

# BLINDFL: Vertical Federated Machine Learning without Peeking into Your Data

Fangcheng Fu\*  
Peking University  
ccchengff@pku.edu.cn

Huanran Xue  
Tencent Inc.  
huanranxue@tencent.com

Yong Cheng  
Tencent Inc.  
peterycheng@tencent.com

Yangyu Tao  
Tencent Inc.  
brucetao@tencent.com

Bin Cui\*<sup>†</sup>  
Peking University  
bin.cui@pku.edu.cn

## ABSTRACT

Due to the rising concerns on privacy protection, how to build machine learning (ML) models over different data sources with security guarantees is gaining more popularity. Vertical federated learning (VFL) describes such a case where ML models are built upon the private data of different participated parties that own disjoint features for the same set of instances, which fits many real-world collaborative tasks. Nevertheless, we find that existing solutions for VFL either support limited kinds of input features or suffer from potential data leakage during the federated execution. To this end, this paper aims to investigate both the functionality and security of ML modes in the VFL scenario.

To be specific, we introduce BLINDFL, a novel framework for VFL training and inference. First, to address the functionality of VFL models, we propose the federated source layers to unite the data from different parties. Various kinds of features can be supported efficiently by the federated source layers, including dense, sparse, numerical, and categorical features. Second, we carefully analyze the security during the federated execution and formalize the privacy requirements. Based on the analysis, we devise secure and accurate algorithm protocols, and further prove the security guarantees under the ideal-real simulation paradigm. Extensive experiments show that BLINDFL supports diverse datasets and models efficiently whilst achieves robust privacy guarantees.

## CCS CONCEPTS

• **Computing methodologies** → **Machine learning**; • **Security and privacy** → **Privacy-preserving protocols**.

## KEYWORDS

Vertical Federated Learning, Data Privacy

\*School of Computer Science & Key Lab of High Confidence Software Technologies (MOE), Peking University

<sup>†</sup>Institute of Computational Social Science, Peking University (Qingdao), China

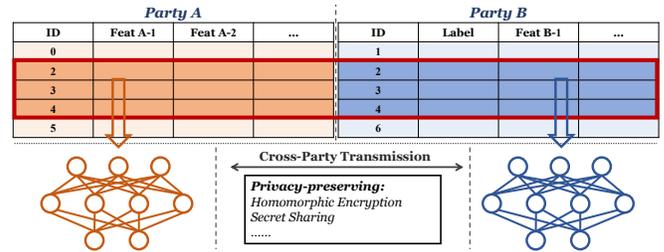
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

SIGMOD '22, June 12–17, 2022, Philadelphia, PA, USA

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9249-5/22/06...\$15.00

<https://doi.org/10.1145/3514221.3526127>



**Figure 1: An illustration of vertical federated learning (VFL). Two parties own disjoint features but have overlapping instances. Privacy-preserving techniques are utilized to protect the data in each party.**

## ACM Reference Format:

Fangcheng Fu, Huanran Xue, Yong Cheng, Yangyu Tao, and Bin Cui. 2022. BLINDFL: Vertical Federated Machine Learning without Peeking into Your Data. In *Proceedings of the 2022 International Conference on Management of Data (SIGMOD '22)*, June 12–17, 2022, Philadelphia, PA, USA. ACM, New York, NY, USA, 21 pages. <https://doi.org/10.1145/3514221.3526127>

## 1 INTRODUCTION

**Background and Motivation.** In recent years, following the explosive surge of data volume and the remarkable success of machine learning (ML) in the whole world, more and more enterprises are thirsty to collect tremendous user data, such as media data, text messages, and daily locations, to build better ML models. Meanwhile, it leads to the increasingly notorious abuse of personal data or even illegal leakage of individual privacy. Thus, the society has raised a growing attention to the protection of data privacy and the supervision on the potential risks to data leakage. Enormous lawful regulations have been established to protect the individual privacy [1, 2, 65]. Consequently, many enterprises are now restricted from collecting a great deal of data for ML tasks.

Owing to such a dilemma of “data shortage”, researchers and data scientists are interested in building ML models with the data of different parties (typically, enterprises or organizations) on the basis of zero data leakage for each individual party. Specifically, a new paradigm called Federated Learning (FL) [33, 34, 37, 45, 69, 74] conveys a possibility to train ML models over multiple data sources that are physically distributed with privacy preservation. In this work, we consider the vertical FL (VFL) scenario. As illustrated in Figure 1, two participated parties own disjoint features (a.k.a.

attributes) but overlap on some instances (a.k.a. samples). The overlapping instances together form a virtually joint dataset, which is vertically partitioned into two parties. *Party B* further holds the ground truth labels. VFL inputs the virtually joint dataset and outputs a federated model that is trained collaboratively.

Our industrial partner runs popular social apps and is able to gather rich user data and precise user profiling. Many collaborators wish to improve the ML ability of their tasks with the help of the data of our industrial partner, such as a Fintech company that hopes to build a more powerful risk model, or an E-commerce company that wishes to make more accurate recommendation<sup>1</sup>. To this end, VFL is a good fit for such kind of cross-enterprise collaboration.

To be formal, the goal of VFL is to unite the features of *Party A* and *Party B*, which are denoted as  $X_A, X_B$ , respectively, to learn a federated model that fits the target labels  $y$  in *Party B*, without any privacy leakage of data in both parties. In this work, we focus on tabular data since it is rare for an image or a sentence to be split and distributed into different parties. Most importantly, the federated model should achieve comparable performance (e.g., accuracy, loss) as the non-federated model trained on collocated datasets.

**Challenges.** We find that the existing VFL solutions can be categorized into two different lines according to how the private features are processed. Nevertheless, the two existing paradigms have complementary strengths and drawbacks.

The first line of works [22, 26, 48, 49] leverages secure multi-party computation (MPC) [76] techniques, such as homomorphic encryption (HE) and secret sharing (SS), to achieve intact privacy guarantees. Typically, they utilize the data outsourcing technique, which is widely used in database services [24, 60], and outsource the datasets to non-colluding servers for ML training or inference. To maintain privacy, all feature values are turned into HE or SS variables when outsourcing so that the servers cannot know the original values. However, such an approach does not fit many real-world datasets. We provide two examples that secretly share features onto two servers in Figure 2. First, for high-dimensional and sparse datasets, the outsourced features become fully dense since the non-zero feature indexes are also private. It prevents us from sparsifying the computation, leading to a performance hazard when the sparsity is high. Second, for categorical features, the embedding lookup operation requires knowing the exact categorical values. However, performing lookup operations on the outsourced values is invalid. As a result, although these MPC-based methods have promising security guarantees, the data outsourcing nature is not suitable for sparse and/or categorical features.

Another line of works follows the split learning paradigm [28, 36, 62, 64, 75, 78, 81], which does not outsource the original datasets so that sparse and categorical features can be handled well. Typically, each party maintains a bottom model in plaintext that extracts the forward activations (a.k.a. hidden representations) using its own private features. The activations of all parties will be exchanged and fed into a top model to make predictions. (See Section 2 for more details.) Nevertheless, the values generated in the bottom models are released in plaintext and would cause data leakage. For

<sup>1</sup>Formally speaking, there can be more than one *Party A*'s. However, most cross-enterprise collaboration follows the two-party setting. Furthermore, our work can be easily generalized to more *Party A*'s. Thus, we only describe the case of only one *Party A* for simplicity, whilst discuss the multi-party setting in our appendix.

**Table 1: Comparison of different works in terms of (i) whether they are suitable for these types of features; (ii) whether they have provable security guarantees under the ideal-real simulation paradigm.**

Paradigm & How Features are Processed	Supported Features			Security Guarantees
	Numerical		Categorical	
	Dense	Sparse		
MPC-based (Data Outsourcing)	✓			✓
Split Learning (Local Bottom Model)	✓	✓	✓	
BLINDFL (this work) (Federated Source Layer)	✓	✓	✓	✓

instance, Figure 2 illustrates a simple example for logistic regression. Although *Party A* cannot get access to *Party B*'s data, it can still infer the labels accurately by analyzing  $X_A W_A$ , because the bottom model (i.e.,  $W_A$ ) is managed by *Party A*. Undoubtedly, such data leakage is forbidden and even illegal in real-world applications. As we will analyze in Section 3, although these works are more flexible to support sparse and/or categorical features, such a design of local bottom models cannot provide provable security guarantees like the MPC-based methods, and thus there exist potential safety hazards.

**Summary of Contributions.** As shown in Table 1, the aforementioned paradigms either support limited feature types or fail to convey promising security guarantees. Motivated by this, we develop a novel VFL framework, namely BLINDFL (pronounced as “blindfold”), to address these challenges. The major contributions of our work are summarized as follows.

**Proposal of BLINDFL.** We propose BLINDFL, a brand new framework for VFL training and inference. BLINDFL keeps the private datasets inside each party without outsourcing and unites the features by an important component called “federated source layer”. By doing so, BLINDFL can support various kinds of input features, including dense, sparse, numerical, and categorical features, whilst achieves security guarantees in the meantime. BLINDFL has been deployed in many productive applications of our industrial partner.

**Analysis of Privacy Requirements.** We anatomize the privacy requirements of VFL training and inference thoroughly. To be specific, we present a comprehensive analysis of the informativeness of all kinds of values generated in the learning process, including forward activations, backward derivatives, model weights, and model gradients, i.e., how it would cause leakage once they are obtained by a party. Upon the analysis, we formulate several privacy requirements — the detailed contents that each party must not get access to during the execution. These privacy requirements shed light on how to judge whether a federated algorithm is secure or not and provide a template to design new algorithm protocols.

**Design of Algorithm Protocols.** Based on the privacy requirements, we devise the federated source layers, the basic building blocks to unite the features from different data sources. The federated source layer leverages the HE and SS techniques to accomplish the aforementioned privacy requirements throughout the algorithm protocols. Two kinds of federated source layers, namely MatMul and Embed-MatMul, are designed for numerical features and categorical features, respectively. We further prove that our algorithm protocols are secure in the presence of semi-honest adversaries

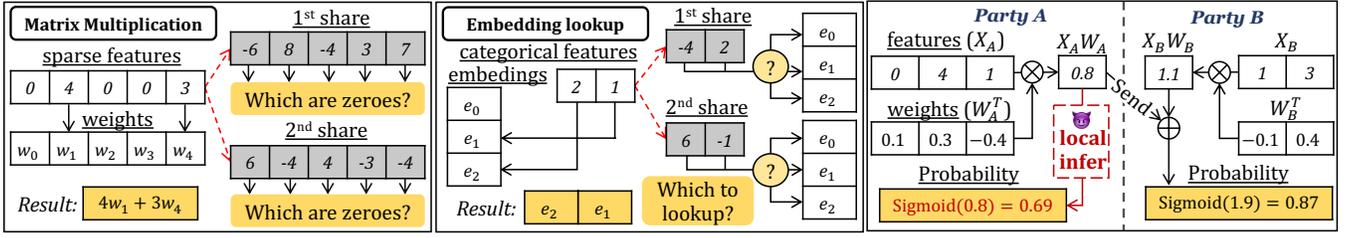


Figure 2: Limitations of existing paradigms. Left-most and middle: Data outsourcing is not suitable for sparse and/or categorical features (illustrated by two examples that secretly share features). Right-most: A simple example of logistic regression following the split learning paradigm. *Party A* can infer the labels accurately since the bottom model (i.e.,  $W_A$ ) is accessible.

that can corrupt up to one party under the ideal-real paradigm. With these two kinds of source layers, BLINDFL is able to support various kinds of features and build diverse VFL models, including generalized linear models (GLMs) and neural networks (NNs).

**Experimental Evaluation.** Comprehensive experiments are conducted to evaluate the effectiveness of BLINDFL. First, empirical results show that BLINDFL is more robust against the semi-honest adversaries and protects the data privacy well. Second, our work outperforms the existing works over 50× in terms of running speed and supports much larger scale of datasets. Third, extensive experiments on a wide range of datasets and models prove that BLINDFL achieves comparable model performance as non-federated learning on collocated datasets, verifying its lossless property.

## 2 PRELIMINARIES

In this section, we briefly introduce the preliminary literature related to our work. For the sake of simplicity, in the rest of the paper, we use the symbol “ $\diamond$ ” to represent an arbitrary party and “ $\oslash$ ” to represent the other party, respectively.

### 2.1 Common ML Ops for Input Features

Given a dataset  $\langle X, y \rangle$  and a loss function  $f$ , where  $X$  is the features and  $y$  is the labels, our goal is to learn a model  $\theta$  that predicts  $\hat{y}$  for  $X$  and minimizes the loss  $f(y, \hat{y})$ . The most prevailing way to solve such a supervised ML problem is to use the first-order gradient optimization, typically, the mini-batch stochastic gradient descent (SGD) and its variants. In each iteration, a mini-batch of instances  $\langle X^{(B)}, y^{(B)} \rangle$  is sampled to calculate the model gradients  $\nabla \theta = \partial f(y^{(B)}, \hat{y}^{(B)}) / \partial \theta$  and the model weights are updated via  $\theta = \theta - \eta \nabla \theta$ , where  $\eta$  is the learning rate (a.k.a. step size). In the rest of the paper, we omit the superscript ( $B$ ) for simplicity.

**Matrix Multiplication.** Matrix multiplication is one of the most common operations in ML. Given the *model weights*  $W \in \mathbb{R}^{IN \times OUT}$ , where  $IN, OUT$  are the input and output dimensionalities, it computes  $Z = XW$  in the forward propagation, which are also known as the *forward activations*. The subsequent modules will take as input  $Z$  to perform their forward propagation routines. During the backward propagation, the *backward derivatives*  $\nabla Z = \partial f / \partial Z$  are propagated from the subsequent modules and the *model gradients* are computed as  $\nabla W = X^T \nabla Z$  according to the chain rule.

**Embedding Lookup.** For categorical features, applying the matrix multiplication is not a common choice in ML since the ordering of categorical values should not matter. In contrast, an embedding

table is usually learned for the categorical inputs<sup>2</sup>. In this work, the embedding table is denoted as  $Q$ . During the forward period, an embedding lookup operation  $E = \text{lkup}(Q, X)$  queries the embedding entries given the categorical indices. The backward propagation computes the model gradients as  $\nabla Q = \text{lkup\_bw}(\nabla E, X)$ . Embedding lookup is usually followed by a matrix multiplication, which outputs  $Z = EW$  given the model weights  $W$ .

### 2.2 Privacy-Preserving Techniques

Security is a widely studied area to protect private data from any adversaries. In this work, we adopt two well-known privacy-preserving techniques — homomorphic encryption and secret sharing — to derive our federated algorithms.

**Homomorphic Encryption.** Homomorphic encryption (HE) [15, 56] describes the cryptographic methods that allow arithmetic computation in the space of ciphers. There are different kinds of HE, such as fully HE, somewhat HE, additive HE, and multiplicative HE [11, 14, 21, 53, 57].

In this work, we focus on additive HE. For instance, Paillier cryptosystem [53] is a well-known additive HE method and has been used in many FL algorithms [10, 19, 25, 71, 75]. Paillier cryptosystem initializes with a key pair  $\langle pk, sk \rangle$ . The public key  $pk$  is utilized in encryption and can be made public to the other party, whilst the secret key (a.k.a. private key)  $sk$  is for decryption and must be kept secret. Given values  $u, v$ , Paillier cryptosystem supports the following types of operations:

- ▷ Encryption:  $\text{Enc}(v, pk) = \llbracket v \rrbracket$ ;
- ▷ Decryption:  $\text{Dec}(\llbracket v \rrbracket, sk) = v$ ;
- ▷ Homomorphic addition:  $\llbracket u \rrbracket + \llbracket v \rrbracket = \llbracket u + v \rrbracket$ ;
- ▷ Scalar addition:  $\llbracket u \rrbracket + v = \llbracket u \rrbracket + \text{Enc}(v, pk) = \llbracket u + v \rrbracket$ ;
- ▷ Scalar multiplication:  $u \llbracket v \rrbracket = \llbracket uv \rrbracket$ .

In the rest of the paper, we assume both parties have generated their own key pairs and exchanged the public keys on initialization. Thus, we omit  $pk$  or  $sk$  and denote  $\llbracket v \rrbracket_\diamond$  as a cipher corresponding to the  $sk$  of *Party*  $\diamond$ .

**Secret Sharing.** As another powerful privacy-preserving technique, secret sharing (SS) [12, 48, 54, 58, 70] breaks a value into pieces of sharing and distributes them to different parties so that none of the parties knows the exact value. For instance, if *Party*  $\diamond$  wishes to secretly share a value  $v$ , it randomly generates its own

<sup>2</sup>Feature engineering techniques such as one-hot encoding or frequency encoding can also be applied to the categorical features, however, we do not discuss them in this work since they are orthogonal.

**Algorithm 1:** The procedure to transform an HE variable  $\llbracket v \rrbracket$  into an SS variable  $\langle \phi, v - \phi \rangle$ .  $v$  is a scalar or a tensor.

```

1 Function HE2SS( $\llbracket v \rrbracket_{\diamond} = None$ ):
2   if  $\llbracket v \rrbracket_{\diamond}$  is None then // Party  $\diamond$  that owns  $sk_{\diamond}$ 
3     Receive  $\llbracket v - \phi \rrbracket_{\diamond}$  and decrypt; return  $v - \phi$ 
4   else // Party  $\diamond$  that does not own  $sk_{\diamond}$ 
5     Randomly generate  $\phi$  with the same shape of  $\llbracket v \rrbracket_{\diamond}$ 
6     Send  $\llbracket v - \phi \rrbracket_{\diamond} = \llbracket v \rrbracket_{\diamond} - \phi$ ; return  $\phi$ 

```

piece of sharing  $v_{\diamond}$  and sends the other piece  $v_{\diamond} = v - v_{\diamond}$  to Party  $\diamond$ . Whenever a party wishes to restore an SS variable, it receives the other piece of sharing and simply adds the two pieces together.

In this work, we focus on two-party additive SS. Given two SS variables  $\langle v_A, v_B \rangle, \langle u_A, u_B \rangle$ , the addition arithmetic can be executed locally as  $\langle u_A + v_A, u_B + v_B \rangle$ . A common way to multiply two SS variables is based on the Beaver triplets. However, each triplet can only be used once and it is time-consuming to generate. We refer interested readers to [5, 55] for more details.

It is worthy to note that HE and SS variables can be transformed easily. For instance, we can transform an HE variable to SS variable via Algorithm 1. Moreover, the operations on HE or SS variables can be easily vectorized. For instance, the matrix multiplication  $\llbracket Z \rrbracket = X \llbracket Y \rrbracket$  is achieved by letting  $\llbracket Z_{ij} \rrbracket = \sum_k X_{ik} \llbracket Y_{kj} \rrbracket$ .

**Security model.** Like many previous works [19, 49, 66, 67, 71], we consider the semi-honest (a.k.a. honest-but-curious) security model, where all parties honestly execute the protocols, whilst the curious parties try to analyze the others' data through any information obtained during the protocols. We assume that a polynomial-time adversary can corrupt one of the two parties during the execution. When analyzing our protocols, we follow the ideal-real paradigm, which is defined as follows.

*Definition 2.1.* ([23, 41]) Let  $\Pi$  be a real-world protocol and  $\mathcal{F}$  be an ideal functionality. We say  $\Pi$  securely realizes  $\mathcal{F}$  if for every adversary  $\mathcal{A}$  attacking the real interaction, there exist a probabilistic polynomial-time simulator  $\mathcal{S}$  attacking the ideal interaction, such that any environment on any input cannot tell apart the real interaction from the ideal interaction, except with negligible probability (in the security parameter  $\kappa$ ).

### 2.3 ML over Different Data Sources

Under the VFL setting, training data come from different sources. The most challenging problem is how to build ML models over different data sources to achieve comparable model performance as non-federated learning on collocated data, whilst guarantee the data privacy of each participated party meanwhile. A series of works are developed to tackle this problem. We divide them into two categories according to how they process the features.

**MPC-based solutions (data outsourcing).** The first kind of works outsources the data of all parties to one or more non-colluding servers for ML training or inference. Before data outsourcing, they transform the feature values into HE or SS variables and leverage the arithmetic properties to perform ML ops. For instance, SecureML [49] secretly shares the features and models onto two servers. Then the matrix multiplication is performed as  $\langle X \rangle \langle W \rangle$ .

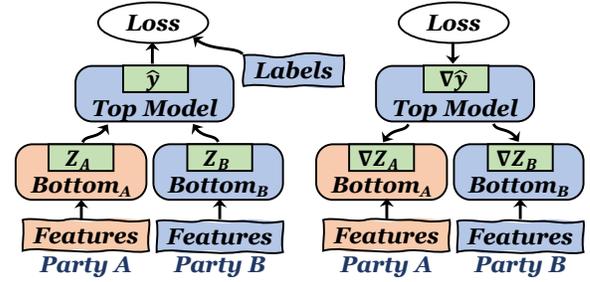


Figure 3: The bottom-to-top architecture of split learning.

Similarly, there are also works that encrypt and send the features to a single server, and utilize the arithmetic properties of HE to carry out the computation [22, 26]. With the help of the HE and SS techniques, none of the servers can reveal the original values, hence these works usually achieve provable security guarantees under the ideal-real paradigm. To apply these works under the VFL setting, we can let each party work as one server and outsource the features among the parties for ML computation.

**Split learning (local bottom models).** The second category follows the split learning paradigm [62, 64]. As shown in Figure 3, each party holds its own data rather than outsourcing. Each party is associated with a bottom model that works as the feature extractor. Party B further manages a top model that makes the final predictions. The forward activations  $Z_A$  and backward derivatives  $\nabla Z_A$  are exchanged between the parties. With such a design of bottom models, split learning is very flexible to support various kinds of features. For instance, to support sparse matrix multiplication, the bottom model of Party  $\diamond$  computes  $X_{\diamond} W_{\diamond}$  using sparsified computation, where  $W_{\diamond}$  is the model weights, and the top model of Party B aggregates  $Z = X_A W_A + X_B W_B$ . To support categorical features, each party manages an embedding table in the bottom model and performs the lookup operation locally. Thus, several VFL algorithms are proposed in such a pattern [28, 36, 75, 78, 81].

## 3 ANATOMY OF EXISTING PARADIGMS

In this section, we anatomize the existing paradigms from the perspective of supported features and security guarantees, respectively.

**The MPC-based solutions.** We first consider the MPC-based solutions, which outsource the features using privacy-preserving techniques to carry out the ML computation.

**Supported features.** As introduced in Section 1, data outsourcing is not suitable for sparse and/or categorical features. On the one hand, sparse datasets will become fully dense, since the outsourced variables should not reveal whether the original values are zeroes or not. On the other hand, after the categorical values are transformed into HE or SS variables, the embedding lookup operation cannot be performed. However, many ML tasks adopt feature engineering techniques such as hashing, quantile binning, and Cartesian product, which make many real-world datasets high-dimensional and sparse. In addition, categorical features are also widely used in practical applications to learn latent representation for better model performance. Therefore, these outsourcing-based solutions are not suitable for many real-world datasets.

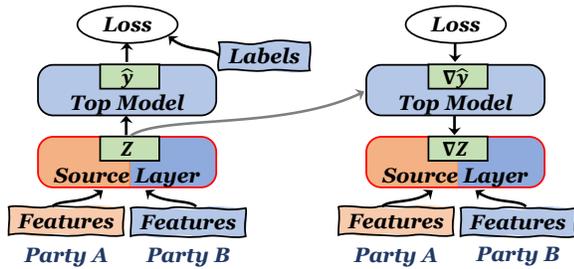


Figure 4: Overview of the forward (left) and backward (right) propagation of BLINDFL.

**Security guarantees.** Despite the aforementioned limitations, the MPC-based solutions usually achieve very promising security guarantees. For instance, SecureML [49] securely realizes the ideal functionality of ML training in the presence of semi-honest adversaries. As a result, methods in this line have a strong privacy preservation ability.

**Split learning.** Next, we consider the methods in the split learning paradigm, where each party is able to process the features locally via a bottom model.

**Supported features.** Compared with the outsourcing-based methods, split learning is more suitable for sparse and/or categorical features since all data are kept local inside each party. During the ML tasks, each party can easily identify whether a feature value is zero or perform the lookup operation given a categorical value.

**Security guarantees.** Since the features are visible to the owner party, it is an intuitive idea to let each party process its own features locally and only share the intermediate results between parties. However, although the raw data are not exposed directly, the intermediate results are also informative and would cause severe data leakage. To be specific, we identify two kinds of leakage:

- **Label leakage.** On the one hand, since *Party A* owns the bottom model, it is able to compute  $Z_A$  individually, which would leak the labels. Take Figure 2 as an example. *Party A* can steal the labels by analyzing  $Z_A = X_A W_A$  with high confidence. Many existing works are vulnerable to this problem [25, 75, 78]. On the other hand, as analyzed by Li et al. [36] and verified by our experiments, the backward derivatives  $\nabla Z_A$  can reveal almost all training labels. The essential idea is that logistic loss produces opposite directions for different labels, so the directions of derivatives reflect the label information inevitably.
- ***Party A*'s feature leakage.** Since the forward activations  $Z_A$  are originated from  $X_A$  and independent from  $X_B$ , *Party B* can analyze a certain level of feature similarities of  $X_A$ . In other words, if the features of two instances  $X_A^{(i)}, X_A^{(j)}$  are very similar, the corresponding activations  $Z_A^{(i)}, Z_B^{(j)}$  would also be very close.

With the hope of avoiding data leakage, some of the existing works try to enhance privacy via MPC techniques such as HE and SS [75, 81]. Nevertheless, these works fail to address all the leakage cases. For instance, Yang et al. [75] propose to encrypt the derivatives in the backward propagation to avoid the label leakage from derivatives. However, they do not consider the label leakage from  $Z_A =$

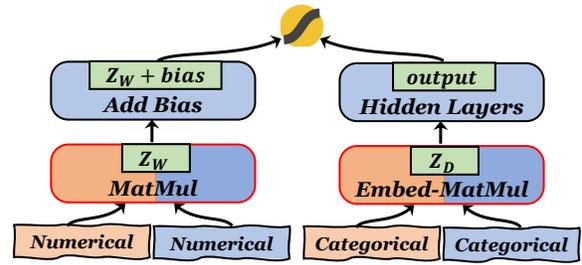


Figure 5: Example of the wide and deep (WDL) model tailored for BLINDFL. There are two source layers – MatMul handles the sparse, numerical features (the wide part) and Embed-MatMul handles the categorical ones (the deep part).

$X_A W_A$  as shown in Figure 2. Thus, *Party A* can still make accurate guesses on the labels.

There are also methods that try to add random noises to perturb these sensitive values. For instance, Li et al. [36] propose to protect the labels by adding noises to  $\nabla Z_A$ . Nevertheless, we do not consider these methods in this work for three reasons. First, there is a tradeoff between protection ability and model quality. When strong privacy is required, the model accuracy will be harmed significantly. Second, these methods cannot provide a formal security guarantee like the MPC-based ones, which is expected in our work. Third, to the best of our knowledge, there is no such method that addresses all these leakage cases together.

In fact, the root cause of such potential data leakage is the design of local bottom models. For instance, as long as *Party A* owns the bottom model, it is able to compute  $Z_A$  anyway, leading to label leakage. To this end, the design of local bottom models cannot provide security guarantees under the ideal-real paradigm. How to build VFL models on non-outsourced data with provable security guarantees needs careful re-investigation.

**Our solution.** In order to support various kinds of features and guarantee privacy preservation in the meantime, this work draws the strengths of both categories. First, our work does not outsource the original datasets. In contrast, we follow the split learning paradigm to keep the private data inside each party. Second, we develop the federated source layers to collaboratively process the features of all parties. Unlike the bottom models, each party is not able to process the features individually in our design. Furthermore, our federated source layers achieve provable security guarantees. Putting them together, our work enjoys the flexibility to support various kinds of features whilst being privacy preserving.

## 4 OVERVIEW AND PRIVACY FORMULATION

### 4.1 Overview

Figure 4 depicts the overview of BLINDFL, which can be decoupled into two parts. First, the federated source layer works as the basic building block that unites the features of both parties. Given the output of the source layer, denoted as  $Z$ , a non-federated sub-module in *Party B*, namely the top model, plays the role of classifier/predictor<sup>3</sup>.

<sup>3</sup>The top model can also be a federated module. However, a non-federated top model is more common in practice to address efficiency. Therefore, this work mainly focuses

During the backward process, *Party B* computes the loss function via the final predictions and ground truth labels, and then back propagates along the top model to obtain the backward derivatives  $\nabla Z$ . Finally, a federated procedure will be executed to update the model weights of the source layer. It is worthy to note that our framework differs from the model architecture of split learning. Our source layer requires all parties to collaboratively execute the learning process and outputs *only* the aggregated results, i.e.,  $Z$ . Whilst in split learning, each party can process the features individually and obtain the unaggregated values from its bottom model.

In practice, since the top model can be an arbitrary sub-module that minimizes the loss between predictions and labels, we concentrate on the source layers. Specifically, we consider two kinds of source layers in this work for different types of input features.

The first kind of source layers, namely MatMul, aims at the numerical feature values and computes a matrix multiplication in the forward propagation, i.e.,  $Z = X_A W_A + X_B W_B$ , where  $W_A, W_B$  are the model weights for *Party A* and *Party B*, respectively. During the backward propagation, a federated procedure is executed to update  $W_\diamond$  by the model gradients  $\nabla W_\diamond = X_\diamond^T \nabla Z$  for each party.

For categorical features, we devise a more complex source layer called Embed-MatMul that fuses the embedding lookup operation and matrix multiplication, i.e.,  $Z = E_A W_A + E_B W_B$ , where  $E_\diamond = \text{lkup}(Q_\diamond, X_\diamond)$  is the lookup operation given the embedding table  $Q_\diamond$ . During the federated backward propagation, the backward derivatives and model gradients are computed as

$$\nabla W_\diamond = E_\diamond^T \nabla Z, \quad \nabla E_\diamond = \nabla Z W_\diamond^T, \quad \nabla Q_\diamond = \text{lkup\_bw}(\nabla E_\diamond, X_\diamond).$$

With these two kinds of source layers, we can derive various VFL models, including generalized linear models and neural networks. For instance, for logistic regression (LR), there is a MatMul source layer with  $OUT = 1$ , whilst the top model adds the bias term and computes the sigmoid function, i.e.,  $\hat{y} = \text{sigmoid}((X_A W_A + X_B W_B) + bias)$ . For another example, as shown in Figure 5, the wide and deep (WDL) model [9] consists of two source layers, one for the sparse, numerical features and the other for the categorical fields.

## 4.2 Privacy to Matter

As discussed in Section 3, although datasets are kept local inside each party, the values generated in the learning process are also informative and would cause data leakage. Thus, before stepping into the design and analysis of the proposed source layers, we would like to discuss the privacy of all kinds of values, including forward activations, backward derivatives, model weights, and model gradients, respectively. In particular, we wish to conclude several guidelines about what contents must not be accessible to a specific party.

**Privacy of forward activations.** Undoubtedly, forward activations have a strong relationship to the ground truth labels since they are learned to fit the labels. For instance, as illustrated in Figure 2, in the MatMul source layer,  $X_A W_A$  can be directly used to make predictions on the labels, so *Party A* should have zero knowledge of them. Consequently, we make the requirement that ① *Party A is not allowed to obtain any forward activations*. However, as we have analyzed in Section 3, since several existing VFL solutions fail to

fulfill this requirement, *Party A* can easily reveal a large fraction of labels, dampening the significance of privacy preservation.

In addition to labels, forward activations inevitably contain sensitive feature information since they are originated from features. As aforementioned, *Party A* is already prohibited from obtaining any forward activations, so we only need to consider whether *Party B* could guess the features of *Party A* via forward activations. We divide forward activations into three kinds: (i) those solely dependent on  $X_A$  (e.g.,  $E_A, X_A W_A$ ), (ii) those dependent on both  $X_A, X_B$  (e.g.,  $Z, \hat{y}$ ), and (iii) those solely dependent on  $X_B$  (e.g.,  $E_B, X_B W_B$ ). For the first kind, as discussed in Section 3, we notice that they would disclose a certain level of feature similarity between different instances. For instance, in the Embed-MatMul layer, once *Party B* obtains  $E_A$ , it realizes whether two instances are equal on some categorical fields by comparing the embedding entries. Therefore, in order to protect the features, ② *we forbid any forward activations that are solely dependent on the features of Party A to be disclosed to Party B*. The second kind aggregates  $X_A, X_B$  for the top model. Since the goal of VFL is to output the inference results to *Party B*, they should be accessible<sup>4</sup>. Furthermore, when we analyze the security guarantees of our work in Section 5.3 and Section 6.3, we will formally prove that these values will not reveal the private data of *Party A*. For the third kind, since they are independent on  $X_A$ , we do not make strict requirements. Instead, we will analyze whether they violate Req ② under specific scenarios. For instance, as we will discuss in Section 5 and Section 6, because *Party B* can derive  $X_A W_A$  (or  $E_A W_A$ ) if it gets access to  $X_B W_B$  (or  $E_B W_B$ ), we will restrict *Party B* from obtaining these forward activations to ensure the privacy of  $X_A$  when designing our source layers.

**Privacy of backward derivatives.** Similarly, in order to avoid the leakage to labels, ③ *Party A is prohibited from accessing any backward derivatives*, e.g.,  $\nabla Z, \nabla E_A$ , since they are originated from the ground truth labels and prediction outputs. In Section 7, we will empirically show that backward derivatives can disclose almost all labels to *Party A*, which makes the privacy preservation in vain. Therefore, Req ③ is vital for designing the VFL algorithms.

Although it is non-trivial to precisely extract the relevance between backward derivatives and features, we can define the security of backward derivatives according to forward activations. In fact, backward derivatives depicts the differences between forward activations and the ideal optimum. Thus, their informativeness regarding features bind together. Motivated as such, we make a symmetric requirement that ④ *if any forward activations are solely dependent on the features of Party A, then the corresponding backward derivatives should also be kept secret from Party B as well*.

**Privacy of model weights and model gradients.** Although seemingly irrelevant to features or labels, the privacy of models is also meaningful. For one thing, given the fact that the values of model weights depict the feature importance upon the learning tasks, it will cause a leakage once a party knows the model information of the other party. For another, model gradients could cause leakage as discussed in recent studies [20, 83]. Although these leakages sound to be task-specific, we still make a tough requirement that ⑤ *the model weights and gradients must be hidden from the*

on non-federated top models whilst discusses how to adapt our federated source layers to federated top models in our appendix.

<sup>4</sup>As described in Section 4.1, our work can adapt to federated top models, where the second kind of forward aggregations are inaccessible. Thus, we focus on the federated source layer in this work and do not restrict *Party B* from them.

**Table 2: Summary of the restrictions for MatMul, i.e., the contents that each party must not get access to.**

	Party A	Party B
Model	$W_A, W_B$	
Forward	$Z, X_A W_A, X_B W_B$	$X_A W_A, X_B W_B$
Backward	$\nabla Z, \nabla W_A, \nabla W_B$	$\nabla W_A$

other party, including the signs and magnitudes of all coordinates, in order to avoid the potential risk of privacy leakage from models. Furthermore, ⑥ we do not allow Party A to obtain the model weights or gradients of its own, even if the sign or magnitude of each coordinate. Otherwise, Party A would infer the labels by analyzing its own feature contributions via the signs or magnitudes.

Obviously, our privacy requirements address the possible leakage in split learning as discussed in Section 3. They provide us with a template to derive the specific restrictions of each party when designing the federated source layers in Section 5 and Section 6.

## 5 MATMUL FEDERATED SOURCE LAYER

Since matrix multiplication is one of the most essential arithmetics in ML, deriving a safe and accurate federated MatMul layer is vital to the VFL paradigm.

### 5.1 Anatomy of Privacy Requirements

To achieve promising privacy guarantees, we follow the analysis in Section 4.2 to derive what kinds of values would cause data leakage and make restrictions on them, i.e., these values should not be accessible to specific parties. We summarize the restrictions in Table 2 and discuss the reasons below:

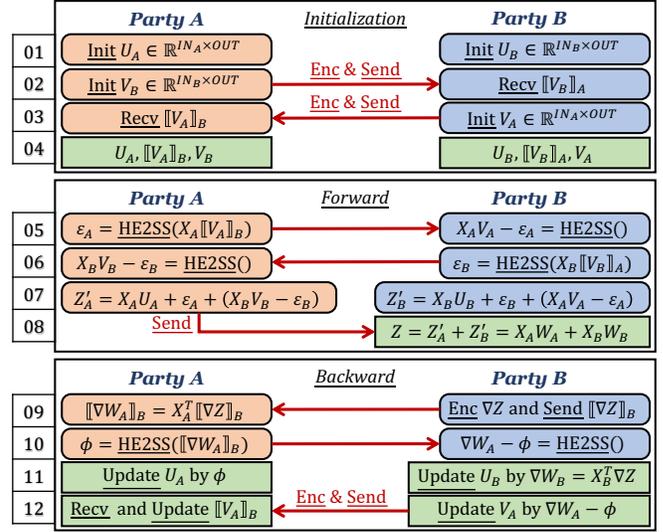
- ▶ First, to avoid label leakage, Party A must not get access to forward activations  $Z$ ,  $X_A W_A$ ,  $X_B W_B$ , backward derivatives  $\nabla Z$ , model weights, and model gradients of its own  $W_A$ ,  $\nabla W_A$ , as analyzed in Req ①③⑥.
- ▶ Second, as analyzed in Req ②, to avoid Party A's feature leakage, Party B must not get access to  $X_A W_A$ , since they are merely the linear transformation of  $X_A$ , which would reveal a certain level of feature similarities. In fact, as  $Z = X_A W_A + X_B W_B$ , this also implies Party B must not get access to  $X_B W_B$  and  $W_B$ .
- ▶ Finally, as analyzed in Req ⑤, we should restrict each party from obtaining the model weights or gradients of the other party  $W_\diamond, \nabla W_\diamond$ .

Revisiting these restrictions, it is worthy to note that the existing split learning based approaches are insecure since the bottom models are typically  $W_A, W_B$ . It gives us a lesson that although features can be maintained locally, we shall not let each party directly learn a bottom model to process the feature individually. In the following subsections, we describe our algorithm protocol for the MatMul source layer and analyze the security guarantees.

### 5.2 Algorithm Protocol

The protocol of our MatMul source layer is presented in Figure 6 and the details are walked through below.

**Initialization.** To begin with, since both parties must not get access to the model weights, we follow the MPC-based methods to



**Figure 6: Our MatMul source layer. All cross-party transmission (red arrows) are protected by HE or SS.**

leverage the SS technique to break model weights onto both parties. To be specific, we secretly share the model weights by  $W_\diamond = U_\diamond + V_\diamond$ , where  $U_\diamond$  and  $V_\diamond$  are held by different parties so that none of them knows what  $W_\diamond$  exactly is. To achieve this goal, on initialization, Party A initializes  $U_A$  for itself and  $V_B$  for Party B (left hand side of Line 1-2). Furthermore, the encrypted version  $\llbracket V_B \rrbracket_A$  is sent to Party B for future use (Line 3). Party B executes a symmetric routine to initialize  $U_B, V_A$ .

**Forward propagation.** Since the model weights are secretly shared, the results are also broken into four parts, i.e.,  $Z = X_A U_A + X_A V_A + X_B U_B + X_B V_B$ . Among them,  $X_\diamond U_\diamond$  can be computed by one party alone, whilst  $X_\diamond V_\diamond$  requires decryption. To be specific, there are three steps in the forward propagation. Since the first two steps are symmetric in both parties, we only describe those in Party A.

- (1) (Line 5-6) Party A computes  $\llbracket X_A V_A \rrbracket_B$  and transforms it into an SS variable  $\langle \epsilon_A, X_A V_A - \epsilon_A \rangle$ ; receives and decrypts the piece of sharing  $X_B V_B - \epsilon_B$  from Party B.
- (2) (Line 7) Party A computes  $Z'_A = X_A U_A + \epsilon_A + (X_B V_B - \epsilon_B)$ .
- (3) (Line 8) Finally, Party B sums into  $Z = Z'_A + Z'_B$ .

It is worthy to note that all the random obfuscation values ( $\epsilon_A, \epsilon_B$ ) are eliminated to achieve the lossless property, i.e.,

$$\begin{aligned} Z &= (X_A U_A + \epsilon_A + (X_B V_B - \epsilon_B)) + (X_B U_B + \epsilon_B + (X_A V_A - \epsilon_A)) \\ &= X_A U_A + X_A V_A + X_B U_B + X_B V_B = X_A W_A + X_B W_B. \end{aligned}$$

**Backward propagation.** For Party B, the model gradients can be computed via  $\nabla W_B = X_B^T \nabla Z$  (right hand side of Line 11). For Party A, we have to leverage the power of HE and SS since  $X_A$  and  $\nabla Z$  are kept secret by different parties. To be specific, we first send Party A the encrypted derivatives  $\llbracket \nabla Z \rrbracket_B$  to compute the encrypted model gradients  $\llbracket \nabla W_A \rrbracket_B$ , which will then be transformed into the SS variable  $\langle \phi, \nabla W_A - \phi \rangle$  (Line 9-10). Furthermore, to prohibit any party from obtaining  $\nabla W_A$  in plaintext, we do not restore the SS variable. Instead, we update the secretly shared model weights



we desiderate a new protocol for our Embed-MatMul source layer with promising security guarantees.

## 6.2 Algorithm Protocol

Considering that both parties must not get access to the embedding tables, we apply SS to the embedding tables similarly, i.e.,  $Q_\diamond = S_\diamond + T_\diamond$ , where  $S_\diamond$  and  $T_\diamond$  are managed by different parties. As a result, the embedding entires are also broken into two parts, i.e.,  $E_\diamond = \text{lkup}(Q_\diamond, X_\diamond) = \text{lkup}(S_\diamond, X_\diamond) + \text{lkup}(T_\diamond, X_\diamond)$ . Therefore, each party cannot obtain  $E_\diamond$  alone since it has zero knowledge on  $T_\diamond$ . We present our Embed-MatMul federated source layer in Figure 7 and describe the sketch of routines below.

**Initialization.** Similar to the MatMul source layer, *Party A* prepares  $S_A, U_A$  for itself and  $T_B, V_B$  for *Party B*. Then the encrypted SS pieces  $\llbracket T_B \rrbracket_A, \llbracket U_A \rrbracket_A, \llbracket V_B \rrbracket_A$  are sent to *Party B* for future use. *Party B* executes a symmetric routine.

**Forward propagation.** We divide the forward process into two stages, for Embed and MatMul, respectively.

The first stage retrieves the secretly shared embedding entries. Since embedding lookup requires the exact values of  $X_\diamond$ , we perform the lookup operation over the encrypted table  $\llbracket T_\diamond \rrbracket_\diamond$  in each party and convert it into an SS variable (Line 5-6). Finally, by combining with the lookup results on the rest piece  $S_\diamond$ , we successfully break the embedding entries in an SS manner, i.e.,  $\langle \psi_\diamond, E_\diamond - \psi_\diamond \rangle$  (Line 7).

The second stage performs two matrix multiplication, i.e.,

$$\text{(Line 8)} \quad Z'_{1,A} + Z'_{1,B} = \psi_A(U_A + V_A) + \psi_B(U_B + V_B),$$

$$\text{(Line 9)} \quad Z'_{2,A} + Z'_{2,B} = (E_B - \psi_B)(V_B + U_B) + (E_A - \psi_A)(V_A + U_A),$$

using the same routine in the forward propagation of our MatMul layer (Figure 6). Finally, the forward outputs can be computed via  $Z = Z'_{1,A} + Z'_{1,B} + Z'_{2,A} + Z'_{2,B} = E_A W_A + E_B W_B$  (Line 10-11).

**Backward propagation.** Next, we describe the backward process, which is also made up of two stages.

The first stage takes in charge of the backward process of MatMul, which updates model weights  $W_\diamond$  by  $\nabla W_\diamond = E_\diamond^T \nabla Z$ . Similar to the backward process of the MatMul source layer, *Party B* encrypts  $\nabla Z$  to protect the labels (Line 12). Upon receiving  $\llbracket \nabla Z \rrbracket_B$ , *Party A* computes  $\llbracket \psi_A^T \nabla Z \rrbracket_B$  via homomorphic arithmetics, which is then transformed into an SS variable, i.e.,  $\langle \phi, \psi_A^T \nabla Z - \phi \rangle$  (Line 13). By computing  $(E_A - \psi_A)^T \nabla Z + (\psi_A^T \nabla Z - \phi) = \nabla W_A - \phi$  in *Party B* (Line 14), we secretly share  $\nabla W_A$  onto both parties, i.e.,  $\langle \phi, \nabla W_A - \phi \rangle$ . A similar routine works for  $\nabla W_B$  as well (Line 15-16). Consequently, none of the parties gets access to  $\nabla W_A$  or  $\nabla W_B$ , which obeys Table 3 so that the model weights can be updated accurately without any information leakage.

The second stage is for the embedding tables. Since the backward operation  $\text{lkup\_bw}(\cdot, \cdot)$  requires to know the exact values of  $X_\diamond$ , we determine to perform it over the encrypted backward derivatives  $\llbracket E_\diamond \rrbracket_\diamond$ . Therefore, both parties first compute  $\llbracket \nabla E_A \rrbracket_B, \llbracket \nabla E_B \rrbracket_A$  via homomorphic arithmetics, respectively (Line 21), and then performs the backward operation  $\llbracket \nabla Q_\diamond \rrbracket_\diamond = \text{lkup\_bw}(\llbracket E_\diamond \rrbracket_\diamond, X_\diamond)$ , which is eventually transformed into SS variables i.e.,  $\langle \rho_\diamond, \nabla Q_\diamond - \rho_\diamond \rangle$  (Line 22-23). Finally, the secretly shared embedding tables are updated in the SS manner (Line 24-26).

```

1 class MatMulOpB(torch.autograd.Function):
2     ... # defines MatMulFw and MatMulBw of Party B
3
4 class MatMulSourceB(FederatedModule):
5     def __init__(self, in_A, in_B, out):
6         super(MatMulSourceB, self).__init__()
7         U_B, enc_V_B, V_A = ... # initialization
8         # federated parameters will be registered to
9         # the federated module automatically
10        self.W_A = FederatedParameter(V_A)
11        self.W_B = FederatedParameter(U_B, enc_V_B)
12
13    def forward(self, X_B):
14        return MatMulOpB.apply(...)
15
16 # user-defined LR model of Party B
17 class LRPartyB(FederatedModule):
18    def __init__(self, in_A, in_B, out):
19        super(LRPartyB, self).__init__()
20        self.source = MatMulSourceB(in_A, in_B, out)
21        self.bias = torch.nn.Parameter(torch.zeros(out))
22
23    def forward(self, X_B):
24        return self.source(X_B) + self.bias
25
26 # training routine of Party B
27 data_loader = ... # Data preparation
28 model = LRPartyB(...)
29 criterion = torch.nn.BCEWithLogitsLoss()
30 fed_optimizer = FederatedSGD(...)
31 for X_B, y in data_loader:
32     output = model(X_B) # forward propagation
33     fed_optimizer.zero_grad()
34     loss = criterion(output, y)
35     loss.backward() # backward propagation
36     fed_optimizer.step() # update paramters

```

Figure 8: Code snippets of the LR algorithm in BLINDFL.

## 6.3 Security Analysis

Indisputably, our protocol in Figure 7 accomplishes all the privacy requirements in Table 3 and the results are lossless as well. To be formal, we identify two ideal functionalities  $\mathcal{F}_{\text{EmbedMatMulFw}}$  and  $\mathcal{F}_{\text{EmbedMatMulBw}}$  for the forward and backward propagation, respectively. In  $\mathcal{F}_{\text{EmbedMatMulFw}}$ , each party inputs the private mini-batch data ( $X_\diamond$ ) and the secretly shared and/or encrypted models (e.g.,  $S_\diamond, \llbracket V_\diamond \rrbracket_\diamond$ ). *Party A* outputs nothing whilst *Party B* outputs  $Z$ . In  $\mathcal{F}_{\text{EmbedMatMulBw}}$ , each party inputs the private mini-batch data and the secretly shared and/or encrypted models, and *Party B* further inputs  $\nabla Z$ . Then, each party outputs the updated models. Then, We present the security guarantees of our Embed-MatMul source layer in Theorem 6.1.

**THEOREM 6.1.** *The protocol of Embed-MatMul source layer securely realizes the ideal functionalities  $\mathcal{F}_{\text{EmbedMatMulFw}}$  and  $\mathcal{F}_{\text{EmbedMatMulBw}}$  in the presence of a semi-honest adversary that can corrupt one party.*

Again, we further analyze the security of  $Z, \nabla Z$  since they are released to *Party B* when the top model is non-federated, which is given in Theorem 6.2.

**THEOREM 6.2.** *Given  $Z, \nabla Z$  in the Embed-MatMul source layer, there are infinite possible values for  $X_A, Q_A, W_A, E_A, E_A W_A$ .*

The proofs are deferred to the appendix due to the space constraint, where we further analyze the security when a federated top model follows our Embed-MatMul source layer.

## 7 IMPLEMENTATION AND EVALUATION

### 7.1 Implementation and Experimental Setup

**Implementation.** We implement BLINDFL on top of PyTorch and the GMP library to be efficient and user-friendly.

**Cryptography Acceleration.** We employ the GMP library to develop an efficient library for the Paillier cryptosystem and support high-performance parallel processing with OpenMP. To vectorize the homomorphic operations, we introduce an abstraction called CryptoTensor, which supports fruitful primitives for both dense and sparse computation of encrypted tensors such as matrix multiplication and scatter addition.

**PyTorch Integration.** As shown in Figure 8, BLINDFL follows the API-style of PyTorch to be user-friendly. The forward and backward procedures of federated source layers are implemented as autograd operations to achieve automatic differentiation. All models are derived from FederatedModule, a wrapper for PyTorch Module that automatically registers the FederatedParameter (e.g.,  $U_\circ, V_\circ$ ). The training routines are the same as non-federated learning, except that we define a FederatedOptimizer to update the secretly shared model weights. BLINDFL has been integrated into the productive pipelines of our industrial partner and deployed to many real-world collaborative applications.

**Experimental Setup.** All experiments are conducted on two servers equipped with 96 cores, 375GB of RAM, and 10Gbps of network bandwidth. For each experiment, we make five runs and report the mean and standard deviation.

**Models.** Since features should be partitioned under the VFL setting, whilst an image or a sentence would not be split and owned by different parties, the inputs of VFL models are usually tabular datasets rather than image or text datasets. Thus, we do not conduct experiments on CV or NLP models in this work. To be specific, We conduct experiments on five widely-used models, which are LR, multinomial LR (MLR), MLP, WDL [9], and DLRM [50].

**Datasets.** As listed in Table 4, we use six public datasets<sup>5</sup> and one industrial advertising dataset. We evenly divide the features for the two parties. Since we focus on the algorithm protocols, we assume datasets are already aligned by the private set intersection (PSI) technique [8, 13, 16], which is a general data preprocessing in VFL [74]. In Section 8, we will discuss how to extend our work when datasets cannot be aligned by PSI.

**Protocols.** As we will show in Section 7.4, BLINDFL achieves comparable model performance as non-federated learning on collocated datasets. Thus, we use the same set of hyper-parameters with learning rate as 0.05, batch size as 128, and embedding dimension as 8. For each model, we train for 10 epochs with momentum SGD, where the momentum value is set as 0.9.

### 7.2 Privacy Preservation

We first conduct experiments to empirically evaluate the robustness of privacy preservation of our work. Since the MPC-based methods can achieve provable security guarantees, we only compare with the split learning based methods in this section.

**Protection for Forward Activations.** As discussed in Section 4.2, once *Party A* gets access to the activations, there is a potential label

Table 4: Description of datasets used in our experiments.

Dataset	#Instances (train/test)	#Features & Avg #nnz	#Classes
<i>a9a</i>	32K/16K	123, 14	2
<i>w8a</i>	50K/15K	300, 12	2
<i>connect-4</i>	50K/17K	126, 42	3
<i>news20</i>	16K/4K	62K, 80	20
<i>higgs</i>	8M/3M	28, 28	2
<i>avazu-app</i>	13M/2M	1M, 14	2
<i>industry</i>	100M/8M	10M, 12	2

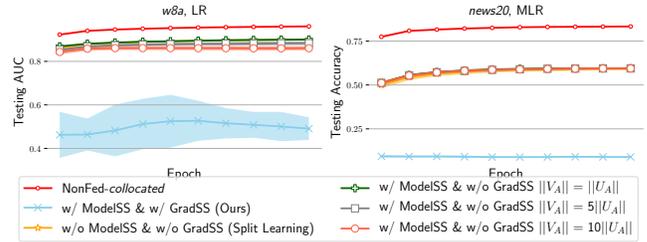


Figure 9: Predicting labels with  $X_A W_A$  or  $X_A U_A$ . We run experiments with or without the SS technique on model weights and gradients (denoted as ModelSS and GradSS).

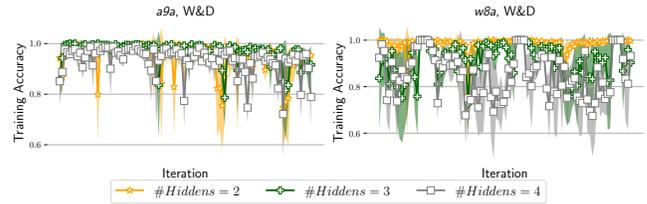


Figure 10: Predicting labels with  $\nabla E_A$ . We vary the number of hidden layers after the embedding table.

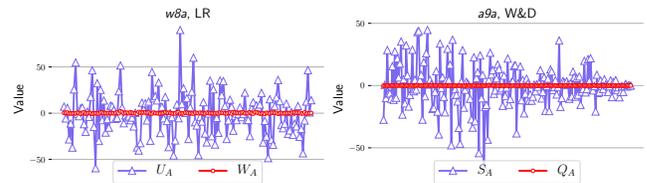


Figure 11: The comparison of model weights ( $W_A, Q_A$ ) and the SS pieces ( $U_A, S_A$ ) on the same coordinate.

leakage. However, many existing works allow *Party A* to manage  $W_A$  when implementing the LR or MLR models [25, 75, 78]. To empirically evaluate the leakage problem, we implement LR and MLR following the split learning paradigm and try to see how accurate it will be if *Party A* predicts the labels via  $X_A W_A$ . For our work, we let *Party A* predict with  $X_A U_A$ . The results are shown in Figure 9.

When *Party A* owns  $W_A$ , although the AUC or accuracy metrics of predicting with  $X_A W_A$  are lower than those of predicting with  $X_A W_A + X_B W_B$ , it is still unsatisfactory — *Party A* is able to make accurate predictions for a large portion of data. For instance, *Party*

<sup>5</sup><https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>

A predicting with  $X_A W_A$  on the *w8a* dataset achieves around 0.9 of AUC on the testing set, which is only 0.06 lower than that of predicting with  $X_A W_A + X_B W_B$ . In contrast, by using our MatMul source layer (Figure 6), the predictions with  $X_A U_A$  are purely random guesses (the AUC metric is around 0.5). It proves that the label information can be protected from the forward activations well with our MatMul source layer since we do not allow *Party A* to manage a local bottom model.

Readers might wonder the necessity of applying SS to model gradients. In other words, if we have already sheltered model weights via  $W_A = U_A + V_A$  on initialization, can we let *Party A* get access to  $\nabla W_A$  and update  $U_A$  directly, rather than updating both  $U_A, V_A$  in the SS manner? To answer this question, we conduct experiments where  $W_A$  is secretly shared on initialization but *Party A* updates  $U_A$  via  $\nabla W_A$ . Furthermore, we amplify the scale of  $V_A$  to better obscure  $W_A$  from *Party A*. Nevertheless, as shown in Figure 9, even though the model weights are secretly shared on initialization, *Party A* can still obtain accurate predictions. Increasing the scale of  $V_A$  is of no use given the slight drop in the AUC/accuracy metrics. This is reasonable because if  $V_A$  keeps unchanged throughout the training phase, we can assume  $\mathbb{E}[X_A V_A] = C$ , where  $C$  is some constant, hence we have  $\mathbb{E}[X_A U_A] = \mathbb{E}[X_A W_A] - C$ . As a result, *Party A* is still capable of making biased predictions, which inevitably leads to the label leakage.

To conclude, compared with the existing works that follow the local bottom model design, BLINDFL is robust to any kind of adversaries that try to learn the private data via forward activations, thanks to our federated source layers.

**Protection for Backward Derivatives.** As discussed in Section 4.2, the backward derivatives are very informative and would leak the training labels. To verify this, we implement the WDL model following the split learning paradigm, i.e., *Party A* owns the embedding table  $Q_A$  in the bottom model and obtains  $\nabla E_A$  in the backward propagation, which is done in several existing works [28, 36, 81]. Then, we let *Party A* predicts the labels via  $\nabla E_A$ .

The results in Figure 10 reveals a horrible fact — *Party A* accurately predicts the labels of almost the entire training datasets via the derivatives  $\nabla E_A$ . Undoubtedly, this disobeys the goal of privacy preservation. Although it focuses on the training data, *Party A* could fit a new model with  $X_A$  and the leaked labels, and utilize the new model to make predictions for more data, leading to a more severe level of label leakage. Moreover, we note that this method attacks successfully regardless of how far  $\nabla E_A$  are away from the labels (i.e., the number of hidden layers between the embedding table and the loss function).

Theoretically speaking, for binary-classification tasks, the backward derivatives for positive and negative instances ought to have opposite directions since they contribute oppositely to the model. Thus, *Party A* can utilize this property and compute the cosine similarity of two derivatives to see whether they have an opposite direction. By doing so, *Party A* achieves an incredible success rate. It gives us a lesson that any algorithms allowing *Party A* to obtain any backward derivatives would be vulnerable from label leakage.

On the contrary, we thoroughly analyze the security of backward derivatives and propose corresponding privacy requirements to forbid data leakage from backward derivatives. For instance, in our Embed-MatMul source layer, *Party A* does not own the embedding

**Table 5: Averaged training time cost of one mini-batch (in seconds). We only record the time cost of matrix multiplication for a fair comparison. The standard deviation for all numbers are smaller than 10% of the mean. SecureML without client aided fails to support the high-dimensional datasets *news20*, *avazu-app*, and *industry*.**

Dataset & Sparsity	Model	Time Cost/Batch (in seconds)		
		BLINDFL	SecureML	SecureML (Client-aided)
<i>a9a</i> (88.72%)	LR	0.018	0.567	0.003
<i>w8a</i> (96.12%)	LR	0.021	1.214	0.016
<i>connect-4</i> (66.67%)	MLP	1.114	5.703	0.008
<i>higgs</i> (Dense)	LR	0.066	0.178	0.002
<i>news20</i> (99.87%)	MLR	1.817	> 1800	0.364
<i>avazu-app</i> (99.99%)	LR	0.038	OOM	4.727
<i>industry</i> (99.99%)	LR	0.034	OOM	47.083

table in plaintext and only gets access to the encrypted derivatives  $\llbracket \nabla E_A \rrbracket_B$ , so BLINDFL is robust to any adversaries that try to peek the private data through backward derivatives.

**Protection for Models.** As we have discussed in Section 4.2, the magnitudes or signs of model weights or gradients express the feature importance upon the learning tasks, so we forbid the models to be obtained. To achieve this goal, we expect the SS pieces of model weights and gradients do not reveal the magnitudes or signs of the ground truth values.

To illustrate the effectiveness of our protection for models, we plot the values of SS pieces and the ground truth values in Figure 11. Due to the space constraints, we only plot the model weights  $W_A, S_A$  and the sharing pieces  $U_A, Q_A$ , whilst similar results are observed for other model weights and gradients. Figure 11 shows that the difference on each coordinate is random and sufficiently large so that both the magnitudes or signs of the ground truth values are inaccessible. As a result, each party cannot infer any information by analyzing the feature importance.

### 7.3 Efficiency

Apart from BLINDFL, the MPC-based methods can also achieve robust security guarantees via data outsourcing. However, unlike these methods, our work is able to utilize the sparsity of features to speed up the computation. To verify this, we conduct experiments to compare the efficiency of BLINDFL and the existing MPC-based counterparts. In particular, we consider SecureML [49], which secretly shares the model weights and feature values via ABY [12], and performs matrix multiplication with the help of Beaver triplets. Moreover, as introduced by Mohassel and Zhang [49], SecureML has a client-aided variant that offloads the generation of Beaver triplets to a non-colluding third party so that no cryptographic operation will be involved in the ML tasks. We compare BLINDFL with both variants of SecureML in terms of the training time cost. Furthermore, we only record the time cost of matrix multiplication and exclude that of non-linear activation functions (e.g., Sigmoid, ReLU) for a fair comparison. The results are given in Table 5.

**BLINDFL vs. SecureML.** We first discuss the performance without client aided on different kinds of datasets, respectively.

First, although SecureML can support the low-dimensional and sparse datasets, i.e., the *a9a*, *w8a*, and *connect-4* datasets, it cannot

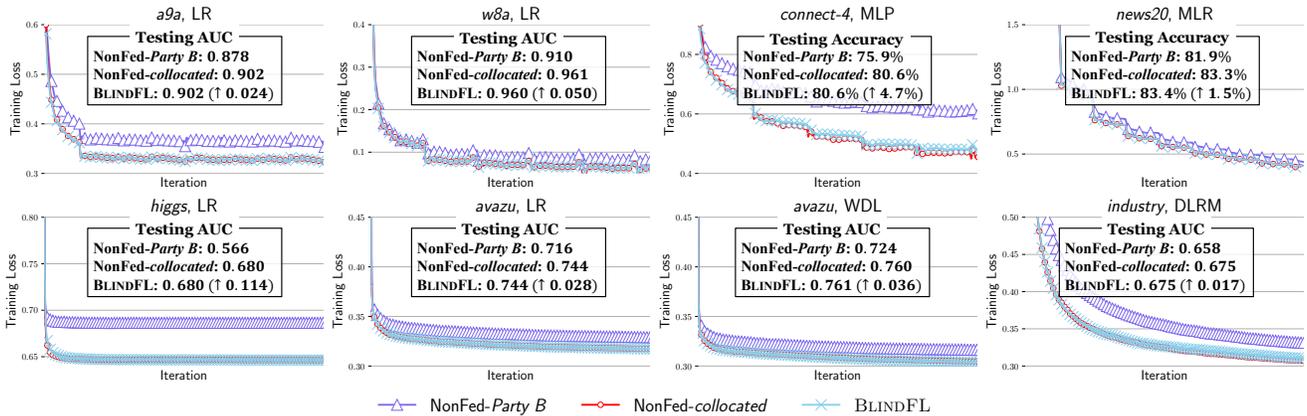


Figure 12: Training loss in terms of iterations and testing AUC or accuracy metrics. The standard deviation for all evaluation metrics are smaller than 5% of the mean. The model performance of BLINDFL is comparable to that of non-federated learning on collocated datasets, and better than that of non-federated learning on Party B’s features only.

take the sparsity into account, leading to umpteenth redundant computations on the zero values. On the contrary, by keeping the data inside each party, BLINDFL supports sparse matrix multiplication and therefore outperforms SecureML by a large extent. Furthermore, BLINDFL gains more improvements when the datasets are sparser, and can be over 50× faster than SecureML.

Second, for the *higgs* dataset, which is low-dimensional and dense, the gap between BLINDFL and SecureML is smaller since there are no zero values. Nevertheless, our work still performs better and gives about 5× of acceleration. The reason is that we keep an encrypted version of  $\llbracket V_\diamond \rrbracket_\diamond$  on Party  $\diamond$  to reduce an extra communication round, whilst SecureML needs to generate new Beaver triplets for each iteration. Thus, our MatMul source layer also supports dense features well.

Third, SecureML fails to support the *news20*, *avazu-app*, and *industry* datasets due to their high dimensionalities (either fails to accomplish the computation within a reasonable time or runs out of memory). In contrast, BLINDFL handles these high-dimensional and sparse datasets well since we do not need to store or process the zero values. To the best of our knowledge, none of the existing works are able to support such a scale of datasets whilst guarantee promising privacy.

**BLINDFL vs. Client-aided SecureML.** The client-aided SecureML runs much faster when the dimensionality is not high as no cryptography operations are needed, whilst BLINDFL cannot be accelerated in this way. Nevertheless, for the *avazu-app* and *industry* datasets, which are extremely high-dimensional and sparse, BLINDFL outperforms the client-aided SecureML. This is reasonable since SecureML has to process all the dimensions whilst BLINDFL only needs to consider the non-zero ones. Furthermore, since there is no such a non-colluding third party in practice, the MPC-based methods would be even slower. As a result, BLINDFL is more suitable for the sparse datasets in real-world VFL applications.

To summarize, BLINDFL can outperform the MPC-based counterparts by an order of magnitude on sparse datasets and supports datasets with higher dimensionalities. Moreover, as we will illustrate in Section 7.4, BLINDFL can also support the embedding lookup

operation for categorical inputs, whilst there are no MPC-based methods that are designed for the embedding lookup operations to the best of knowledge. As a result, BLINDFL is more powerful as it gains much better efficiency and functionality.

**Scalability.** In general, the time cost of the federated source layers dominates the total running time since the cryptography operations are time-expensive whilst the top model is a non-federated module. Thus, the running time of our work is proportional to the output dimensionality of the source layer, and does not change significantly w.r.t. the number of layers. Due to the space constraint, we defer the experimental results in our appendix and do not discuss the details here.

## 7.4 Model Ability

As introduced in Section 1, the VFL models are expected to be lossless, i.e., the performance of VFL models should be (i) better than non-federated learning on the features of Party B only (denoted as NonFed-Party B); and (ii) comparable to non-federated learning on collocated features of both parties (denoted as NonFed-collocated). To assess the lossless property of BLINDFL, we conduct experiments on extensive datasets and models and show the results in Figure 12.

On all experiments, the convergence of BLINDFL is consistent to that of NonFed-collocated and is better than that of NonFed-Party B. The evaluation metrics on testing sets have similar observations as well — the testing AUC or accuracy metrics of BLINDFL and NonFed-collocated are comparable, which are better than those of NonFed-Party B with the help of the extra features from Party A. The experimental results verify the lossless property of BLINDFL. In practice, it is forbidden to collect the data from different parties together for ML training or inference. As a result, BLINDFL is more superior as it unites the features in different parties to enhance the model ability whilst guarantees the data privacy of all parties.

**Summary.** To conclude, BLINDFL outperforms the existing works by providing more robust privacy guarantees and faster running speed. Furthermore, BLINDFL achieves comparable model ability as collocated learning for a variety of models and datasets. Consequently, BLINDFL is a good fit for a wide range of VFL applications.

## 8 RELATED WORKS

**Vertical FL (VFL).** With the ever-evolving concerns on data privacy, how to build ML models over different data sources with privacy preservation is gaining more popularity. VFL considers the case where different parties jointly train ML models over the partitioned features, which fits numerous real-world cross-enterprise collaborations [10, 19, 71, 74, 75]. We categorize the existing works into two lines based on how they process the private features.

**The MPC-based methods.** Since cryptographic methods can deliver promising security guarantees, it is a straightforward idea to apply them to the private data directly. A number of works encrypt features via HE and feed the encrypted features to ML models [4, 22, 26, 63, 77]. To perform both addition and multiplication, these works adopt somewhat HE or fully HE, which are extremely time-consuming. Other works adopt SS to protect the data. Bogdanov et al. [6] proposed a data analysis framework based on additive SS. Makri et al. [43] trained the support vector machine models for image classification. SecureML [49] supports various ML models via the ABY scheme [12]. ABY3 [48] extends ABY to the three-party scenarios. SecureNN [66, 67] optimizes for three-party neural networks. Nevertheless, although these works convey considerable privacy guarantees, the generic MPC framework involves sophisticated protocols and time-consuming computation primitives to achieve zero knowledge disclosure. Worse still, outsourcing the datasets is not suitable for sparse or categorical features.

**Split learning.** Rather than data outsourcing, the split learning paradigm [62, 64] keeps the private data inside each party and leverages a local bottom model to process the features. Many VFL algorithms are designed in this way [28, 36, 75, 78, 81]. Since features are visible to the owner party, it is very flexible to support various kinds of features. However, as discussed in Section 3, the design of local bottom models faces several data leakage problems and fails to convey provable security guarantees. Although some existing works try to apply some cryptography operations to enhance privacy, they still suffer from data leakage. For instance, Yang et al. [75] tried to avoid label leakage from derivatives via HE and SS. However, *Party B* can still infer the labels as depicted in Figure 2. In addition, noisy perturbation (or differential privacy) is also utilized to blur the intermediate results [36]. Nevertheless, adding noises inevitably affects the model accuracy. In practice, it is non-trivial to figure out a suitable amount of noises that conveys considerable security whilst produces desirable model accuracy. Furthermore, all these works cannot provide security guarantees under the ideal-real paradigm as the MPC-based methods.

**Data alignment assumption.** Although most of the VFL algorithms assume the instances of different parties have been aligned via the PSI technique [8, 13, 16], there are also works that focus on the cases where data cannot be aligned due to stronger privacy requirements [42, 61]. However, our work can be extended to these cases easily. For instance, Liu et al. [42] proposed that only *Party B* can obtain the intersection whilst *Party A* cannot. Then, for instances outside the intersection, *Party B* sets their derivatives as zeros so that model gradients will not be affected. This method can be applied to our work by tweaking Line 9 of Figure 6 and Line 12 of Figure 7. For another example, in [61], both parties can only obtain the union instead of intersection. Then, they propose

to generate synthetic features and labels for instances outside the intersection. Obviously, this technique can also be integrated with our work during the data preprocessing phase.

**Reconstruction-based inference attacks.** There are also many works that try to break the privacy of ML models via more advanced attacks. For instance, membership inference attacks (MIAs) [27, 59] try to infer whether an instance was used to train a given model. However, in our work, since the federated source layers are not released in plaintext, none of the parties (or any other attackers) can apply MIAs to attack the models. Deep leakage from gradients (DLG) is also a popular kind of attack methods [20, 83] which recover the private data from model gradients. However, model gradients are not disclosed to others in our algorithm protocol.

**Horizontal FL (HFL).** Besides VFL, HFL [33, 34, 45] is also one of the most popular fields of FL, where all parties have disjoint instance sets for the same features, i.e., the datasets are horizontally partitioned. There has been a wide range of works studying HFL [7, 35, 38]. The most essential privacy-preserving technique in this field is differential privacy (DP) [3, 68], which is adopted to obscure the model weights or gradients so that the individual data become indistinguishable. Our work differs from these works in objection since we consider the VFL scenario.

**Other related studies.** Beyond FL, there is also an arousing interest in various kinds of federated computing. For instance, the private set intersection technique [8, 13, 16] secretly extracts and aligns the joint set of private tables in different parties. Apart from the “join” operation, many other operations over private databases are studied [29, 40, 82]. Achieving privacy preservation from the hardware perspective is also a popular topic. With the help of the trusted execution environments (TEEs) such as Intel SGX and AMD memory encryption [32, 44], we can get rid of umpteenth expensive cryptographic arithmetics when implementing privacy-preserving ML algorithms [52]. Nevertheless, the available memory for TEE is too small to support a large data volume and it requires extra hardware supports.

## 9 CONCLUSION AND FUTURE DIRECTIONS

This work proposed BLINDFL, a brand new VFL framework. To address the functionality of VFL models, we designed the federated source layers to unite the data from different data sources. To protect data privacy, we analyzed the privacy requirements in-depth and carefully devised safe and accurate algorithm protocols. Experimental results show that BLINDFL is able to support various kinds of features efficiently and achieve promising privacy guarantees.

Beyond this work, we wish to strengthen BLINDFL in two directions. First, it is worthy to study how to extract feature interactions between parties to support more ML models (e.g., factorization machines). Second, how to apply adaptive optimizers (e.g., Adam) to secretly shared model gradients is also an interesting topic. We will leave the exploration of these topics as our future works.

## ACKNOWLEDGMENTS

This work is supported by National Natural Science Foundation of China (NSFC) (No. 61832001), PKU-Tencent joint research Lab, and Beijing Academy of Artificial Intelligence (BAAI). Bin Cui is the corresponding author.

## REFERENCES

- [1] 2018. California Consumer Privacy Act (CCPA). <https://oag.ca.gov/privacy/ccpa>.
- [2] 2021. Virginia Consumer Data Protection Act (CDPA). <https://lis.virginia.gov/cgi-bin/legp604.exe?212+sum+HB2307>.
- [3] Martin Abadi, Andy Chu, Ian J. Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep Learning with Differential Privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, CCS 2016*. ACM, 308–318.
- [4] Yoshinori Aono, Takuya Hayashi, Le Trieu Phong, and Lihua Wang. 2016. Scalable and Secure Logistic Regression via Homomorphic Encryption. *IACR Cryptol. ePrint Arch.* 2016 (2016), 111.
- [5] Donald Beaver. 1991. Efficient Multiparty Protocols Using Circuit Randomization. In *Advances in Cryptology - CRYPTO '91, 11th Annual International Cryptology Conference*, Vol. 576. Springer, 420–432.
- [6] Dan Bogdanov, Sven Laur, and Jan Willemson. 2008. Sharemind: A Framework for Fast Privacy-Preserving Computations. In *13th European Symposium on Research in Computer Security, ESORICS 2008*, Vol. 5283. Springer, 192–206.
- [7] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir Ivanov, Chloé Kiddon, Jakub Konečný, Stefano Mazzocchi, Brendan McMahan, Timon Van Overveldt, David Petrou, Daniel Ramage, and Jason Roselander. 2019. Towards Federated Learning at Scale: System Design. In *Proceedings of Machine Learning and Systems 2019, MLSys 2019*. mlsys.org.
- [8] Hao Chen, Kim Laine, and Peter Rindal. 2017. Fast Private Set Intersection from Homomorphic Encryption. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS 2017*. ACM, 1243–1255.
- [9] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishikesh Aradhya, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Isipir, Rohan Anil, Zakaria Haque, Lichan Hong, Vishal Jain, Xiaobing Liu, and Hemal Shah. 2016. Wide & Deep Learning for Recommender Systems. In *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems, DLRS@RecSys 2016*. ACM, 7–10.
- [10] Kewei Cheng, Tao Fan, Yilun Jin, Yang Liu, Tianjian Chen, Dimitrios Papadopoulos, and Qiang Yang. 2021. SecureBoost: A Lossless Federated Learning Framework. *IEEE Intell. Syst.* 36, 6 (2021), 87–98.
- [11] Jung Hee Cheon, Andrey Kim, Miran Kim, and Yong Soo Song. 2017. Homomorphic Encryption for Arithmetic of Approximate Numbers. In *23rd International Conference on the Theory and Applications of Cryptology and Information Security, ASIACRYPT 2017*, Vol. 10624. Springer, 409–437.
- [12] Daniel Demmler, Thomas Schneider, and Michael Zohner. 2015. ABY - A Framework for Efficient Mixed-Protocol Secure Two-Party Computation. In *22nd Annual Network and Distributed System Security Symposium, NDSS 2015*. The Internet Society.
- [13] Changyu Dong, Liqun Chen, and Zikai Wen. 2013. When private set intersection meets big data: an efficient and scalable protocol. In *Proceedings of the 2013 ACM SIGSAC Conference on Computer and Communications Security, CCS 2013*. ACM, 789–800.
- [14] Junfeng Fan and Frederik Vercauteren. 2012. Somewhat Practical Fully Homomorphic Encryption. *IACR Cryptol. ePrint Arch.* 2012 (2012), 144.
- [15] Caroline Fontaine and Fabien Galand. 2007. A Survey of Homomorphic Encryption for Nonspecialists. *EURASIP J. Inf. Secur.* 2007 (2007).
- [16] Michael J. Freedman, Kobbi Nissim, and Benny Pinkas. 2004. Efficient Private Matching and Set Intersection. In *International Conference on the Theory and Applications of Cryptographic Techniques, EUROCRYPT 2004*, Vol. 3027. Springer, 1–19.
- [17] Fangcheng Fu, Yuzheng Hu, Yihan He, Jiawei Jiang, Yingxia Shao, Ce Zhang, and Bin Cui. 2020. Don't Waste Your Bits! Squeeze Activations and Gradients for Deep Neural Networks via TinyScript. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020*, Vol. 119. PMLR, 3304–3314.
- [18] Fangcheng Fu, Jiawei Jiang, Yingxia Shao, and Bin Cui. 2019. An Experimental Evaluation of Large Scale GBDT Systems. *Proc. VLDB Endow.* 12, 11 (2019), 1357–1370.
- [19] Fangcheng Fu, Yingxia Shao, Lele Yu, Jiawei Jiang, Huanran Xue, Yangyu Tao, and Bin Cui. 2021. VF<sup>2</sup>Boost: Very Fast Vertical Federated Gradient Boosting for Cross-Enterprise Learning. In *Proceedings of the 2021 International Conference on Management of Data, SIGMOD 2021*. ACM, 563–576.
- [20] Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. 2020. Inverting Gradients - How easy is it to break privacy in federated learning?. In *Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020*.
- [21] Craig Gentry. 2009. Fully homomorphic encryption using ideal lattices. In *Proceedings of the 41st Annual ACM Symposium on Theory of Computing, STOC 2009*. ACM, 169–178.
- [22] Ran Gilad-Bachrach, Nathan Dowlin, Kim Laine, Kristin E. Lauter, Michael Naehrig, and John Wernsing. 2016. CryptoNets: Applying Neural Networks to Encrypted Data with High Throughput and Accuracy. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016*, Vol. 48. JMLR.org, 201–210.
- [23] Oded Goldreich. 2004. *The Foundations of Cryptography - Volume 2: Basic Applications*. Cambridge University Press.
- [24] Hakan Hacigümüs, Sharad Mehrotra, and Balakrishna R. Iyer. 2002. Providing Database as a Service. In *18th IEEE International Conference on Data Engineering, ICDE 2002*. IEEE Computer Society, 29–38.
- [25] Stephen Hardy, Wilko Henecka, Hamish Ivey-Law, Richard Nock, Giorgio Patrini, Guillaume Smith, and Brian Thorne. 2017. Private federated learning on vertically partitioned data via entity resolution and additively homomorphic encryption. *CoRR abs/1711.10677* (2017). arXiv:1711.10677
- [26] Ehsan Hesamifard, Hassan Takabi, and Mehdi Ghasemi. 2017. CryptoDL: Deep Neural Networks over Encrypted Data. *CoRR abs/1711.05189* (2017). arXiv:1711.05189
- [27] Hongsheng Hu, Zoran Salcic, Gillian Dobbie, and Xuyun Zhang. 2021. Membership Inference Attacks on Machine Learning: A Survey. *CoRR abs/2103.07853* (2021). arXiv:2103.07853
- [28] Yaochen Hu, Di Niu, Jianming Yang, and Shengping Zhou. 2019. FDML: A Collaborative Machine Learning Framework for Distributed Features. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2019*. ACM, 2232–2240.
- [29] Tianxi Ji, Pan Li, Emre Yilmaz, Erman Arday, Yanfang Ye, and Jinyuan Sun. 2021. Differentially Private Binary- and Matrix-Valued Data Query: An XOR Mechanism. *Proc. VLDB Endow.* 14, 5 (2021), 849–862.
- [30] Jiawei Jiang, Bin Cui, Ce Zhang, and Fangcheng Fu. 2018. DimBoost: Boosting Gradient Boosting Decision Tree to Higher Dimensions. In *Proceedings of the 2018 International Conference on Management of Data, SIGMOD 2018*. ACM, 1363–1376.
- [31] Jiawei Jiang, Fangcheng Fu, Tong Yang, and Bin Cui. 2018. SketchML: Accelerating Distributed Machine Learning with Data Sketches. In *Proceedings of the 2018 International Conference on Management of Data, SIGMOD 2018*. ACM, 1269–1284.
- [32] David Kaplan, Jeremy Powell, and Tom Woller. 2016. AMD memory encryption. *White paper* (2016).
- [33] Jakub Konečný, H. Brendan McMahan, Daniel Ramage, and Peter Richtárik. 2016. Federated Optimization: Distributed Machine Learning for On-Device Intelligence. *CoRR abs/1610.02527* (2016). arXiv:1610.02527
- [34] Jakub Konečný, H. Brendan McMahan, Felix X. Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. 2016. Federated Learning: Strategies for Improving Communication Efficiency. *CoRR abs/1610.05492* (2016). arXiv:1610.05492
- [35] Mathias Lécuyer, Riley Spahn, Kiran Vodrahalli, Roxana Geambasu, and Daniel Hsu. 2019. Privacy accounting and quality control in the sage differentially private ML platform. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles, SOSP 2019*. ACM, 181–195.
- [36] Oscar Li, Jiankai Sun, Xin Yang, Weihao Gao, Hongyi Zhang, Junyuan Xie, Virginia Smith, and Chong Wang. 2022. Label Leakage and Protection in Two-party Split Learning. In *International Conference on Learning Representations, ICLR 2022*.
- [37] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. 2020. Federated Learning: Challenges, Methods, and Future Directions. *IEEE Signal Process. Mag.* 37, 3 (2020), 50–60.
- [38] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. 2020. Federated Optimization in Heterogeneous Networks. In *Proceedings of Machine Learning and Systems 2020, MLSys 2020*. mlsys.org.
- [39] Yang Li, Yu Shen, Wentao Zhang, Yuanwei Chen, Huaqun Jiang, Mingchao Liu, Jiawei Jiang, Jinyang Gao, Wentao Wu, Zhi Yang, Ce Zhang, and Bin Cui. 2021. OpenBox: A Generalized Black-box Optimization Service. In *Proceedings of the 27th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2021*. ACM, 3209–3219.
- [40] Gang Liang and Sudarshan S. Chawathe. 2004. Privacy-Preserving Inter-database Operations. In *Second Symposium on Intelligence and Security Informatics, ISI 2004*, Vol. 3073. Springer, 66–82.
- [41] Yehuda Lindell. 2017. How to Simulate It - A Tutorial on the Simulation Proof Technique. In *Tutorials on the Foundations of Cryptography*. Springer International Publishing, 277–346.
- [42] Yang Liu, Xiong Zhang, and Libin Wang. 2020. Asymmetrical Vertical Federated Learning. *CoRR abs/2004.07427* (2020). arXiv:2004.07427
- [43] Eleftheria Makri, Dragos Rotaru, Nigel P. Smart, and Frederik Vercauteren. 2019. EPIC: Efficient Private Image Classification (or: Learning from the Masters). In *The Cryptographers' Track at the RSA Conference 2019, CT-RSA 2019*, Vol. 11405. Springer, 473–492.
- [44] Frank McKeen, Ilya Alexandrovich, Alex Berenzon, Carlos V. Rozas, Hisham Shafi, Vedvyas Shanbhogue, and Uday R. Savagaonkar. 2013. Innovative instructions and software model for isolated execution. In *The Second Workshop on Hardware and Architectural Support for Security and Privacy, HASP 2013*. ACM, 10.
- [45] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. 2017. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, AISTATS 2017*, Vol. 54. PMLR, 1273–1282.
- [46] Xupeng Miao, Xiaonan Nie, Yingxia Shao, Zhi Yang, Jiawei Jiang, Lingxiao Ma, and Bin Cui. 2021. Heterogeneity-Aware Distributed Machine Learning Training via Partial Reduce. In *Proceedings of the 2021 International Conference on*

- Management of Data, SIGMOD 2021*. ACM, 2262–2270.
- [47] Xupeng Miao, Hailin Zhang, Yining Shi, Xiaonan Nie, Zhi Yang, Yangyu Tao, and Bin Cui. 2021. HET: Scaling out Huge Embedding Model Training via Cache-enabled Distributed Framework. *Proc. VLDB Endow.* 15, 2 (2021), 312–320.
- [48] Payman Mohassel and Peter Rindal. 2018. ABY<sup>3</sup>: A Mixed Protocol Framework for Machine Learning. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, CCS 2018*. ACM, 35–52.
- [49] Payman Mohassel and Yupeng Zhang. 2017. SecureML: A System for Scalable Privacy-Preserving Machine Learning. In *2017 IEEE Symposium on Security and Privacy, SP 2017*. IEEE Computer Society, 19–38.
- [50] Maxim Naumov, Dheevatsa Mudigere, Hao-Jun Michael Shi, Jianyu Huang, Narayanan Sundaraman, Jongsoo Park, Xiaodong Wang, Udit Gupta, Carole-Jean Wu, Alisson G. Azzolini, Dmytro Dzulgakov, Andrey Malleevich, Iliia Cherniavskii, Yinghai Lu, Raghuraman Krishnamoorthi, Ansha Yu, Volodymyr Kondratenko, Stephanie Pereira, Xianjie Chen, Wenlin Chen, Vijay Rao, Bill Jia, Liang Xiong, and Misha Smelyanskiy. 2019. Deep Learning Recommendation Model for Personalization and Recommendation Systems. *CoRR abs/1906.00091* (2019). arXiv:1906.00091
- [51] Xiaonan Nie, Shijie Cao, Xupeng Miao, Lingxiao Ma, Jilong Xue, Youshan Miao, Zichao Yang, Zhi Yang, and Bin Cui. 2021. Dense-to-Sparse Gate for Mixture-of-Experts. *CoRR abs/2112.14397* (2021). arXiv:2112.14397
- [52] Olga Ohrimenko, Felix Schuster, Cédric Fournet, Aastha Mehta, Sebastian Nowozin, Kapil Vaswani, and Manuel Costa. 2016. Oblivious Multi-Party Machine Learning on Trusted Processors. In *25th USENIX Security Symposium, USENIX Security 16*. USENIX Association, 619–636.
- [53] Pascal Paillier. 1999. Public-Key Cryptosystems Based on Composite Degree Residuosity Classes. In *International Conference on the Theory and Applications of Cryptographic Techniques, EUROCRYPT 1999*, Vol. 1592. Springer, 223–238.
- [54] Ignacio Cascudo Pueyo, Ronald Cramer, and Chaoping Xing. 2011. The Torsion-Limit for Algebraic Function Fields and Its Application to Arithmetic Secret Sharing. In *31st Annual Cryptology Conference, CRYPTO 2011*, Vol. 6841. Springer, 685–705.
- [55] Pille Pullonen et al. 2013. Actively secure two-party computation: Efficient Beaver triple generation. *Instructor* (2013).
- [56] Ronald L Rivest, Len Adleman, Michael L Dertouzos, et al. 1978. On data banks and privacy homomorphisms. *Foundations of secure computation* 4, 11 (1978), 169–180.
- [57] Ronald L. Rivest, Adi Shamir, and Leonard M. Adleman. 1978. A Method for Obtaining Digital Signatures and Public-Key Cryptosystems. *Commun. ACM* 21, 2 (1978), 120–126.
- [58] Adi Shamir. 1979. How to Share a Secret. *Commun. ACM* 22, 11 (1979), 612–613.
- [59] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership Inference Attacks Against Machine Learning Models. In *2017 IEEE Symposium on Security and Privacy, SP 2017*. IEEE Computer Society, 3–18.
- [60] Radu Sion. 2007. Secure Data Outsourcing. In *Proceedings of the 33rd International Conference on Very Large Data Bases, VLDB 2007*. ACM, 1431–1432.
- [61] Jiankai Sun, Xin Yang, Yuanshun Yao, Aonan Zhang, Weihao Gao, Junyuan Xie, and Chong Wang. 2021. Vertical Federated Learning without Revealing Intersection Membership. *CoRR abs/2106.05508* (2021). arXiv:2106.05508
- [62] Chandra Thapa, Mahawaga Arachchige Pathum Chamikara, and Seyit Camtepe. 2020. SplitFed: When Federated Learning Meets Split Learning. *CoRR abs/2004.12088* (2020). arXiv:2004.12088
- [63] Tim van Elsloo, Giorgio Patrini, and Hamish Ivey-Law. 2019. SEALion: a Framework for Neural Network Inference on Encrypted Data. *CoRR abs/1904.12840* (2019). arXiv:1904.12840
- [64] Praneeth Vepakomma, Otkrist Gupta, Tristan Swedish, and Ramesh Raskar. 2018. Split learning for health: Distributed deep learning without sharing raw patient data. *CoRR abs/1812.00564* (2018). arXiv:1812.00564
- [65] Paul Voigt and Axel Von dem Bussche. 2017. The eu general data protection regulation (gdpr). *A Practical Guide, 1st Ed., Cham: Springer International Publishing* (2017).
- [66] Sameer Wagh, Divya Gupta, and Nishanth Chandran. 2018. SecureNN: Efficient and Private Neural Network Training. *IACR Cryptol. ePrint Arch.* 2018 (2018), 442.
- [67] Sameer Wagh, Divya Gupta, and Nishanth Chandran. 2019. SecureNN: 3-Party Secure Computation for Neural Network Training. *Proc. Priv. Enhancing Technol.* 2019, 3 (2019), 26–49.
- [68] Kang Wei, Jun Li, Ming Ding, Chuan Ma, Howard H. Yang, Farhad Farokhi, Shi Jin, Tony Q. S. Quek, and H. Vincent Poor. 2020. Federated Learning With Differential Privacy: Algorithms and Performance Analysis. *IEEE Trans. Inf. Forensics Secur.* 15 (2020), 3454–3469.
- [69] Kang Wei, Jun Li, Chuan Ma, Ming Ding, Sha Wei, Fan Wu, Guihai Chen, and Thilina Ranbaduge. 2022. Vertical Federated Learning: Challenges, Methodologies and Experiments. *CoRR abs/2202.04309* (2022). arXiv:2202.04309
- [70] Man Wei, Siwei Sun, Zihao Wei, and Lei Hu. 2021. Unbalanced sharing: a threshold implementation of SM4. *Sci. China Inf. Sci.* 64, 5 (2021).
- [71] Yuncheng Wu, Shaofeng Cai, Xiaokui Xiao, Gang Chen, and Beng Chin Ooi. 2020. Privacy Preserving Vertical Federated Learning for Tree-based Models. *Proc. VLDB Endow.* 13, 11 (2020), 2090–2103.
- [72] Xu Xie, Fei Sun, Zhaoyang Liu, Shiwen Wu, Jinyang Gao, Bolin Ding, and Bin Cui. 2020. Contrastive Learning for Sequential Recommendation. *CoRR abs/2010.14395* (2020). arXiv:2010.14395
- [73] Xu Xie, Fei Sun, Xiaoyong Yang, Zhao Yang, Jinyang Gao, Wenwu Ou, and Bin Cui. 2021. Explore User Neighborhood for Real-time E-commerce Recommendation. In *37th IEEE International Conference on Data Engineering, ICDE 2021*. IEEE, 2464–2475.
- [74] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. 2019. Federated Machine Learning: Concept and Applications. *ACM Trans. Intell. Syst. Technol.* 10, 2 (2019), 12:1–12:19.
- [75] Shengwen Yang, Bing Ren, Xuhui Zhou, and Liping Liu. 2019. Parallel Distributed Logistic Regression for Vertical Federated Learning without Third-Party Coordinator. *CoRR abs/1911.09824* (2019). arXiv:1911.09824
- [76] Andrew Chi-Chih Yao. 1982. Protocols for Secure Computations (Extended Abstract). In *23rd Annual Symposium on Foundations of Computer Science, FOCS 1982*. IEEE Computer Society, 160–164.
- [77] Jiawei Yuan and Shucheng Yu. 2014. Privacy Preserving Back-Propagation Neural Network Learning Made Practical with Cloud Computing. *IEEE Trans. Parallel Distributed Syst.* 25, 1 (2014), 212–221.
- [78] Qingsong Zhang, Bin Gu, Cheng Deng, and Heng Huang. 2021. Secure Bilevel Asynchronous Vertical Federated Learning with Backward Updating. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021*. AAAI Press, 10896–10904.
- [79] Wentao Zhang, Jiawei Jiang, Yingxia Shao, and Bin Cui. 2020. Snapshot boosting: a fast ensemble framework for deep neural networks. *Sci. China Inf. Sci.* 63, 1 (2020), 112102.
- [80] Wentao Zhang, Yu Shen, Zheyu Lin, Yang Li, Xiaosen Li, Wen Ouyang, Yangyu Tao, Zhi Yang, and Bin Cui. 2022. PaSca: a Graph Neural Architecture Search System under the Scalable Paradigm. *CoRR abs/2203.00638* (2022). arXiv:2203.00638
- [81] Yifei Zhang and Hao Zhu. 2020. Additively Homomorphical Encryption based Deep Neural Network for Asymmetrically Collaborative Machine Learning. *CoRR abs/2007.06849* (2020). arXiv:2007.06849
- [82] Wenting Zheng, Ankur Dave, Jethro G. Beekman, Raluca Ada Popa, Joseph E. Gonzalez, and Ion Stoica. 2017. Opaque: An Oblivious and Encrypted Distributed Analytics Platform. In *14th USENIX Symposium on Networked Systems Design and Implementation, NSDI 2017*. USENIX Association, 283–298.
- [83] Ligeng Zhu, Zhijian Liu, and Song Han. 2019. Deep Leakage from Gradients. In *Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019*. 14747–14756.

## A SECURITY ANALYSIS

In this section, we provide the definition of our ideal functionalities and proofs for the theorems in Section 5.3 and Section 6.3.

### A.1 Proofs for Section 5.3

**THEOREM 5.1.** *The protocol of MatMul source layer securely realizes the ideal functionalities  $\mathcal{F}_{\text{MatMulFw}}$  and  $\mathcal{F}_{\text{MatMulBw}}$  in the presence of a semi-honest adversary that can corrupt one party.*

**PROOF.** We formally provide the definition and prove the security of the two idea functionalities in Lemma A.1 and A.2, respectively. Putting them together, we complete the proof for Theorem 5.1.  $\square$

Forward propagation of the MatMul source layer  $\mathcal{F}_{\text{MatMulFw}}$

**Inputs:**

- ▷ *Party A* inputs features  $X_A$ , secretly shared and/or encrypted models  $U_A, \llbracket V_A \rrbracket_B$ , and keys  $sk_A, pk_B$ ;
- ▷ *Party B* inputs features  $X_B$ , secretly shared and/or encrypted models  $U_B, \llbracket V_B \rrbracket_A$ , and keys  $pk_A, sk_B$ .

**Outputs:**

- ▷ *Party A* outputs nothing;
- ▷ *Party B* outputs  $Z = X_A W_A + X_B W_B$ , where  $W_\circ = U_\circ + V_\circ$ .

**LEMMA A.1.** *The protocol  $\Pi_{\text{MatMulFw}}$  in Line 5-8 of Figure 6 securely realizes  $\mathcal{F}_{\text{MatMulFw}}$  in the presence of a semi-honest adversary that can corrupt one party.*

**PROOF.** First, *Party A* only receives one message by invoking the  $\Pi_{\text{HE2SS}}$  protocol (Line 6), whilst the other values are computed locally. Therefore, the view of *Party A* can be perfectly simulated by simulating the  $\mathcal{F}_{\text{HE2SS}}$  functionality as discussed in Lemma A.5.

Second, *Party B* receives two messages in the protocol (Line 5 and Line 8). Similarly, the first message can be simulated by simulating the  $\mathcal{F}_{\text{HE2SS}}$  functionality. Denote the simulated versions of  $X_A V_A - \varepsilon_A$  and  $\varepsilon_B$  are  $(X_A V_A - \varepsilon_A)^*$  and  $\varepsilon_B^*$ , respectively. Then, we can simulate the second message, i.e.  $Z'_A$ , by computing  $Z'^*_A = Z - Z'^*_B$ , where  $Z'^*_B = X_B U_B + \varepsilon_B^* + (X_A V_A - \varepsilon_A)^*$ . Since both  $Z'_A$  (or  $Z'^*_A$ ) and  $Z^*_A$  (or  $Z'^*_B$ ) represent one piece of sharing of the output  $Z$ , they have the same probability distribution. As a result, we perfectly simulate the view of *Party B*.  $\square$

Backward propagation of the MatMul source layer  $\mathcal{F}_{\text{MatMulBw}}$

**Inputs:**

- ▷ *Party A* inputs features  $X_A$ , secretly shared and/or encrypted models  $U_A, \llbracket V_A \rrbracket_B$ , and keys  $sk_A, pk_B$ ;
- ▷ *Party B* inputs features  $X_B$ , derivatives  $\nabla Z$ , secretly shared and/or encrypted models  $U_B, \llbracket V_B \rrbracket_A$ , and keys  $pk_A, sk_B$ .

**Outputs:**

- ▷ *Party A* outputs  $\phi$ ;
- ▷ *Party B* outputs  $\nabla W_A - \phi, \nabla W_B$ , where  $\nabla W_\circ = X_\circ^T \nabla Z$ .

**LEMMA A.2.** *The protocol  $\Pi_{\text{MatMulBw}}$  in Line 9-12 of Figure 6 securely realizes  $\mathcal{F}_{\text{MatMulBw}}$  in the presence of a semi-honest adversary that can corrupt one party.*

**PROOF.** First, *Party A* receives two messages, which are  $\llbracket \nabla Z \rrbracket_B$  (Line 9) and  $\llbracket \nabla W_A - \phi \rrbracket_B$  (Line 12). We construct a simulator that

randomly picks plaintexts  $\nabla Z^*, (\nabla W_A - \phi)^*$  and encrypts them using  $pk_B$  to obtain  $\llbracket \nabla Z^* \rrbracket_B, \llbracket (\nabla W_A - \phi)^* \rrbracket_B$ . Since  $\nabla Z$  is the input and  $\nabla W_A - \phi$  is one piece of random secret,  $\nabla Z^*$  (or  $(\nabla W_A - \phi)^*$ ) and  $\nabla Z$  (or  $\nabla W_A - \phi$ ) share the same probability distribution. Furthermore, without *Party B*'s secret key  $sk_B$ , the ciphertexts  $\llbracket \nabla Z \rrbracket_B$  (or  $\llbracket \nabla W_A - \phi \rrbracket_B$ ) and  $\llbracket \nabla Z^* \rrbracket_B$  (or  $\llbracket (\nabla W_A - \phi)^* \rrbracket_B$ ) are computationally indistinguishable from the perspective of *Party A*. The other values in the view of *Party A* are computed locally. For instance, the simulator can simulate  $\llbracket W_A \rrbracket_B$  by computing  $\llbracket W_A^* \rrbracket_B = X_A^T \llbracket \nabla Z^* \rrbracket_B$ . Again, they are computationally indistinguishable from the perspective of *Party A*. Consequently, the view of *Party A* can be simulated perfectly.

Second, *Party B* only receives one messages by invoking the  $\Pi_{\text{HE2SS}}$  protocol (Line 10), whilst the other values are computed locally. Therefore, the view of *Party B* can be perfectly simulated by simulating the  $\mathcal{F}_{\text{HE2SS}}$  functionality as discussed in Lemma A.5.  $\square$

**THEOREM 5.2.** *Given  $Z, \nabla Z$  in the MatMul source layer, there are infinite possible values for  $X_A, W_A, X_A W_A$ .*

**PROOF.** First, we consider the linear equation  $Z = Z_A + Z_B$  where  $Z_A = X_A W_A, Z_B = X_B W_B$ . With only  $Z \in \mathbb{R}^{BS \times OUT}$  being known to *Party B*, there are  $BS \times OUT$  known values in total. However, since both  $Z_A, Z_B$  are unknown, which make up to be  $2 \times BS \times OUT$  variables in the equation. As a result, there must be infinite possible solution to the equation, which means there are infinite possible values for  $Z_A = X_A W_A$ . Then, for any possible  $X_A W_A$ , given an arbitrary invertible matrix  $M \in \mathbb{R}^{IN_A \times IN_A}$ , we have  $(X_A M^{-1})(M W_A) = X_A W_A$ . Consequently, there are infinite possible values for  $X_A, W_A$ , since both  $X_A W_A$  and  $M$  are arbitrary.

For  $\nabla Z$ , since it is computed locally on *Party B*, no extra information related to  $X_A, W_A, X_A W_A$  are provided. Thus, there are still infinite possible values for  $X_A, W_A, X_A W_A$ .  $\square$

### A.2 Proofs for Section 6.3

**THEOREM 6.1.** *The protocol of Embed-MatMul source layer securely realizes the ideal functionalities  $\mathcal{F}_{\text{EmbedMatMulFw}}$  and  $\mathcal{F}_{\text{EmbedMatMulBw}}$  in the presence of a semi-honest adversary that can corrupt one party.*

**PROOF.** We formally provide the definition and prove the security of the two idea functionalities in Lemma A.3 and A.4, respectively. Putting them together, we complete the proof for Theorem 6.1.  $\square$

Forward propagation of the Embed-MatMul source layer  $\mathcal{F}_{\text{EmbedMatMulFw}}$

**Inputs:**

- ▷ *Party A* inputs features  $X_A$ , secretly shared and/or encrypted models  $S_A, \llbracket T_A \rrbracket_B, T_B, U_A, \llbracket V_A \rrbracket_B, \llbracket U_B \rrbracket_B, V_B$ , and keys  $sk_A, pk_B$ ;
- ▷ *Party B* inputs features  $X_B$ , secretly shared and/or encrypted models  $S_B, \llbracket T_B \rrbracket_A, T_A, U_B, \llbracket V_B \rrbracket_A, \llbracket U_A \rrbracket_A, V_A$ , and keys  $pk_A, sk_B$ .

**Outputs:**

- ▷ *Party A* outputs nothing;
- ▷ *Party B* outputs  $Z = E_A W_A + E_B W_B$ , where  $E_\circ = \text{lkup}(Q_\circ, X_\circ), Q_\circ = S_\circ + T_\circ, W_\circ = U_\circ + V_\circ$ .

**LEMMA A.3.** *The protocol  $\Pi_{\text{EmbedMatMulFw}}$  in Line 5-11 of Figure 7 securely realizes  $\mathcal{F}_{\text{EmbedMatMulFw}}$  in the presence of a semi-honest adversary that can corrupt one party.*

PROOF. First, *Party A* receives three messages in total. In Line 6, *Party A* receives one message by the  $\Pi_{\text{HE2SS}}$  protocol. In Line 7-8, *Party A* invokes the  $\text{MatMulFw}$  routine twice, and each receives one message by the  $\Pi_{\text{HE2SS}}$  protocol, as discussed in our proof of Lemma A.1. Whilst all other values are computed locally. Obviously, the view of *Party A* can be perfectly simulated by simulating the  $\mathcal{F}_{\text{HE2SS}}$  functionality.

Second, *Party B* receives four messages in the protocol. Symmetrical as *Party A*, the first three messages are received via the  $\Pi_{\text{HE2SS}}$  protocol. We can simulate these three messages by simulating the  $\mathcal{F}_{\text{HE2SS}}$  functionality. The fourth message is received in Line 11. Denote the simulated versions of  $Z'_{1,B}$  and  $Z'_{2,B}$  are  $Z^*_{1,B}$  and  $Z^*_{2,B}$ , respectively. Then, we can simulate the fourth message, i.e.  $Z'_A$ , by computing  $Z^*_A = Z - Z^*_B$ , where  $Z^*_B = Z^*_{1,B} + Z^*_{2,B}$ . Since both  $Z'_A$  (or  $Z'_B$ ) and  $Z^*_A$  (or  $Z^*_B$ ) represent one piece of sharing of the output  $Z$ , they have the same probability distribution. As a result, we perfectly simulate the view of *Party B*.  $\square$

Backward propagation of the Embed-MatMul source layer  $\mathcal{F}_{\text{EmbedMatMulBw}}$

**Inputs:**

- ▷ *Party A* inputs features  $X_A$ , secretly shared activations  $\psi_A, E_B - \psi_B$ , secretly shared and/or encrypted models  $S_A, \llbracket T_A \rrbracket_B, T_B, U_A, \llbracket V_A \rrbracket_B, \llbracket U_B \rrbracket_B, V_B$ , and keys  $sk_A, pk_B$ ;
- ▷ *Party B* inputs features  $X_B$ , derivatives  $\nabla Z$ , secretly shared activations  $E_A - \psi_A, \psi_B$ , secretly shared and/or encrypted models  $S_B, \llbracket T_B \rrbracket_A, T_A, U_B, \llbracket V_B \rrbracket_A, \llbracket U_A \rrbracket_A, V_A$ , and keys  $pk_A, sk_B$ .

**Outputs:**

- ▷ *Party A* outputs  $\phi, \xi, \rho_A, \nabla Q_A - \rho_B$ ;
- ▷ *Party B* outputs  $\nabla W_A - \phi, \nabla W_B - \xi, \nabla Q_A - \rho_A, \rho_B$ , where  $\nabla W_\circ = E_\circ^T \nabla Z, \nabla E_\circ = \nabla Z W_\circ^T, \nabla Q_\circ = \text{lkup\_bw}(\nabla E_\circ, X_\circ)$ .

LEMMA A.4. *The protocol  $\Pi_{\text{EmbedMatMulBw}}$  in Line 12-26 of Figure 7 securely realizes  $\mathcal{F}_{\text{EmbedMatMulBw}}$  in the presence of a semi-honest adversary that can corrupt one party.*

PROOF. First, *Party A* receives six messages, which are  $\llbracket \nabla Z \rrbracket_B, \llbracket \nabla Z V_A^T \rrbracket_B$  (Line 12),  $\llbracket \nabla W_A - \phi \rrbracket_B$  (Line 18),  $\llbracket \nabla W_B - \xi \rrbracket_B$  (Line 19),  $\nabla Q_B - \rho_B$  (Line 23), and  $\llbracket \nabla Q_A - \rho_A \rrbracket_B$  (Line 25), respectively. For  $\nabla Q_B - \rho_B$ , which is received from the  $\Pi_{\text{HE2SS}}$  protocol, we can simulate this message by simulating the  $\mathcal{F}_{\text{HE2SS}}$  functionality as discussed in Lemma A.5. For the other five messages, we construct a simulator that randomly picks plaintexts  $\nabla Z^*, (\nabla Z V_A^T)^*, (\nabla W_A - \phi)^*, (\nabla W_B - \xi)^*, (\nabla Q_A - \rho_A)^*$  and encrypts them using  $pk_B$  to obtain  $\llbracket \nabla Z^* \rrbracket_B, \llbracket (\nabla Z V_A^T)^* \rrbracket_B, \llbracket (\nabla W_A - \phi)^* \rrbracket_B, \llbracket (\nabla W_B - \xi)^* \rrbracket_B, \llbracket (\nabla Q_A - \rho_A)^* \rrbracket_B$ . Since  $\nabla Z$  is the input and  $V_A, \nabla W_A - \phi, \nabla W_B - \xi, \nabla Q_A - \rho_A$  are random secrets,  $\nabla Z^*$  and  $\nabla Z$  (or  $(\nabla Z V_A^T)^*$  and  $\nabla Z V_A^T$ , or  $(\nabla W_A - \phi)^*$  and  $\nabla W_A - \phi$ , or  $(\nabla W_B - \xi)^*$  and  $\nabla W_B - \xi$ , or  $(\nabla Q_A - \rho_A)^*$  and  $\nabla Q_A - \rho_A$ ) share the same probability distribution. Furthermore, without *Party B*'s secret key  $sk_B$ , the ciphertexts  $\llbracket \nabla Z^* \rrbracket_B$  and  $\llbracket \nabla Z \rrbracket_B$  (or  $\llbracket (\nabla W_A - \phi)^* \rrbracket_B$  and  $\llbracket \nabla W_A - \phi \rrbracket_B$ , or  $\llbracket (\nabla W_B - \xi)^* \rrbracket_B$  and  $\llbracket \nabla W_B - \xi \rrbracket_B$ , or  $\llbracket (\nabla Q_A - \rho_A)^* \rrbracket_B$  and  $\llbracket \nabla Q_A - \rho_A \rrbracket_B$ ) are computationally indistinguishable from the perspective of *Party A*. The other values in the view of *Party A* are computed locally. For instance, the simulator can simulate  $\llbracket \nabla E_A \rrbracket_B$  by computing  $\llbracket \nabla E_A^* \rrbracket_B = \llbracket \nabla Z^* \rrbracket_B U_A^T + \llbracket (\nabla Z V_A^T)^* \rrbracket_B$ . Again, they are computationally indistinguishable from the perspective of *Party A*. Consequently, the view of *Party A* can be simulated perfectly.

For *Party B*, there are three messages received by invoking the  $\Pi_{\text{HE2SS}}$  protocol (Line 13, 15, and 22), and three messages that are ciphertexts of secrets (Line 17, 20, 26). For the first three messages, we can simulate them by simulating the  $\mathcal{F}_{\text{HE2SS}}$  functionality. For the last three messages, we simulate them in a similar way as simulating those of *Party A*, i.e., by randomly picking ciphertexts and encrypting them using  $pk_A$ . Again, since the messages are ciphertexts of random secrets, the simulation shares the same probability distribution as the original view and they are computationally indistinguishable from the perspective of *Party B*. As a result, the simulator perfectly simulates the view of *Party B*.  $\square$

THEOREM 5.2. *Given  $Z, \nabla Z$  in the Embed-MatMul source layer, there are infinite possible values for  $X_A, Q_A, W_A, E_A, E_A W_A$ .*

PROOF. First, we consider the linear equation  $Z = Z_A + Z_B$  where  $Z_A = E_A W_A, Z_B = E_B W_B$ . With only  $Z \in \mathbb{R}^{BS \times OUT}$  being known to *Party B*, there are  $BS \times OUT$  known values in total. However, since both  $Z_A, Z_B$  are unknown, which make up to be  $2 \times BS \times OUT$  variables in the equation. As a result, there must be infinite possible solution to the equation, which means there are infinite possible values for  $Z_A = X_A E_A$ . Then, for any possible  $X_A E_A$ , given an arbitrary invertible matrix  $M \in \mathbb{R}^{IN_A \times IN_A}$ , we have  $(E_A M^{-1})(M W_A) = E_A W_A$ . Consequently, there are infinite possible values for  $E_A, W_A$ , since both  $E_A W_A$  and  $M$  are arbitrary. Similarly, for any possible  $E_A$ , there are also infinite possible values for  $X_A, Q_A$ .

For  $\nabla Z$ , since it is computed locally on *Party B*, no extra information related to  $X_A, Q_A, W_A, E_A, E_A W_A$  are provided. Thus, there are still infinite possible values for  $X_A, Q_A, W_A, E_A, E_A W_A$ .  $\square$

### A.3 Proofs for Common Sub-Routines

Algorithm 1 and Algorithm 2 are two commonly used sub-routines in our work. In this subsection, we prove that the algorithm protocols securely realize the ideal functionalities as defined below.

Transformation from HE variables to SS variables  $\mathcal{F}_{\text{HE2SS}}$

**Inputs:**

- ▷ *Party  $\diamond$*  inputs a ciphertext  $\llbracket v \rrbracket_\diamond$  and the other party's public key  $pk_\diamond$ ;
- ▷ *Party  $\bar{\diamond}$*  inputs the secret key  $sk_\diamond$ .

**Outputs:**

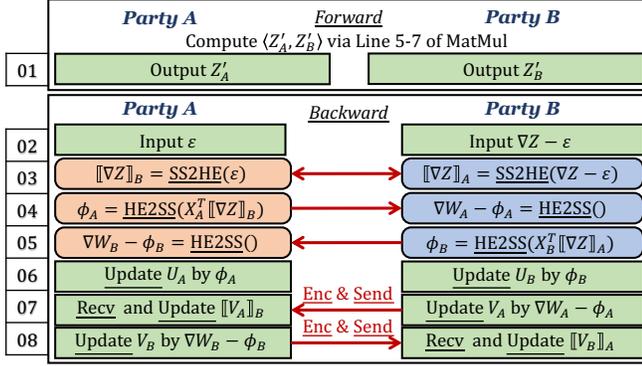
- ▷ *Party  $\diamond$*  outputs a random value  $\phi$ ;
- ▷ *Party  $\bar{\diamond}$*  outputs  $v - \phi$ .

LEMMA A.5. *The protocol  $\Pi_{\text{HE2SS}}$  in Algorithm 1 securely realizes  $\mathcal{F}_{\text{HE2SS}}$  in the presence of a semi-honest adversary that can corrupt one party.*

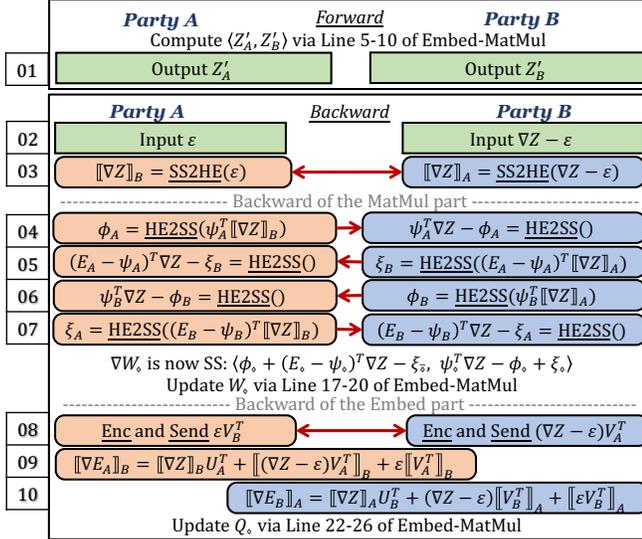
PROOF. Since *Party  $\diamond$*  receives no messages from *Party  $\bar{\diamond}$* , it is easy to know that our protocol is secure against the corruption of *Party  $\diamond$* . Hence, we only need to discuss how to simulate the view of *Party  $\bar{\diamond}$* . To do so, the simulator encrypts  $v - \phi$  with  $pk_B$  to obtain  $\llbracket v - \phi \rrbracket_\diamond^*$ . Here we differentiate  $\llbracket v - \phi \rrbracket_\diamond^*$  and the  $\llbracket v - \phi \rrbracket_\diamond$  in Algorithm 1 since they are from different times of encryption. Obviously, they share the same probability distribution and are computationally indistinguishable from the perspective of *Party  $\bar{\diamond}$* . As a result, it perfectly simulates the view of *Party  $\bar{\diamond}$* .  $\square$

**Algorithm 2:** The procedure to transform an SS variable  $\langle v_\diamond, v_\diamond \rangle$  into an HE variable  $\llbracket v \rrbracket$  where  $v = v_\diamond + v_\diamond$ .

- 1 **Function** SS2HE( $v_\diamond$ ):
- 2 Enc and Send the SS piece of this party  $v_\diamond$  using  $pk_\diamond$
- 3 Recv the encrypted SS piece of the other party  $\llbracket v_\diamond \rrbracket_\diamond$
- 4 **return**  $\llbracket v \rrbracket_\diamond = \llbracket v_\diamond \rrbracket_\diamond + v_\diamond$



**Figure 13:** The MatMul source layer followed by an SS-based top model. All cross-party transmission (red arrows) are protected by HE or SS.



**Figure 14:** The Embed-MatMul source layer followed by an SS-based top model. All cross-party transmission (red arrows) are protected by HE or SS.

Transformation from SS variables to HE variables  $\mathcal{F}_{\text{SS2HE}}$

**Inputs:**

- ▷ Party  $\diamond$  inputs a piece of SS  $v_\diamond$  and keys  $pk_\diamond, pk_\diamond$ ;
- ▷ Party  $\diamond$  inputs a piece of SS  $v_\diamond$  and keys  $pk_\diamond, pk_\diamond$ .

**Outputs:**

- ▷ Party  $\diamond$  outputs a ciphertext  $\llbracket v \rrbracket_\diamond$ ;
- ▷ Party  $\diamond$  outputs a ciphertext  $\llbracket v \rrbracket_\diamond$ , where  $v = v_\diamond + v_\diamond$ .

LEMMA A.6. The protocol  $\Pi_{\text{SS2HE}}$  in Algorithm 2 securely realizes  $\mathcal{F}_{\text{SS2HE}}$  in the presence of a semi-honest adversary that can corrupt one party.

PROOF. Since both parties have symmetric routines, we only describe how to simulate the view of Party  $\diamond$ . The only messages received in the protocol is  $\llbracket v_\diamond \rrbracket_\diamond$ . To simulate it, the simulator computes  $\llbracket v_\diamond \rrbracket_\diamond^* = \llbracket v \rrbracket_\diamond - v_\diamond$ . Here we differentiate  $\llbracket v_\diamond \rrbracket_\diamond^*$  and the  $\llbracket v_\diamond \rrbracket_\diamond$  in Algorithm 2 since they are from different times of encryption. Obviously, they have the same probability distribution since  $v = v_\diamond + v_\diamond$ , and they are computationally indistinguishable from the perspective of Party  $\diamond$ . As a result, the view of Party  $\diamond$  can be perfectly simulated.  $\square$

## B EXTENSION TO FEDERATED TOP MODELS

Although the top models are usually non-federated to address efficiency in practice, we can also use a federated top model so that Party B cannot get access to  $Z$  or  $\nabla Z$ , strengthening the security guarantees. In this section, we discuss how to adapt our work to the federated top model. Without loss of generality, we assume the top model utilizes the SS technique (e.g., SecureML).

First, for the forward propagation, our source layers should output an SS variable  $\langle Z'_A, Z'_B \rangle$  satisfying  $Z'_A + Z'_B = Z$ . In fact, this can be easily done with our algorithm protocols, as described below.

- ▷ For the MatMul source layer, we have the fact that  $Z'_\diamond = X_\diamond U_\diamond + \varepsilon_\diamond + X_\diamond V_\diamond - \varepsilon_\diamond$ , where  $\varepsilon_\diamond, \varepsilon_\diamond$  are generated by different parties. Hence,  $\langle Z'_A, Z'_B \rangle$  forms an SS variable of  $Z$  logically. As shown in Line 1 of Figure 13, Party A and Party B return  $Z'_A$  and  $Z'_B$ , respectively. By doing so, we can easily adapt the MatMul source layer to the federated top model.
- ▷ Similarly, for the Embed-MatMul source layer, we can also obtain the SS variable  $\langle Z'_A, Z'_B \rangle$  satisfying  $Z = Z'_A + Z'_B$  via Line 5-10 of Figure 7. Thus, we directly output them to the federated top model.

Second, for the backward propagation, Party B no longer gets access to  $\nabla Z$ . Instead, both parties take as input an SS variable  $\langle \varepsilon, \nabla Z - \varepsilon \rangle$ . For the MatMul source layer, the major difference is that Party B cannot compute model gradients  $\nabla W_B = X_B^T \nabla Z$  alone. To tackle this problem, we turn the SS variable into an HE variable  $\llbracket \nabla Z \rrbracket_A$  for Party B via Algorithm 2 (Line 3 of Figure 13). Next, Party B proceeds to compute the encrypted model gradients  $\llbracket \nabla W_B \rrbracket_A = X_B^T \llbracket \nabla Z \rrbracket_A$ , which are then transformed into an SS variable  $\langle \nabla W_B - \phi_B, \phi_B \rangle$  (Line 4-5 of Figure 13). Eventually, the update will be performed based on the secretly shared mode gradients. For the Embed-MatMul source layer, the routines are more complex. We refer interested readers to Figure 14 for more details.

To summarized, the federated source layers proposed in this work can also be integrated with federated top models.

### B.1 Security Analysis

In this subsection, we provide the security analysis of BLINDFL when the top model is also a federated module. Without loss of generality, we assume the model weights of the federated top model are also secretly shared and/or encrypted. Furthermore, since our work focuses on the security guarantees of federated source layers,

we assume the federated top model securely realizes the following ideal functionality for simplicity.

Forward and backward propagation of SS-based top model  $\mathcal{F}_{\text{TopSS}}$

**Inputs:**

- ▷ *Party A* inputs secretly shared activations  $Z'_A$ , secretly shared and/or encrypted models of the top model, and keys  $sk_A, pk_B$ ;
- ▷ *Party B* inputs inputs secretly shared activations  $Z'_B$ , labels  $y$ , secretly shared and/or encrypted models of the top model, and keys  $pk_A, sk_B$ .

**Outputs:**

- ▷ *Party A* outputs  $\varepsilon$ ;
- ▷ *Party B* outputs  $\nabla Z - \varepsilon$ , where  $\nabla Z$  denotes the backward derivatives of  $Z = Z'_A + Z'_B$ .

Then, we identify the ideal functionality of ML training of BLINDFL with an SS-based top model, and provide the security guarantees in Theorem B.1.

ML Training of BLINDFL with an SS-based top model  $\mathcal{F}_{\text{ML}}$

**Inputs:**

- ▷ *Party A* inputs features  $X_A$  and keys  $sk_A, pk_B$ ;
- ▷ *Party B* inputs features  $X_B$ , labels  $y$ , and keys  $pk_A, sk_B$ .

**Outputs:**

- ▷ Both parties output the trained federated source layers and *Party B* further outputs the trained federated top model (model weights are secretly shared or encrypted).

**THEOREM B.1.** *Assume the federated top model  $\Pi_{\text{TopSS}}$  securely realizes  $\mathcal{F}_{\text{TopSS}}$  in the presence of a semi-honest adversary that can corrupt one party. The training of BLINDFL using the MatMul and Embed-MatMul federated source layers securely realizes  $\mathcal{F}_{\text{ML}}$  in the presence of a semi-honest adversary that can corrupt one party.*

**PROOF.** Since all values are secretly shared and/or encrypted during the training process, it is actually a hybrid interaction of the  $\Pi_{\text{MatMulSSFw}}$ ,  $\Pi_{\text{MatMulSSBw}}$ ,  $\Pi_{\text{Embed-MatMulSSFw}}$ ,  $\Pi_{\text{Embed-MatMulSSBw}}$ , and  $\Pi_{\text{TopSS}}$  protocols. Therefore, simulation can be done using the hybrid argument. Combing Lemma B.2, B.3, B.4, B.5, and the assumption on the federated top model, we prove Theorem B.1.  $\square$

Forward propagation of the MatMul source layer followed by an SS-based top model  $\mathcal{F}_{\text{MatMulSSFw}}$

**Inputs:**

- ▷ *Party A* inputs features  $X_A$ , secretly shared and/or encrypted models  $U_A, \llbracket V_A \rrbracket_B$ , and keys  $sk_A, pk_B$ ;
- ▷ *Party B* inputs features  $X_B$ , secretly shared and/or encrypted models  $U_B, \llbracket V_B \rrbracket_A$ , and keys  $pk_A, sk_B$ .

**Outputs:**

- ▷ *Party A* outputs  $Z'_A$ ;
- ▷ *Party B* outputs  $Z'_B$ , where  $Z'_A + Z'_B = X_A W_A + X_B W_B$ ,  $W_\circ = U_\circ + V_\circ$ .

**LEMMA B.2.** *The protocol  $\Pi_{\text{MatMulSSFw}}$  in Line 1 of Figure 13 securely realizes  $\mathcal{F}_{\text{MatMulSSFw}}$  in the presence of a semi-honest adversary that can corrupt one party.*

**PROOF.** This lemma can be proved similarly as Lemma A.1. Thus, we would like to refer readers to the proof of Lemma A.1.  $\square$

Backward propagation of the MatMul source layer followed by an SS-based top model  $\mathcal{F}_{\text{MatMulSSBw}}$

**Inputs:**

- ▷ *Party A* inputs features  $X_A$ , secretly shared derivatives  $\varepsilon$ , secretly shared and/or encrypted models  $U_A, \llbracket V_A \rrbracket_B$ , and keys  $sk_A, pk_B$ ;
- ▷ *Party B* inputs features  $X_B$ , secretly shared derivatives  $\nabla Z - \varepsilon$ , secretly shared and/or encrypted models  $U_B, \llbracket V_B \rrbracket_A$ , and keys  $pk_A, sk_B$ .

**Outputs:**

- ▷ *Party A* outputs  $\phi_A, \nabla W_B - \phi_B$ ;
- ▷ *Party B* outputs  $\nabla W_A - \phi_A, \phi_B$ , where  $\nabla W_\circ = X_\circ^T \nabla Z$ .

**LEMMA B.3.** *The protocol  $\Pi_{\text{MatMulSSBw}}$  in Line 2-8 of Figure 13 securely realizes  $\mathcal{F}_{\text{MatMulSSBw}}$  in the presence of a semi-honest adversary that can corrupt one party.*

**PROOF.** Since both parties have symmetric routines, we only describe how to simulate the view of *Party A*. There are three messages received in the protocol. For the first two messages, i.e.,  $\llbracket \nabla Z \rrbracket_B, \nabla W_B - \phi_B$  (Line 3 and 5), they can be simulated by simulating the  $\Pi_{\text{SS2HE}}$  and  $\Pi_{\text{HE2SS}}$  protocols, respectively. For the third message, i.e.,  $\llbracket \nabla W_A - \phi_A \rrbracket_B$  (Line 7), we construct a simulator that randomly picks plaintexts  $(\nabla W_A - \phi_A)^*$  and encrypts them using  $pk_B$  to obtain  $\llbracket (\nabla W_A - \phi_A)^* \rrbracket_B$ . Since  $\nabla W_A - \phi_A$  is one piece of random secret,  $(\nabla W_A - \phi_A)^*$  and  $\nabla W_A - \phi_A$  share the same probability distribution. Furthermore, without *Party B*'s secret key  $sk_B$ ,  $\llbracket \nabla W_A - \phi_A \rrbracket_B$  and  $\llbracket (\nabla W_A - \phi_A)^* \rrbracket_B$  are computationally indistinguishable from the perspective of *Party A*. The other values in the view of *Party A* are computed locally. For instance, denoting  $\llbracket \nabla Z^* \rrbracket_B$  as the simulated version of  $\llbracket \nabla Z \rrbracket_B$ , the simulator can simulate  $\llbracket W_A \rrbracket_B$  by computing  $\llbracket W_A^* \rrbracket_B = X_A^T \llbracket \nabla Z^* \rrbracket_B$ . Again, they are computationally indistinguishable from the perspective of *Party A*. Consequently, the view of *Party A* can be simulated perfectly.  $\square$

Forward propagation of the Embed-MatMul source layer followed by an SS-based top model  $\mathcal{F}_{\text{EmbedMatMulSSFw}}$

**Inputs:**

- ▷ *Party A* inputs features  $X_A$ , secretly shared and/or encrypted models  $S_A, \llbracket T_A \rrbracket_B, T_B, U_A, \llbracket V_A \rrbracket_B, \llbracket U_B \rrbracket_B, V_B$ , and keys  $sk_A, pk_B$ ;
- ▷ *Party B* inputs features  $X_B$ , secretly shared and/or encrypted models  $S_B, \llbracket T_B \rrbracket_A, T_A, U_B, \llbracket V_B \rrbracket_A, \llbracket U_A \rrbracket_A, V_A$ , and keys  $pk_A, sk_B$ .

**Outputs:**

- ▷ *Party A* outputs  $Z'_A$ ;
- ▷ *Party B* outputs  $Z'_B$ , where  $Z'_A + Z'_B = E_A W_A + E_B W_B$ ,  $E_\circ = \text{lkup}(Q_\circ, X_\circ)$ ,  $Q_\circ = S_\circ + T_\circ$ ,  $W_\circ = U_\circ + V_\circ$ .

**LEMMA B.4.** *The protocol  $\Pi_{\text{EmbedMatMulSSFw}}$  in Line 1 of Figure 14 securely realizes  $\mathcal{F}_{\text{EmbedMatMulSSFw}}$  in the presence of a semi-honest adversary that can corrupt one party.*

**PROOF.** This lemma can be proved similarly as Lemma A.3. Thus, we would like to refer readers to the proof of Lemma A.3.  $\square$

Backward propagation of the Embed-MatMul source layer followed by an SS-based top model  $\mathcal{F}_{\text{EmbedMatMulSSBw}}$

**Inputs:**

- ▷ *Party A* inputs features  $X_A$ , secretly shared derivatives  $\varepsilon$ , secretly shared activations  $\psi_A, E_B - \psi_B$ , secretly shared and/or encrypted models  $S_A, \llbracket T_A \rrbracket_B, T_B, U_A, \llbracket V_A \rrbracket_B, \llbracket U_B \rrbracket_B, V_B$ , and keys  $sk_A, pk_B$ ;

**Algorithm 3:** The MatMul source layer under the multi-party setting. We assume there are  $M$  Party A's.

```

1 Function MultiPartyMatMulInit( $IN_{A(i)}, IN_B, OUT$ ):
2   if is Party B then
3     Initialize  $U_B \in \mathbb{R}^{IN_B \times OUT}$ 
4     foreach  $i \leftarrow 1$  to  $M$  do
5       Initialize  $V_{A(i)} \in \mathbb{R}^{IN_{A(i)} \times OUT}$ 
6       Enc and Send  $V_{A(i)}$ ; Recv  $\llbracket V_{B(i)} \rrbracket_{A(i)}$ 
7     return  $U_B, \llbracket V_{B(i)} \rrbracket_{A(i)}, V_{A(i)}$ 
8   else // Party A(i)
9     Initialize  $U_{A(i)} \in \mathbb{R}^{IN_{A(i)} \times OUT}, V_{B(i)} \in \mathbb{R}^{IN_B \times OUT}$ 
10    Enc and Send  $V_{B(i)}$ ; Recv  $\llbracket V_{A(i)} \rrbracket_B$ 
11    return  $U_{A(i)}, \llbracket V_{A(i)} \rrbracket_B, V_{B(i)}$ 
12 Function MultiPartyMatMulFw( $X_\circ, U_\circ, \llbracket V_\circ \rrbracket_\delta$ ):
13   if is Party B then
14     foreach  $i \leftarrow 1$  to  $M$  do
15       /* Line 5-8 of Figure 6 */
16        $Z_i = \text{MatMulFw}(X_B, U_B/M, \llbracket V_{B(i)} \rrbracket_{A(i)})$ 
17     return  $Z = \sum_i Z_i$ 
18   else // Party A(i)
19     /* Line 5-8 of Figure 6 */
20      $\text{MatMulFw}(X_A, U_{A(i)}, \llbracket V_{A(i)} \rrbracket_B)$ 
21     return null
22 Function MultiPartyMatMulBw( $X_\circ, U_\circ, V_\circ, \nabla Z$ ):
23   if is Party B then
24     Enc  $\nabla Z$  into  $\llbracket \nabla Z \rrbracket_B$ 
25     foreach  $i \leftarrow 1$  to  $M$  do
26       Send  $\llbracket \nabla Z \rrbracket_B$ 
27        $\nabla W_{A(i)} - \phi_{A(i)} = \text{HE2SS}(\nabla Z)$ 
28       Update  $V_{A(i)}$  via  $\nabla W_{A(i)} - \phi_{A(i)}$ , Enc and Send
29        $V_{A(i)}$ 
30     Update  $U_B$  via  $\nabla W_B = X_B^T \nabla Z$ 
31   else // Party A(i)
32     Recv  $\llbracket \nabla Z \rrbracket_B$ 
33      $\phi_{A(i)} = \text{HE2SS}(X_{A(i)}^T \llbracket \nabla Z \rrbracket_B)$ 
34     Recv  $\llbracket V_{A(i)} \rrbracket_B$ 

```

▷ Party B inputs features  $X_B$ , secretly shared derivatives  $\nabla Z - \varepsilon$ , secretly shared activations  $E_A - \psi_A, \psi_B$ , secretly shared and/or encrypted models  $S_B, \llbracket T_B \rrbracket_A, T_A, U_B, \llbracket V_B \rrbracket_A, \llbracket U_A \rrbracket_A, V_A$ , and keys  $pk_A, sk_B$ .

**Outputs:**

▷ Party A outputs  $\phi, \xi, \rho_A, \nabla Q_A - \rho_B$ ;  
 ▷ Party B outputs  $\nabla W_A - \phi, \nabla W_B - \xi, \nabla Q_A - \rho_A, \rho_B$ , where  $\nabla W_\circ = E_\circ^T \nabla Z, \nabla E_\circ = \nabla Z W_\circ^T, \nabla Q_\circ = \text{lkup\_bw}(\nabla E_\circ, X_\circ)$ .

**LEMMA B.5.** *The protocol  $\Pi_{\text{EmbedMatMulSSBw}}$  in Line 2-10 of Figure 7 securely realizes  $\mathcal{F}_{\text{EmbedMatMulSSBw}}$  in the presence of a semi-honest adversary that can corrupt one party.*

**PROOF.** Since both parties have symmetric routines, we only describe how to simulate the view of Party A. There are eight messages received in total.

First, for the four messages received in Line 3, 5, 6, and 8 of Figure 14, they are received by invoking the  $\Pi_{\text{HE2SS}}$  and  $\Pi_{\text{SS2HE}}$

protocols, so we can simulate them by simulating the  $\mathcal{F}_{\text{HE2SS}}$  and  $\mathcal{F}_{\text{SS2HE}}$  functionalities.

The  $\Pi_{\text{EmbedMatMulSSBw}}$  protocol in Figure 14 also involves parts of the  $\Pi_{\text{EmbedMatMulBw}}$  protocol in Figure 7, where four messages are also received (Line 18, 19, 23, and 25 of Figure 7). For the message received in Line 23 of Figure 7, we can simulate it by simulating the  $\mathcal{F}_{\text{HE2SS}}$  functionalities. For the other three messages received in Line 18, 19, and 25 of Figure 7, we construct a simulator that randomly picks plaintexts  $(\nabla W_A - \phi)^*, (\nabla W_B - \xi)^*, (\nabla Q_A - \rho_A)^*$  and encrypts them using  $pk_B$  to obtain  $\llbracket (\nabla W_A - \phi)^* \rrbracket_B, \llbracket (\nabla W_B - \xi)^* \rrbracket_B, \llbracket (\nabla Q_A - \rho_A)^* \rrbracket_B$ . Since  $\nabla W_A - \phi, \nabla W_B - \xi, \nabla Q_A - \rho_A$  are random secrets,  $(\nabla W_A - \phi)^*$  and  $\nabla W_A - \phi$  (or  $(\nabla W_B - \xi)^*$  and  $\nabla W_B - \xi$ , or  $(\nabla Q_A - \rho_A)^*$  and  $\nabla Q_A - \rho_A$ ) share the same probability distribution. Furthermore, without Party B's secret key  $sk_B$ , the ciphertexts  $\llbracket (\nabla W_A - \phi)^* \rrbracket_B$  and  $\llbracket \nabla W_A - \phi \rrbracket_B$  (or  $\llbracket (\nabla W_B - \xi)^* \rrbracket_B$  and  $\llbracket \nabla W_B - \xi \rrbracket_B$ , or  $\llbracket (\nabla Q_A - \rho_A)^* \rrbracket_B$  and  $\llbracket \nabla Q_A - \rho_A \rrbracket_B$ ) are computationally indistinguishable from the perspective of Party A. Therefore, they can be perfectly simulated.

The other values in the view of Party A are computed locally. For instance, denoting the simulated version of  $\llbracket \nabla Z \rrbracket_B, \llbracket (\nabla Z - \varepsilon) V_A^T \rrbracket_B$  as  $\llbracket \nabla Z \rrbracket_B^*, \llbracket (\nabla Z - \varepsilon) V_A^T \rrbracket_B^*$ , respectively, the simulator can simulate  $\llbracket E_A \rrbracket_B$  by computing  $\llbracket E_A \rrbracket_B^* = \llbracket \nabla Z \rrbracket_B^* U_A^T + \llbracket (\nabla Z - \varepsilon) V_A^T \rrbracket_B^* + \varepsilon \llbracket V_A^T \rrbracket_B$ . Again, they are computationally indistinguishable from the perspective of Party A. Consequently, the view of Party A can be simulated perfectly.  $\square$

## C EXTENSION TO MULTI-PARTY LEARNING

Although we focus on the two-party setting throughout this paper, the proposed federated source layers of BLINDFL can also be generalized to the multi-party setting where there are two or more Party A's. Here we provide an example of how to adapt the MatMul source layer to more than one Party A's. The adaption for Embed-MatMul can be achieved similarly.

The multi-party MatMul source layer is presented in Algorithm 3. The idea is to secretly share the model weights of each Party A with Party B one by one and let all Party A's execute the same routines. To be specific, for the  $i$ -th Party A (denoted as Party A(i)), the model weights are secretly shared as  $W_{A(i)} = U_{A(i)} + V_{A(i)}$ , where  $V_{A(i)}$  is managed by Party B. For Party B, the model weights are broken into more pieces, i.e.,  $W_B = U_B + \sum_i V_{B(i)}$ , where  $V_{B(i)}$  is managed by Party A(i).

Consequently, the goal of forward propagation is to compute the following results:

$$\begin{aligned}
 Z &= \sum_i X_{A(i)} W_{A(i)} + X_B W_B \\
 &= \sum_i X_{A(i)} (U_{A(i)} + V_{A(i)}) + X_B (U_B + \sum_i V_{B(i)}) \\
 &= \sum_i (X_{A(i)} U_{A(i)} + X_{A(i)} V_{A(i)} + X_B \frac{U_B}{M} + X_B V_{B(i)}),
 \end{aligned}$$

where  $M$  is the number of Party A's. This can be reckoned as performing federated matrix multiplication between Party B and each Party A(i), respectively, and summing up the  $M$  intermediate results to achieve the final results. The detailed routines are presented in Line 12-19 of Algorithm 3.

Table 6: Averaged training time cost of one mini-batch (in seconds). We only record the time cost of matrix multiplication for a fair comparison.

Dataset & Sparsity	Model	Time Cost/Batch (Seconds)		
		BLINDFL	SecureML	SecureML (Client-aided)
<i>fmnist</i> (Dense)	MLP	16.741	38.913	0.006

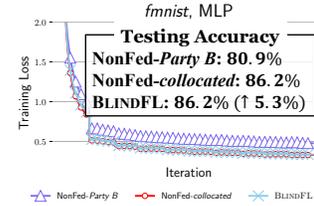


Figure 15: Training loss in terms of iterations and testing AUC or accuracy metrics.

Table 7: Scalability w.r.t. the output dimensionality of the MatMul source layer (*connect-4*, 3-layer MLP).

Hidden Dim	32	64	128	256
Related Time Cost	1.00×	1.91×	3.94×	8.06×
Validation Accuracy	79.8%	80.6%	81.8%	82.3%

For the backward propagation, *Party B* can still compute model gradients  $\nabla W_B = X_B^T \nabla Z$  on its own. For *Party A(i)*, it only needs to achieve the secretly shared model gradients  $\langle \phi_{A(i)}, \nabla W_{A(i)} - \phi_{A(i)} \rangle$  via a similar routine as Figure 6. Line 20-31 of Algorithm 3 shows the detailed procedure.

Finally, it is worthy to note that here we only provide a simple method to extend BLINDFL to the multi-party setting. It is interesting to explore how to optimize the algorithm protocols with smaller computation complexity and fewer communication rounds. We would like to leave them as our future work.

## D MORE EXPERIMENTS

### D.1 Experiments on Fashion MNIST

In Section 7, we conduct our experiments on tabular datasets. To evaluate whether our work also supports more kinds of datasets, we conduct experiments on the Fashion MNIST (*fmnist*) dataset in this section. To be specific, we split each image into two  $14 \times 28$  subfigures to simulate the feature partitioning in VFL. Then, we train the MLP models over the partitioned dataset. As shown in Table 6 and Figure 15, the results are consistent with those on other datasets. First, BLINDFL runs faster than SecureML without client aided but slower than the client-aided SecureML. Second, BLINDFL has a similar convergence as non-federated learning on collocated

Table 8: Scalability w.r.t. the number layers (*connect-4*, MLP).

# Layers	3	4	5	6
Related Time Cost	1.00×	1.01×	1.02×	1.02×
Validation Accuracy	80.6%	80.9%	81.0%	80.6%

data (NonFed-*collocated*), and achieves better performance than non-federated learning on the features of *Party B* only (NonFed-*Party B*), verifying the lossless property of our work.

### D.2 Scalability

In this section, we train the MLP models on the *connect-4* dataset to empirically evaluate the scalability of our work.

**Scalability w.r.t. #hidden dim.** We first fix the number of hidden layers as 3 and vary the output dimensionality of the first layer, i.e., the MatMul source layer. (The output dimensionalities of the second and third layers are fixed as 16 and 3, respectively.) As shown in Table 7, when the first layer is wider, the validation accuracy increases slightly, whilst the training time cost increases almost proportionally. This is because the number of cryptography operations in the MatMul source layer also increases correspondingly.

**Scalability w.r.t. #layers.** Then, we assess the effect of the number of layers. To do so, we fix the output dimensionalities of the first and the last but one layers as 64 and 16, and vary the number of layers by inserting 32-unit layers in the middle. For instance, in the 4-layer case, the output dimensionalities of all layers are 64, 32, 16, and 3, respectively. As shown in Table 8, the number of layers does not affect the training time cost significantly. This is reasonable since the top model is a non-federated module and the major overhead is the MatMul source layer.