, K. Laine, and P. Rindal, "Fast private set intersection from homomorphic encryption," in *Proc. of CCS*, 2017, pp. 1243–1255.
- [76] V. Kolesnikov, N. Matania, B. Pinkas, M. Rosulek, and N. Trieu, "Practical multi-party private set intersection from symmetric-key techniques," in *Proc. of CCS*, 2017, pp. 1257–1272.
- [77] B. Pinkas, T. Schneider, and M. Zohner, "Scalable private set intersection based on OT extension," ACM Transactions on Privacy and Security, vol. 21, no. 2, pp. 7:1–7:35, 2018.
- [78] E. D. Cristofaro, P. Gasti, and G. Tsudik, "Fast and private computation of cardinality of set intersection and union," in *Proc. of CANS*, 2012, pp. 218–231.
- [79] R. Egert, M. Fischlin, D. Gens, S. Jacob, M. Senker, and J. Till-manns, "Privately computing set-union and set-intersection cardinality via bloom filters," in *Proc. of ACISP*, 2015, pp. 413–430.
- [80] E. Fenske, A. Mani, A. Johnson, and M. Sherr, "Distributed measurement with private set-union cardinality," in *Proc. of CCS*, 2017, pp. 2295–2312.
- [81] G. Mezzour, A. Perrig, V. Gligor, and P. Papadimitratos, "Privacy-preserving relationship path discovery in social networks," in *Proc. of CANS*, 2009, pp. 189–208.
- [82] M. Nagy, E. D. Cristofaro, A. Dmitrienko, N. Asokan, and A. Sadeghi, "Do I know you?: efficient and privacy-preserving common friend-finder protocols and applications," in *Proc. of ACSAC*, 2013, pp. 159–168.
- [83] P. Baldi, R. Baronio, E. D. Cristofaro, P. Gasti, and G. Tsudik, "Countering GATTACA: efficient and secure testing of fully-sequenced human genomes," in *Proc. of CCS*, 2011, pp. 691–702.
- [84] A. Narayanan, N. Thiagarajan, M. Lakhani, M. Hamburg, and D. Boneh, "Location privacy via private proximity testing," in *Proc. of NDSS*, 2011.
- [85] B. Pinkas, T. Schneider, G. Segev, and M. Zohner, "Phasing: Private set intersection using permutation-based hashing," in *Proc. of USENIX Security*, 2015, pp. 515–530.
- [86] A. Mani, T. Wilson-Brown, R. Jansen, A. Johnson, and M. Sherr, "Understanding tor usage with privacy-preserving measurement," in *Proc. of IMC*, 2018, pp. 175–187.
- [87] C. Clifton, M. Kantarcioglu, J. Vaidya, X. Lin, and M. Y. Zhu, "Tools for privacy preserving distributed data mining," ACM Sigkdd Explorations Newsletter, vol. 4, no. 2, pp. 28–34, 2002.
- [88] L. Kissner and D. X. Song, "Privacy-preserving set operations," in *Proc. of CRYPTO*, 2005, pp. 241–257.
- [89] K. Frikken, "Privacy-preserving set union," in *Proc. of ACNS*, 2007, pp. 237–252.
- [90] J. H. Seo, J. H. Cheon, and J. Katz, "Constant-round multi-party private set union using reversed laurent series," in *Proc. of PKC*, 2012, pp. 308–412.
- [91] J. Hong, J. W. Kim, J. Kim, K. Park, and J. H. Cheon, "Constant-round privacy preserving multiset union," *Bulletin of the Korean Mathematical Society*, vol. 50, no. 6, pp. 1799–1816, 2013.
- [92] V. Kolesnikov, M. Rosulek, N. Trieu, and X. Wang, "Scalable private set union from symmetric-key techniques," IACR Cryptology ePrint Archive, Report 2019/776, 2019, https://eprint.iacr.org/2019/776.
- [93] D. Many, M. Burkhart, and X. Dimitropoulos, "Fast private set operations with sepia," Communication Systems Group, ETH Zürich, Switzerland, Tech. Rep. TIK-Report No. 345, 2012.
- [94] A. Davidson and C. Cid, "An efficient toolkit for computing private set operations," in *Proc. of ACISP*, 2017, pp. 261–278.
- [95] A. Miyaji and K. Shishido, "Efficient and quasi-accurate multiparty private set union," in *Proc. of SMARTCOMP*, 2018, pp. 309–314.
- [96] T. Kivinen and M. Kojo, "More Modular Exponential (MODP) Diffie-Hellman groups for Internet Key Exchange (IKE)," https://www.ietf.org/ rfc/rfc3526.txt, May 2003.

[97] E. Barker and J. Kelsey, "NIST Special Publication 800-90A Revision 1: Recommendation for Random Number Generation Using Deterministic Random Bit Generators," https://csrc.nist.gov/publications/detail/sp/800-90a/rev-1/final, Jun. 2015.