

PackVFL: Efficient HE Packing for Vertical Federated Learning

Liu Yang

Hong Kong University of Science and
Technology
Hong Kong, China
lyangau@cse.ust.hk

Shuwei Cai

Hong Kong University of Science and
Technology (Guangzhou)
Guangzhou, China
scaiak@connect.hkust-gz.edu.cn

Di Chai

Hong Kong University of Science and
Technology
Hong Kong, China
dchai@cse.ust.hk

Junxue Zhang

Han Tian
Yilun Jin
Hong Kong University of Science and
Technology
Hong Kong, China
jzhangcs.htianab.yilun.jin@connect.ust.hk

Kun Guo

Fuzhou University
Fuzhou, China
gukn@fzu.edu.cn

Kai Chen

Qiang Yang
Hong Kong University of Science and
Technology
Hong Kong, China
kaichen.qyang@cse.ust.hk

Abstract

As an essential tool of secure distributed machine learning, vertical federated learning (VFL) based on homomorphic encryption (HE) suffers from severe efficiency problems due to data inflation and time-consuming operations. To this core, we propose PackVFL, an efficient VFL framework based on packed HE (PackedHE), to accelerate the existing HE-based VFL algorithms. PackVFL packs multiple cleartexts into one ciphertext and supports single-instruction-multiple-data (SIMD)-style parallelism. We focus on designing a high-performant matrix multiplication (MatMult) method since it takes up most of the ciphertext computation time in HE-based VFL. Besides, devising the MatMult method is also challenging for PackedHE because a slight difference in the packing way could predominantly affect its computation and communication costs. Without domain-specific design, directly applying SOTA MatMult methods is hard to achieve optimal.

Therefore, we make a three-fold design: 1) we systematically explore the current design space of MatMult and quantify the complexity of existing approaches to provide guidance; 2) we propose a hybrid MatMult method according to the unique characteristics of VFL; 3) we adaptively apply our hybrid method in representative VFL algorithms, leveraging distinctive algorithmic properties to further improve efficiency. As the batch size, feature dimension and model size of VFL scale up to large sizes, PackVFL consistently delivers enhanced performance. Empirically, PackVFL propels existing VFL algorithms to new heights, achieving up to a 51.52× end-to-end speedup. This represents a substantial 34.51× greater speedup compared to the direct application of SOTA MatMult methods.

PVLDB Reference Format:

Liu Yang, Shuwei Cai, Di Chai, Junxue Zhang, Han Tian, Yilun Jin, Kun Guo, Kai Chen, and Qiang Yang. PackVFL: Efficient HE Packing for Vertical Federated Learning. PVLDB, 14(1): XXX-XXX, 2020.
doi:XX.XX/XXX.XX

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.
Proceedings of the VLDB Endowment, Vol. 14, No. 1 ISSN 2150-8097.
doi:XX.XX/XXX.XX

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at [URL_TO_YOUR_ARTIFACTS](#).

1 Introduction

Vertical federated learning (VFL) [69, 70] based on homomorphic encryption (HE) [3] gradually becomes a trend of secure distributed machine learning among data silos [16, 25, 49, 58]. VFL solves the problem: multiple data owners, *e.g.*, companies or institutions, hold vertically partitioned data. They want to collaboratively train models without leaking the original data. During the process of HE-based VFL, raw data are maintained locally. Only intermediate results for calculating model updates are encrypted by HE [78] and exchanged among parties for further cryptographic computation. HE-based VFL is crucial for legally enriching ML features in both academia and industry.

However, state-of-the-art (SOTA) HE-based VFL algorithms [16, 18, 70, 77] suffer from severe efficiency issues, which hinder VFL’s wider application. On the one hand, their adopted HE methods, *e.g.*, Paillier [57], largely inflate data size to more than 40× [74], leading to communication and memory overheads. On the other hand, the adopted HE methods introduce time-consuming cryptographic operations, dominating the training process. We dig into these operations and find that matrix multiplication (MatMult) is the efficiency bottleneck. Taking the VFL-LinR algorithm [70] as an example, MatMult gradually occupies the majority of cryptographic computation time, up to 99.23%, shown in Tab. 1.

Our Contribution. In this paper, we propose PackVFL, an efficient VFL framework to accelerate the existing HE-based VFL algorithms. In detail, PackVFL is based on packed homomorphic encryption (PackedHE) [9, 10, 19, 24] and packs multiple cleartexts into one ciphertext to alleviate the data inflation problem and support parallel computation, *i.e.*, single-instruction-multiple-data (SIMD), over the packed values. More specifically, we focus on devising a high-performant MatMult method tailored for the VFL scenario. Empirically, we show the superiority of PackVFL,

Batch size	VFL-LinR [70]	Naively applying PackedHE
2	0.014s (10.14%)	0.037s (75.47%)
8	0.059s (29.29%)	0.227s (95.37%)
32	0.459s (65.19%)	0.834s (98.73%)
128	5.487s (87.66%)	4.079s (99.76%)
512	82.39s (96.71%)	20.73s (99.94%)
2048	1419s (99.23%)	111.1s (99.98%)

Table 1: Time consumed by cryptographic operations in one training batch. Numbers in parentheses represent the proportion of MatMult.

with 33.30 \times , 3.22 \times , 51.52 \times speedup over SOTA HE-based VFL algorithms, *i.e.*, VFL-LinR [70], CAESAR [16], VFL-NN [77] algorithms, respectively, in the end-to-end experiment.

1.1 Our Techniques

As the most time-consuming cryptographic operation in HE-based VFL, MatMult is also the design challenge of PackVFL. When PackedHE MatMult is adopted in the VFL scenario, the way of packing cleartexts primarily affects its computation and communication cost. Without domain-specific design, MatMult still seriously slows down the training process. For instance, we utilize the naive PackedHE MatMult method [27] in VFL-LinR and find it even slower than Paillier with small batch size, *e.g.*, 32, as shown in Tab. 1. The reason is the naive (row-order) method contains too many rotation operations [27], causing a significant computation overhead. We provide a detailed explanation in §3.

Besides, SOTA PackedHE MatMult methods proposed by [28, 31, 35, 50, 76] are also sub-optimal for the VFL scenario without domain specific designs. Among them, the methods of GAZELLE [35], DELPHI [50] and GALA [76] cause extra computation overheads, while the methods of Cheetah [31] and Iron [28] result in extensive communication overheads in VFL. Based on this fact, we raise the question: *Can we design a PackedHE MatMult method that is ideally tailored for the current HE-based VFL algorithms?* As we will show, the answer is yes. Our new MatMult method outperforms SOTA PackedHE MatMult methods theoretically in §4 and empirically in §7.

1.1.1 Systematical Exploration of Design Space. As one of the design emphases of PackedHE, the MatMult method has been continuously undergoing updates and iterations [27, 28, 31, 35, 50, 76]. Therefore, we provide a comprehensive analysis of the current design space and quantify the complexity of existing approaches to guide our design of PackVFL in §3. Based on the main idea of packing cleartexts, we divide them into slot packing methods [27, 35, 50, 60, 67, 76] and coefficient packing methods [28, 31]. Taking the MatMult between cleartext matrix X and ciphertext vector $\llbracket y \rrbracket$ as an example, we compare their methodologies. Slot packing methods follow the standard encoding procedure of PackedHE. It arranges which cleartexts of X are encoded together into a plaintext and their order. On the contrary, coefficient packing methods deviate from the standard procedure and directly map the cleartexts of X to specific coefficients of plaintext. The slot and coefficient packing methods have pros and cons considering computation and communication complexities. *Thus, the challenge lies in selecting the appropriate design path and further driving domain-specific innovations.*

1.1.2 MatMult Design for VFL Characteristics. To overcome the above challenge, we summarize three characteristics of the VFL’s required MatMult operation and design a hybrid MatMult method correspondingly in §4. The first characteristic is that the two operands of VFL MatMult are held by geo-distributed parties. One operand is encrypted by one party and transmitted to another for MatMult computation. After that, the resulting ciphertext is sent back. Through theoretical and empirical analyses of computation and communication costs, we opt for the slot packing concept. As a result, we propose *PackVFL’s diagonal method* as the core component of our hybrid approach, detailed in §4.1. The second characteristic is wide-range operand size. In VFL, the operand size of MatMult is related to batch size, feature dimension, and model architecture, with a wide range from small to large. Hence, we design delicate *input packing* and *input partitioning* techniques to further improve efficiency, correspondingly for the small- and large-operand situation in §4.2. The third characteristic is that the MatMult result is passively decrypted after transmission to the party with secret key. No extra operation is conducted. Due to this characteristic, we design the *lazy rotate-and-sum (RaS)* mechanism as the last component of our hybrid method to eliminate the remaining time-consuming ciphertext operation (rotation [27]) contained at the end of MatMult in §4.3.

Till now, our hybrid method has already shown its superiority over SOTA PackedHE MatMult methods, discussed in §4. For the computation comparison in §7.2, our hybrid method performs the best with the highest 846 \times and 1.24 \times speedup over the naive method [27] and the SOTA method of GALA [76], respectively.

1.1.3 Adaption Design to SOTA VFL Algorithms. To further improve the end-to-end training efficiency, we dig into the secure protocols of three representative HE-based VFL algorithms, *i.e.*, VFL-LinR [70], CAESAR [16], VFL-NN [77]¹ for illustration, and adaptively apply our MatMult method in them with three extra delicate mechanisms, leveraging their distinctive algorithmic properties, in §5. More detailedly, we design the *multiplication level reduction* mechanism, which modifies CAESAR’s original computation process to use more efficient PackedHE parameters in §5.2. To make the lazy RaS component of our hybrid method feasible with this modification, we also design the *cleartext inverse RaS* mechanism. In §5.3, focusing on VFL-NN, we encounter a new MatMult scenario between matrices $X\llbracket Y \rrbracket$. Rather than adhering to the conventional approach to diagonally encode the cleartext matrix X , we change our mind to diagonally encode the ciphertext matrix $\llbracket Y \rrbracket$, significantly enhancing efficiency. Moreover, we develop the *transposed matrices’ diagonal conversion* mechanism, enabling conversion between two encrypted transposed matrices, further halving the communication cost.

Together with the adaption optimization, PackVFL has more speedup over VFL-LinR [70], CAESAR [16] and VFL-NN [77] than the SOTA method of GALA by 9.02 \times , 0.91 \times and 33.1 \times respectively. More details are shown in §7.1. Besides, we also show that none of our designs harms the accuracy of VFL algorithms in §7.1.

¹Our focus on VFL-LinR, CAESAR, and VFL-NN is due to their foundational status and widespread recognition in HE-based VFL applications [69].

1.1.4 PackVFL’s Innovations and New Insights. As a cross-disciplinary effort bridging federated learning (FL) and cryptography, PackVFL makes significant contributions to both fields. For the FL community, PackVFL stands as one of the pioneering works to demonstrate the superiority of PackedHE over Paillier for VFL. We provide a counter-intuitive result that PackedHE is more suitable for VFL with our elaborate design. The VFL’s community has the intuition that PackedHE is much slower than Paillier since PackedHE is more complicated and supports more cryptographic operations. We are the first to devise an efficient MatMult method tailored to meet VFL’s specific requirements, showcasing its effectiveness across various VFL algorithms.

For the cryptography community, our proposed MatMult method exhibits potential advantages over state-of-the-art slot packing approaches, not only in VFL but also in other related domains such as secure model inference [48]. Additionally, we extend the MatMult scenario from matrix-vector $X[\mathbf{y}]$ to matrix-matrix $X[\mathbf{Y}]$ and design PackedHE techniques of the *cleartext inverse RaS* and *transposed matrices’ diagonal conversion*, which is unprecedented in earlier works.

1.2 Other Related Work

Several existing works [34, 46, 56, 60, 65, 67] also utilize PackedHE to construct federated algorithms. However, [56, 65] provide no design details about MatMult. [34, 46] designed multiplication between ciphertext vectors, which differs from our setting. [60, 67] involve SOTA PackedHE MatMult designs. But, their target MatMult operation is a single-party continuous MatMult operation without interaction, which is different from the requirements VFL. Therefore, they cannot be adopted in the VFL scenario. Besides, PackedHE hardware acceleration works, e.g., [59, 62, 75], can also be utilized to improve our performance further since PackVFL is a high-level application of PackedHE and changes no basic cryptographic operations.

2 Background and Preliminaries

2.1 Vertical Federated Learning

Vertical federated learning (VFL) [69, 70] takes a significant part of cross-silo distributed machine learning, whose participants are business companies [25, 49]. Companies usually maintain different portraits for common users, i.e., vertical data partition [16]. VFL enables secure model training over these abundant distributed features to achieve better prediction accuracy, complying with laws and regulations such as GDPR². Therefore, VFL holds substantial real-world importance. For example, [23] utilized VFL to help prevent COVID-19 [32].

2.1.1 HE-based VFL is a Widely-Adopted and Practical Solution. During the VFL training process, intermediate results instead of original data are exchanged among parties. To further protect intermediate results from disclosing privacy [78], VFL adopts various privacy-preserving methods, e.g., differential privacy (DP) [22], secret sharing (SS) [6], and homomorphic encryption (HE) [3]. *HE-based VFL [13, 16, 18, 70, 77] turns out to be a practical solution since*

²GDPR is a regulation in EU law on data protection and privacy in the European Union and the European Economic Area. <https://gdpr.eu/>.

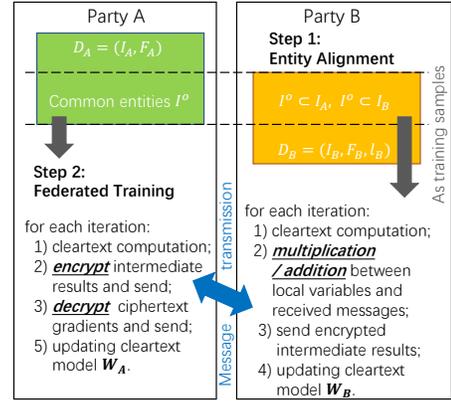


Figure 1: Illustration of HE-based VFL. Some VFL protocols may contain a trusted third party, which we omit to simplify.

DP-based VFL [66] and SS-based VFL [51, 52] suffer from either severe accuracy loss [69]³ or large communication overhead [16]. HE could provide the homomorphism between calculation over ciphertext and calculation over cleartext. HE-based VFL encrypts the transmitted messages, conducts computations on the received ciphertext, and decrypts the results to update the model.

Fig. 1 illustrates a two-party HE-based VFL scenario. Party A (P_A) owns a private dataset $\mathcal{D}_A = (\mathcal{I}_A, F_A)$, where \mathcal{I} stands for the sample identifiers and F_A represents the features. Party B (P_B) holds $\mathcal{D}_B = (\mathcal{I}_B, F_B, I_B)$, where I_B means labels. The first step of VFL is entity alignment, e.g., private set intersection (PSI) [41], to securely find the samples with common identifiers between two parties and use them as training data. The second step is federated training. At each iteration, P_A and P_B first conduct cleartext calculation over local data, then exchange encrypted ciphertext for multiplication, e.g., matrix multiplication (MatMult), addition, and decrypt gradients to update model weights W_A, W_B .

2.1.2 Paillier’s Batching is Limited for VFL. Paillier [3] is VFL’s most commonly adopted HE approach. It inflates data size and introduces time-consuming operations due to limited support for SIMD operations. *Although Paillier’s batching techniques [33, 74] tried to fill several cleartexts into one ciphertext, they do not support the MatMult operation.* As shown in Tab. 1, MatMult dominates the cryptographic operations of the VFL training process with up to 99.23% proportion. Therefore, Paillier’s batching techniques cannot be utilized in VFL. Besides, Paillier is based on the decisional composite residuosity assumption, which is vulnerable to post-quantum attacks [61]. Yet, Paillier remains the predominant HE technique in VFL.

2.2 Packed Homomorphic Encryption

Packed homomorphic encryption (PackedHE), e.g., BGV [10], BFV [24], CKKS [19], and multiparty BFV [53], naturally supports more functional and powerful SIMD-style batching than Paillier. As shown in

³DP is more adept at defending the adversarial attacks [39] and membership attacks [8], which are orthogonal to the privacy leakage during the training process.

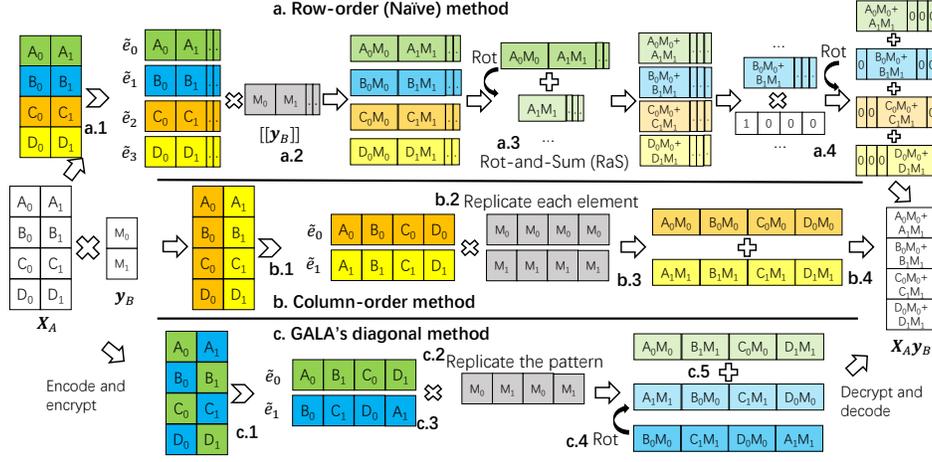


Figure 2: Illustration of slot packing methods, i.e., the row-order (naïve) method, column-order method, and our proposed generalized diagonal method, for PackedHE MatMult operation $X_A \llbracket y_B \rrbracket = \llbracket X_A y_B \rrbracket$. We set $m = N' = 4, n = 2$.

Fig. 1, PackedHE with the naïve MatMult method has already outperformed Paillier as the number of encrypted cleartexts increases. Moreover, PackedHE is based on a variant of learning with errors (LWE), i.e., ring LWE (RLWE) [47], which is quantum-resilient [55].

LWE encrypts a cleartext vector $x \in \mathbb{Z}_q^N$ into a ciphertext $\llbracket x \rrbracket = (c_0, c_1) = (x - As + e, A)$, where N stands for the vector length, q is the cleartext modulus (large prime number), $A \in \mathbb{Z}_q^{N \times N}$ is uniformly sampled, $s \in \mathbb{Z}_q^N$ represents the secret key, and $e \in \mathbb{Z}_q^N$ is small noise to make the problem hard. The decryption can only be conducted with $s: c_0 + c_1 s = x - As + e + As = x + e \approx x$. One disadvantage of LWE is that the size of matrix A is quadratic with the vector length N , which will cause large communication costs to transmit ciphertext.

PackedHE utilizes RLWE to solve this problem. For a cleartext vector $x \in \mathbb{C}_{N/2}$, before encryption, PackedHE encodes it to plaintext, i.e., an integer polynomial $\tilde{x} \in \mathbb{Z}_q[X]/(X^N + 1)$. For example, the standard encoding of CKKS utilizes cyclotomic polynomial to construct a one-to-one mapping between the cleartext vector and plaintext polynomial [19]. Then, PackedHE encrypts plaintext \tilde{x} to ciphertext $\llbracket x \rrbracket = (\tilde{c}_0, \tilde{c}_1) = (\tilde{x} - \tilde{a}s + \tilde{e}, \tilde{a})$. Values in the cleartext vector can be regarded to be encoded/encrypted in the corresponding position (slot) of plaintext/ciphertext with the same order. Slot number $N' = N/2$ refers to how many cleartexts can be encrypted in one RLWE ciphertext. The size of polynomial \tilde{a} is linear with N , which is much more efficient than LWE.

We conclude the basic RLWE operations that will be used in the following sections:

- **O1 (Add)**, which represents slot-wise addition, e.g., $\llbracket x \rrbracket + \llbracket y \rrbracket = \llbracket x + y \rrbracket$ or $x + \llbracket y \rrbracket = \llbracket x + y \rrbracket$;
- **O2 (Mult)**, which represents slot-wise multiplication, e.g., $\llbracket x \rrbracket \times \llbracket y \rrbracket = \llbracket x \times y \rrbracket$ or $x \times \llbracket y \rrbracket = \llbracket x \times y \rrbracket$;
- **O3 (Rot)**, which represents a rotation operation that shifts ciphertext slots in sequence. $\text{RotL/R}(\llbracket x \rrbracket, i)$ denotes rotating $\llbracket x \rrbracket$ to the left/right for i position. Rotation is the most costly basic

RLWE operation, which is more than 10× slower than **O1 (Add)** and **O2 (Mult)**;

- **O4 (HstRot)**, which represents hoisting rotation that is a more efficient optimization to conduct multiple rotation operations over the same ciphertext [35]. $\text{HstRotL/R}(\llbracket x \rrbracket, i)$ denotes i -position left/right HstRot of $\llbracket x \rrbracket$.

2.2.1 Designing PackedHE MatMult is Non-Trivial. PackedHE has already achieved good efficiency when conducting basic **O1 (Add)** and **O2 (Mult)** operations with SIMD property. However, there exist difficulties when designing a more complicated PackedHE MatMult method [35]. *The main reason is that PackedHE needs **O3 (Rot)** operations to sum values in different slots of ciphertexts. Since matrix multiplication needs a large number of vector inner sum operations, PackedHE MatMult is inclined to be inefficient without proper designs.*

3 Systematical Exploration of Design Space

The existing works [27, 28, 31, 35, 50, 76] have proposed SOTA designs for PackedHE MatMult. Taking the VFL's MatMult between matrix and vector as an example, we introduce and compare their methods to provide guidance for our design. The VFL's MatMult is formulated as:

$$X_A \llbracket y_B \rrbracket = \llbracket X_A y_B \rrbracket, \quad (1)$$

where $X_A \in \mathbb{R}^{m \times n}$ is held by P_A , $y_B \in \mathbb{R}^{n \times 1}$ is held by P_B . In VFL's process, P_B encrypts y_B as $\llbracket y_B \rrbracket$ and sends it to P_A for MatMult. The resulting ciphertext is transmitted back to P_B for decryption. X_A should also be encoded to plaintexts before MatMult.

We mainly divide the existing MatMult methods of [27, 28, 31, 35, 50, 76] into two categories below and will show their advantages and disadvantages in the following of this section:

- **Slot packing methods** that follow the standard encoding procedure of PackedHE, decide which cleartexts are encoded together into a plaintext and arrange their order [27, 35, 50, 76];

Category	Method	Computation Complexity				Communication complexity	
		# O1 (Add)	# O2 (Mult)	# O3 (Rot)	# O4 (HstRot)	P_B to P_A	P_A to P_B
Slot packing	Naive [27]	$m \log_2 n + m - 1$	$2m$	$m \log_2 n + m - 1$	0	1 RLWE-ct	1 RLWE-ct
	Column-order [27]	$n - 1$	n	0	0	n RLWE-ct	1 RLWE-ct
	GALA [76]'s Diagonal	$\min(m, n) - 1 + \log_2 \lceil \frac{n}{m} \rceil$	$\min(m, n)$	$\min(m, n) - 1$	$\log_2 \lceil \frac{n}{m} \rceil$	1 RLWE-ct	1 RLWE-ct
	PackVFL's Diagonal	$\min(m, n) - 1 + \log_2 \lceil \frac{n}{m} \rceil$	$\min(m, n)$	$\log_2 \lceil \frac{n}{m} \rceil$	$\min(m, n) - 1$	1 RLWE-ct	1 RLWE-ct
	GALA [76]	$\lceil \frac{mn}{N'} \rceil - 1$	$\lceil \frac{mn}{N'} \rceil$	$\lceil \frac{mn}{N'} \rceil - 1$	0	1 RLWE-ct	1 RLWE-ct
	PackVFL	$\lceil \frac{mn}{N'} \rceil - 1$	$\lceil \frac{mn}{N'} \rceil$	0	$\lceil \frac{mn}{N'} \rceil - 1$	1 RLWE-ct	1 RLWE-ct
Coefficient packing	Cheetah [31]	0	1	0	0	1 RLWE-ct	m LWE-ct

Table 2: Complexity of MatMult methods between matrix and vector. Computation complexity is mainly decided by the numbers (#) of involved O3 (Rot) and O4 (HstRot) operations. Communication complexity contains messages from P_B to P_A and from P_A to P_B . We assume N is large enough with no need for matrix/vector partition operation.

- **Coefficient packing methods** that deviate from the standard procedure and focus on designing novel mapping between clear-texts to specific coefficients of plaintext [28, 31].

In the following of this paper, we mention GAZELLE [35], DELPHI [50], GALA [76], Cheetah [31], and Iron [31] to only represent their MatMult methods for simplicity. Since their primary focus is on secure CNN inference, other technical contributions, *e.g.*, secure convolution techniques, are not applicable in our VFL scenario. Besides, we assume that m, n, N are both power of two. In the real-world scenario, m, n are not always the power of two, we can conduct the padding operation over X_A and \mathbf{y}_B with zero [7]. We set $m = 4, n = 2$ for demonstration. Matrix X_A has four rows, *i.e.*, $[A_0, A_1], [B_0, B_1], [C_0, C_1],$ and $[D_0, D_1]$, while column-vector \mathbf{y}_B contains two elements, *i.e.*, M_0 and M_1 . $A_i, B_i,$ and M_i are scalars to illustrate the computation process. Tab. 2 shows both the computation and communication complexity of the involved method. For a clear comparison, we assume N is large enough without needing matrix/vector partition.

3.1 Slot Packing

The fundamental idea of slot packing is deciding the encoding logic of the matrix operand [27]. Among the choices, row-order and column-order encodings are relatively intuitive. The slot packing methods are illustrated in Fig. 2 with $N' = 4$.

3.1.1 Row-Order (Naive) Method. The row-order (naive) method [27] is regarded as the most intuitive solution, which is frequently discussed as a baseline in related works [35, 50, 76]. And we also use the naive method to conduct the motivation experiments in Tab. 1. The efficiency of the naive method is largely slowed down by the required $m \log_2 n + m - 1$ **O3 (Rot)** operations, shown in Tab. 2. In our experiments, **O3 (Rot)** consumes around 90% time of the naive method when $m = n = 128$.

Shown in Fig. 2, the process of the naive method is: 1) each row of matrix X_A is encoded to \tilde{e}_i , respectively. For example, \tilde{e}_0 contains $[A_0, A_1, \cdot, \cdot]$, where we make omissions for vacant slots; 2) we multiply each \tilde{e}_i with ciphertext $\llbracket \mathbf{y}_B \rrbracket$; 3) since the values in one ciphertext cannot be directly summed up in PackedHE, we conduct $\log_2 n$ Rotation-and-Sum (RaS) for each resulting ciphertext, which shifts the ciphertext for some specific position, *i.e.*, **O3 (Rot)**, and adds the new ciphertext back to the old one. For instance, we conduct $\text{RotL}(\llbracket [A_0 M_0, A_1 M_1, \cdot, \cdot] \rrbracket, 1)$, obtain $\llbracket [A_1 M_1, A_0 M_0, \cdot, \cdot] \rrbracket$ and conduct $\llbracket [A_0 M_0, A_1 M_1, \cdot, \cdot] \rrbracket + \llbracket [A_1 M_1, A_0 M_0, \cdot, \cdot] \rrbracket$. After $m \log_2 n$ RaS containing $m \log_2 n$ **O3 (Rot)**, we get the result of each row's

dot product with $\llbracket \mathbf{y}_B \rrbracket$ at the first slot; 4) we want to obtain the expected result of MatMult in one ciphertext. To achieve this goal, we multiply each resulting ciphertext of the previous step with a cleartext indicator vector, which is only set to one at the first position and zero at other positions. Finally, we conduct extra $m - 1$ **O3 (Rot)** to further adjust the maintained values to appropriate slots and sum them together.

3.1.2 Column-Order Method. The column-order method [27] is less studied, even though owning a small computation complexity with no **O3 (Rot)** operations, shown in Tab. 2. However, it has a nearly $n \times$ communication complexity than naive method. Column-order method follows an opposite logic to the naive method. Fig. 2 shows its process: 1) each column of matrix X_A is encoded into \tilde{e}_i . \tilde{e}_0 contains $[A_0, B_0, C_0, D_0]$ and \tilde{e}_1 contains $[A_1, B_1, C_1, D_1]$; 2) the vector \mathbf{y} is encoded and encrypted to $n = 2$ ciphertexts. Each one separately replicates the corresponding element of \mathbf{y}_B , *i.e.*, $\llbracket [M_0, M_0, M_0, M_0] \rrbracket$ and $\llbracket [M_1, M_1, M_1, M_1] \rrbracket$. These n RLWE-ciphertexts are sent from P_B to P_A , causing large communication overhead. By contrast, the naive method only transmits one RLWE-ciphertext (ct); 3) slot-wise multiplication is conducted for each \tilde{e}_i and corresponding ciphertext from P_B ; 4) P_A performs slot-wise addition for results of the previous step and sends the obtained one RLWE-ct back to P_B .

3.1.3 GALA's Diagonal Method. Diagonal methods [35, 50, 76] are popular for slot packing. Their main idea is to place the values required to be added in the same slot of different ciphertexts to reduce the needed **O3 (Rot)**. DELPHI [50] is a GAZELLE [35] variant, assuming the input matrix of MatMult can be known in advance and moving part of the computation to the preprocessing phase. Since this assumption is not applicable in VFL, DELPHI degenerates to GAZELLE. GALA [76] outperforms GAZELLE by making the number of **O3 (Rot)** disproportional to N' and eliminating all **O4 (HstRot)** on $\llbracket \mathbf{y}_B \rrbracket$. Therefore, we illustrate the SOTA diagonal method of GALA.

Fig. 2 shows the process of GALA's diagonal method: 1) each diagonal of matrix X_A is encoded into \tilde{e}_i , following the diagonal order. \tilde{e}_0 contains $[A_0, B_1, C_0, D_1]$ and \tilde{e}_1 contains $[B_0, C_1, D_0, A_1]$; 2) the vector \mathbf{y}_B is encoded and encrypted, replicating the whole pattern, as $\llbracket [M_0, M_1, M_0, M_1] \rrbracket$. This RLWE-ct is sent from P_B to P_A ; 3) slot-wise multiplication is conducted by P_A for each \tilde{e}_i and the received ciphertext; 4) P_A rotates $\llbracket [B_0 M_0, C_1 M_1, D_0 M_0, A_1 M_1] \rrbracket$ to $\llbracket [A_1 M_1, B_0 M_0, C_1 M_1, D_0 M_0] \rrbracket$; 5) P_A conducts slot-wise addition for the result of step (4) and $\llbracket [A_0 M_0, B_1 M_1, C_0 M_0, D_1 M_1] \rrbracket$. The

obtained RLWE-ct is sent back to P_B for decryption. Shown in Tab. 2, GALA’s diagonal method maintains a small communication complexity same as the naive method as well as achieve a smaller computation complexity.

3.2 Coefficient Packing

Coefficient packing methods [28, 31] are recently proposed as SOTA methods, which do not follow previous standard encoding procedures, e.g., [10, 19, 24] and directly map cleartexts to coefficients of a plaintext polynomial. Cheetah [31] is proposed as a MatMult method between matrix and vector. Iron [28] extends it to MatMult between matrices. Their computation is only one multiplication, with no other operation, e.g., **O3 (Rot)**. However, they also suffer from high communication costs when applied in the VFL scenario, shown in Tab. 2. For illustration, we only introduce Cheetah [31], which suits our task, i.e., Eq. 1. The process of Cheetah is: 1) \tilde{X}_A is constructed according to mapping $\tilde{X}_A[i \cdot n + n - 1 - j] = X_A[i, j]$:

$$\begin{aligned} \tilde{X}_A = & A_1X^0 + A_0X^1 + B_1X^2 + B_0X^3 + \\ & C_1X^4 + C_0X^5 + D_1X^6 + D_0X^7 \in \mathbb{Z}_q[X]/(X^{16} + 1); \end{aligned} \quad (2)$$

2) \tilde{y}_B is constructed according to mapping $\tilde{y}_B[i] = \mathbf{y}[i]$:

$$\tilde{y}_B = M_0X^0 + M_1X^1 \in \mathbb{Z}_q[X]/(X^{16} + 1); \quad (3)$$

3) multiplication is conducted between \tilde{X}_A and \tilde{y}_B to obtain ⁴:

$$\begin{aligned} \tilde{X}_A\tilde{y}_B = & \dots + (A_0M_0 + A_1M_1)X^1 + \dots + (B_0M_0 + B_1M_1)X^3 + \\ & \dots + (C_0M_0 + C_1M_1)X^5 + \dots + (D_0M_0 + D_1M_1)X^7 + \\ & \dots \pmod{(X^{16} + 1, q)}; \end{aligned} \quad (4)$$

4) the MatMult results can be found with mapping $(X_A\mathbf{y}_B)[i] = (\tilde{X}_A\tilde{y}_B)[i \cdot n + n - 1]$. Cheetah extracts specific RLWE coefficients to separate m LWE-cts using [17] to protect the privacy of the other unneeded coefficients. In our example, four LWE-cts are extracted from coefficients $A_0M_0 + A_1M_1$, $B_0M_0 + B_1M_1$, $C_0M_0 + C_1M_1$, $D_0M_0 + D_1M_1$ and sent to P_B for decryption. As described in §2, we can calculate that these extracted LWE-cts are $m(1 + N)/2N \times$ larger than the original RLWE-ct, causing Large communication overheads.

4 MatMult Design for VFL Characteristics

With the guidance of comprehensive analysis in §3, we conclude characteristics of VFL’s MatMult operation and design a high-performant hybrid MatMult method accordingly, shown in Fig. 3. The characteristics of VFL’s MatMult are summarized below:

- **C1 (Geo-Distributed Operand)**: the two operands are owned by geo-distributed parties of VFL. The encrypted operand is sent to the other party for MatMult, and the ciphertext result is transmitted back for decryption;
- **C2 (Wide-Range Operand Size)**: as a machine learning technique, VFL often has varying batch sizes or feature dimensions, leading to varying operand sizes of MatMult, from large to small;
- **C3 (Passive Decryption)**: the resulting ciphertexts of MatMult are transmitted to the party owning secret key for pure decryption without extra operations.

⁴We only show the multiplication between polynomials for better illustration, which is same as [31]. In fact, \tilde{y}_B should be encrypted as RLWE-ct before multiplication.

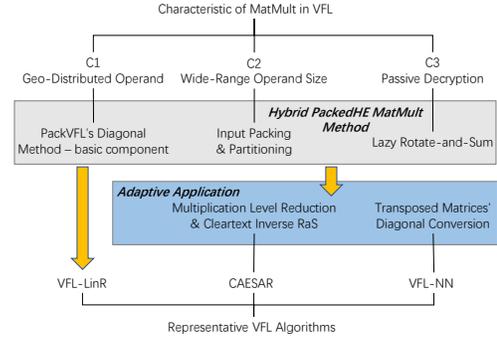


Figure 3: Overview of PackVFL. PackVFL contains two parts: 1) we design a hybrid MatMult method in terms of the characteristics of VFL; 2) we adaptively apply the proposed MatMult method to representative VFL algorithms.

4.1 PackVFL’s Diagonal Method

The **C1 (Geo-Distributed Operand)** characteristic indicates a wide area network (WAN) situation with relatively small bandwidth. Large communication overheads of MatMult could severely slow down the training process. Besides, parties (companies) of VFL usually contain rich computing resources. *Our design principle is: trying to reduce the computation complexity after guaranteeing that communication complexity is small enough.*

4.1.1 Our Choice of Diagonal Method. Considering the **C1 (Geo-Distributed Operand)** characteristic of VFL, we choose the diagonal method as our basic component. More specifically, comparing to the naive method, diagonal method can achieve less computation complexity while maintain the same communication complexity. For comparison with column-order and coefficient packing methods, we conduct an experiment where the bandwidth between P_A and P_B is 50MB/s [33] and $m = n = 512$, $N = 8192$. The transmission time of the column-order method and Cheetah is beyond 10s, which is much larger than the computation overhead (under 1s) of diagonal method. Therefore, we choose the most suitable design path for the VFL MatMult.

However, through the systematical analysis in §3, we observe that the SOTA diagonal method of GALA is not optimal for computation efficiency. The reason is that GALA still contains unnecessary **O3 (Rot)** operations caused by its diagonal encoding way. Therefore, we propose PackVFL’s diagonal method, which adopts a different diagonal encoding to further replace the unnecessary **O3 (Rot)** with more efficient **O4 (HstRot)** operations.

PackVFL’s diagonal method encodes the matrix X_A into:

$$\tilde{e}_i[j] = X_A[j \bmod m, (i + j) \bmod n], \quad (5)$$

where $i = 0, 1, \dots, \min(m, n) - 1$ and $j = 0, 1, \dots, \max(m, n) - 1$. This equation is a general solution for both $m \leq n$ (short-and-wide matrix) and $m > n$ (tall-and-skinny matrix) scenario. Fig. 4 illustrates the tall-and-skinny matrix case of $m > n$. The process is described as: 1) each \tilde{e}_i is generated along the diagonal, following Eq. 5. For example, $[A_0, B_1, C_0, D_1]$ are encoded into \tilde{e}_0 ; 2) the ciphertext $[\mathbf{y}_B]$ is constructed by replicating the whole \mathbf{y}_B for $\frac{m}{n}$

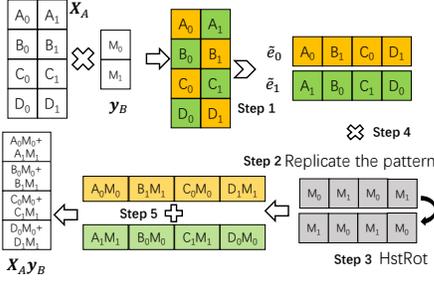


Figure 4: Illustration of PackVFL’s diagonal method.

times; 3) $n - 1$ **O3 (Rot)** need to be conducted on $[\mathbf{y}_B]$ to obtain $n - 1$ new ciphertexts. We use **O4 (HstRot)** to replace them for better efficiency; 4) multiplication between \tilde{e}_i and corresponding rotated ciphertext is performed; 5) we add the resulting ciphertexts to get $[\mathbf{X}_A \mathbf{y}_B]$. For the short-and-wide matrix case $m \leq n$, the construction of $[\mathbf{y}_B]$ needs no replication operation and additional RaS operations should be conducted. Readers could refer to Fig. 5a ($m = n$) and Fig. 5b ($m < n$) for more details.

Comparing to GALA in Tab. 2, PackVFL’s diagonal method substitutes $\min(m, n) - 1 - \log_2 \lceil \frac{n}{m} \rceil$ **O3 (Rot)** with the same number of **O4 (HstRot)** operations generally, achieving the optimal computation complexity. $\lceil \cdot \rceil$ stands for ceiling function. In the case of $m \geq n$, PackVFL’s diagonal method could substitute all **O3 (Rot)** with **O4 (HstRot)**, while, in the case of $m < n$, our diagonal method still needs $\log_2 \frac{n}{m}$ **O3 (Rot)** for the final RaS operations.

4.2 Input Packing and Partitioning

The characteristic **C2 (Wide-Range Operand Size)** indicates two scenarios: small size operands ($m, n < N'$) and large size operand ($m > N'$ or $n > N'$). The challenges of efficiently supporting small-and large-size operands are different. The former is to effectively leverage vacant slots for high-performant SIMD, while the latter one is to minimize redundant overheads across blocks as many as possible. We design input packing and partitioning, respectively.

4.2.1 Input Packing. For scenario $m, n < N'$, we pack multiple matrix diagonals into one plaintext polynomial to fully utilize the slots. In the beginning, we obtain the transformed matrix $X' \in \mathbb{R}^{m' \times N'}$ with N'/n sub-matrices vertically concatenated:

$$X' [i \bmod m', j + n \lfloor \frac{i}{m'} \rfloor] = X [i, j], \quad (6)$$

where $m' = \lfloor \frac{mn}{N'} \rfloor$, $\lfloor \cdot \rfloor$ stands for floor function, $i = 0, 1, \dots, m - 1$ and $j = 0, 1, \dots, n - 1$. Then, we conduct diagonal encoding for X' , which is slightly different from Eq. 5. We illustrate it in Fig. 5a to compute $X_A^T \llbracket d \rrbracket$. After we move the lower part of matrix X_A^T to the right and form a transformed matrix at step 2, the diagonal encoding is performed in parts. For example, at step 3, B_0 is placed after A_3 , not at the last slot of \tilde{e}_1 , while D_0 is set at the last slot. On the other side, the encrypted vector $\llbracket d \rrbracket$ is expanded by replicating its pattern as $\llbracket [M_0, M_1, M_2, M_3, M_0, M_1, M_2, M_3] \rrbracket$, at step 1.

4.2.2 Input Partitioning. For the scenarios $m > N'$ or $n > N'$, we split the original matrix/vector into multiple blocks in slot size N' and design tricks to further improve efficiency.

Large Operand Size is Common in VFL. VFL’s operand sizes are influenced by factors such as batch size, feature dimension, and the number of model parameters. These elements are increasingly scaling up in real-world applications. For instance, within the realm of machine learning, VFL is progressively adopting larger batch sizes — such as 15,000 [36], 4,096 [29], and 8,192 [42] — to enhance parallelism. Besides, the era of big data sees a surge in datasets with high feature dimensions, such as 22,283, 19,993 [43], reflecting a growing trend [73]. In addition, as GPT [11] gains widespread attention, large models, such as those with 65 billion [63] and 130 billion parameters [72], are becoming a focal point of research.

In the case $m > N' \geq n$, we horizontally split the original matrix X into $\frac{m}{N'}$ sub-matrices, then conduct the diagonal method between each sub-matrix and $[\mathbf{y}_B]$ separately. Besides, we only need to complete one group of **O4 (HstRot)** on $[\mathbf{y}_B]$ for all sub-matrices. For the case $m \leq N' < n$, we provide an example in Fig. 5b to compute $X_A \llbracket \langle w_A \rangle_2 \rrbracket$. $\langle \cdot \rangle$ stands for the secret share. At step 1&2, we vertically divide both the original matrix X_A and vector $\langle w_A \rangle_2$ into $\frac{n}{N'}$ blocks. At step 3, we conduct diagonal method for each block. Originally, we will obtain $\frac{n}{N'}$ ciphertext results and send them to P_B for decryption, which vastly increases communication complexity. However, we bring forward the addition operation of results in different blocks at step 4. Hence, we need only one transmission of the aggregated ciphertext.

4.3 Lazy Rotate-and-Sum.

Considering the **C3 (Passive Decryption)** characteristic, we design the Lazy Rotate-and-Sum (RaS) technique. The main idea is to eliminate the remaining ciphertext RaS operations, which contains **O3 (Rot)**, and replace them with cleartext sum-up computation. We provide illustrations in both Fig. 5a and Fig. 5b. Taking Fig. 5a as an example, after step 4, P_A obtains $\llbracket [A_0 M_0 + A_1 M_1, B_1 M_1 + B_2 M_2, A_2 M_2 + A_3 M_3, B_3 M_3 + B_0 M_0, \dots] \rrbracket$. Naively, a set of time-consuming RaS operations was necessary to aggregate intermediate elements within the ciphertext before sending it to P_B for decryption. However, since P_B performs no additional operations on the ciphertext beyond decryption and directly returns the results, the slots of the ciphertext that need decryption still align correctly with the indices of the resulting cleartext vectors. Consequently, we can defer the RaS operations after the decryption process. Ignoring the masking procedures at step 5&6, we will explain them in §5.1. P_A directly sends $\llbracket [A_0 M_0 + A_1 M_1, B_1 M_1 + B_2 M_2, A_2 M_2 + A_3 M_3, B_3 M_3 + B_0 M_0, \dots] \rrbracket$ to P_B and obtains $[A_0 M_0 + A_1 M_1, B_1 M_1 + B_2 M_2, A_2 M_2 + A_3 M_3, B_3 M_3 + B_0 M_0, \dots]$ after decryption. Finally, we could substitute the RaS operations with cleartext $A_0 M_0 + A_1 M_1 + A_2 M_2 + A_3 M_3, B_0 M_0 + B_1 M_1 + B_2 M_2 + B_3 M_3, C_0 M_0 + C_1 M_1 + C_2 M_2 + C_3 M_3, D_0 M_0 + D_1 M_1 + D_2 M_2 + D_3 M_3$.

Combining the above three components of our hybrid MatMult method, we can also compute the complexity of PackVFL for the cases of $m > N'$ or $n > N'$, shown in Tab. 3.

PackVFL vs. GALA. Besides the improvement of diagonal method explained in 4.1, PackVFL also has the following advantages over GALA: 1) PackVFL’s diagonal and input packing method are formally defined with equations. In contrast, GALA introduces its method only using legends and textual descriptions, which could lead to misinterpretation and poor generalizability; 2) GALA lacks mechanisms for input partitioning, whereas PackVFL incorporates

	$m > N' \geq n$	$m \leq N' < n$	$m, n > N'$
# O1 (Add)	$\frac{mn-m}{N'}$	$\frac{nm-N'}{N'}$	$\frac{mnN'-N'^2}{N'^2}$
# O2 (Mult)	$\frac{mn}{N'}$	$\frac{mn}{N'}$	$\frac{mn}{N'}$
# O3 (Rot)	0	0	0
# O4 (HstRot)	$n-1$	$\frac{mn-n}{N'}$	$\frac{nN'-n}{N'}$
P_B to P_A	1 RLWE-ct	$\frac{n}{N'}$ RLWE-ct	$\frac{n}{N'}$ RLWE-ct
P_A to P_B	$\frac{m}{N'}$ RLWE-ct	1 RLWE-ct	$\frac{m}{N'}$ RLWE-ct

Table 3: Complexity of PackVFL when $m > N'$ or $n > N'$.

sophisticated mechanisms that eliminate redundant operations when processing segmented blocks. Our experiments demonstrate a significant improvement in PackVFL when the operand size, e.g., feature dimension or batch size, is large in §7; 3) PackVFL extends the concept of lazy RaS in GALA beyond its original application from SS ciphertext to cleartext, thereby increasing its practical utility and relevance across diverse contexts. Additionally, in the VFL scenario, PackVFL introduces the innovative mechanism, i.e., cleartext inverse RaS, to facilitate the implementation of lazy RaS in §5.2.

5 Adaption Design to SOTA VFL Algorithms

In this section, we provide a comprehensive guidance on integrating the proposed PackedHE MatMult method by explaining how to apply it to representative VFL algorithms, e.g., VFL-LinR [70], CAESAR [16], and VFL-NN [77]. These three algorithms are fundamental and well-recognized by the VFL community [69]. Other algorithms, e.g., [12, 26, 30, 40, 45, 68, 71], which also contain MatMult operations, can similarly apply PackVFL. Accelerating HE-based VFL algorithms, e.g., [18], whose core computation is not MatMult, will be our future work. For CAESAR and VFL-NN, we analyze their unique algorithmic properties and made three adaption designs to improve efficiency further, as shown in Fig. 3:

- **Multiplication Level Reduction**, which optimizes CAESAR’s calculation process to reduce the number of accumulated multiplication over a single ciphertext for more efficient PackedHE parameters;
- **Cleartext Inverse Rotate-and-Sum (RaS)**, which is designed to make our MatMult method feasible for CAESAR after multiplication level reduction;
- **Transposed Matrices’ Diagonal Conversion**, which is a conversion mechanism between diagonal encodings of matrix and its transpose to reduce half of the communication cost in VFL-NN.

5.1 Adaption to VFL-LinR

VFL-LinR [70] securely trains a linear regression model, involving three parties. P_A, P_B hold X_A, X_B, \mathbf{y}_B and train weights θ_A, θ_B with the help of a third-party arbiter P_C . P_A, P_B exchange encrypted intermediate results $\llbracket \mathbf{u}_A \rrbracket, \llbracket \mathbf{d} \rrbracket$ and calculate over them to obtain the encrypted model updates $\llbracket \frac{\delta L}{\delta \theta_A} \rrbracket, \llbracket \frac{\delta L}{\delta \theta_B} \rrbracket$. Then, P_A, P_B add random masks on the updates and send them to P_C for decryption. Fig. 5a illustrates the computation process of P_A , and the operations of P_B could refer to Alg. 1.

Starting from step 2 in Tab. 1 of VFL-LinR [70], P_B receives the $\llbracket \mathbf{u}_A \rrbracket$ from P_A , computes the $\llbracket \mathbf{d} \rrbracket$ with the replicated encoding, and sends it to P_A . Then, P_A, P_B conduct vertical input packing over plaintext matrix X_A^T, X_B^T , and generate $\frac{mn}{N'}$ plaintext $\mathbf{e}_{A,i}, \mathbf{e}_{B,i}$. For

Algorithm 1: Modification to VFL-LinR Algorithm

-
- Input:** P_B receives $\llbracket \mathbf{u}_A \rrbracket$ from P_A (step 2 in Tab. 1 of VFL-LinR [70]);
- Output:** P_A, P_B obtain $\frac{\delta L}{\delta \theta_A}, \frac{\delta L}{\delta \theta_B}$, respectively (step 4 in Tab. 1 of VFL-LinR);
- P_B computes replicated $\llbracket \mathbf{d} \rrbracket = \llbracket \mathbf{u}_A \rrbracket + (\mathbf{u}_B - \mathbf{y})$ and sends it to P_A ;
 - P_A, P_B conduct vertical input packing over plaintext matrix X_A^T, X_B^T , respectively;
 - while** $i = 0, 1, \dots, \frac{mn}{N'} - 1$ **do**
 - P_A, P_B generate plaintext $\tilde{\mathbf{e}}_{A,i}, \tilde{\mathbf{e}}_{B,i}$ with diagonal encodings, respectively;
 - P_A, P_B conduct multiplication $\tilde{\mathbf{e}}_{A,i} \times \text{HstRotL}(\llbracket \mathbf{d} \rrbracket, i), \tilde{\mathbf{e}}_{B,i} \times \text{HstRotL}(\llbracket \mathbf{d} \rrbracket, i)$;
 - end**
 - P_A, P_B conduct addition over the above $\frac{mn}{N'}$ resulted ciphertexts separately;
 - P_A, P_B generate the random mask $\mathbf{r}_A, \mathbf{r}_B$ and add the mask on the result of the previous step;
 - P_A, P_B send the masked ciphertext to P_C ;
 - P_C decrypts the received messages and sends the cleartexts back to P_A, P_B , respectively;
 - P_A, P_B remove $\mathbf{r}_A, \mathbf{r}_B$ and apply lazy RaS to obtain the gradients $\frac{\delta L}{\delta \theta_A}, \frac{\delta L}{\delta \theta_B}$.
-

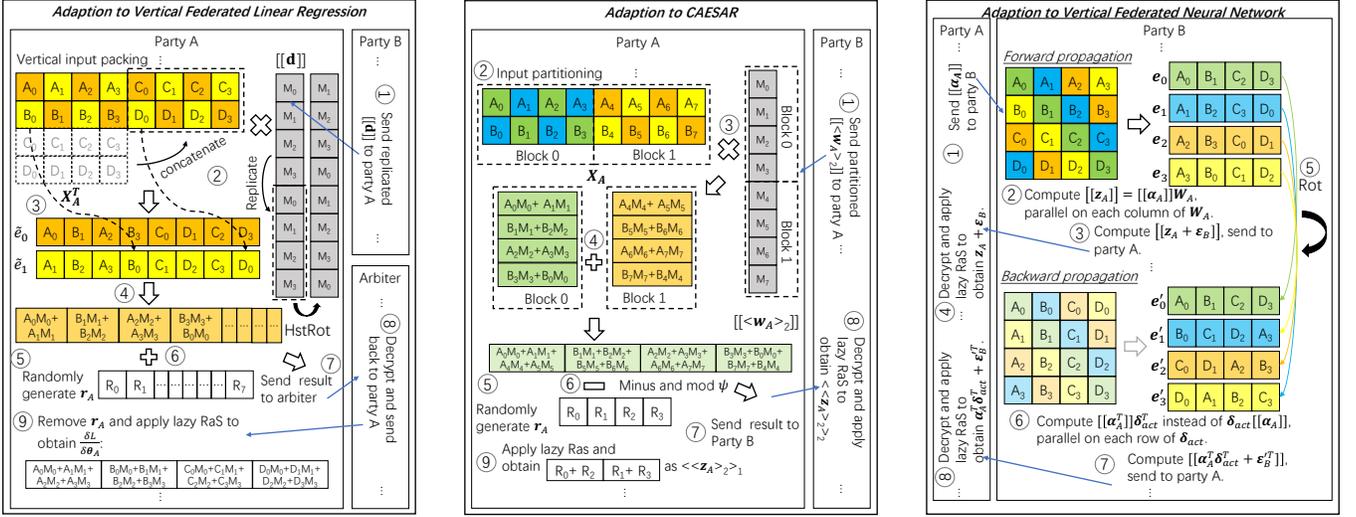
each plaintext, P_A, P_B conduct multiplication with corresponding the rotated ciphertext $\text{HstRotL}(\llbracket \mathbf{d} \rrbracket, i)$. Next, P_A, P_B separately sum up the results. After that, P_A, P_B send the masked ciphertexts to P_C . P_C decrypts them and sends the results back. P_A, P_B remove the masks and apply the lazy RaS to obtain gradients for model update. The process ends at step 4 in Tab. 1 of VFL-LinR.

5.2 Adaption to CAESAR

CAESAR [16] is a vertical federated logistic regression algorithm combining HE and SS to enhance the security, containing two parties. P_A, P_B distributively train on X_A, X_B, \mathbf{y}_B to obtain weights $\mathbf{w}_A, \mathbf{w}_B$. CAESAR is a SOTA end-to-end training algorithm. Our modifications are mainly located in Protocol 1 of CAESAR. Besides, we also make extra designs, i.e., **multiplication level reduction** and **plaintext inverse RaS**, to further improve efficiency.

Taking the prediction calculation process as an example, we show the modification to Protocol 1 of CAESAR. P_A, P_B initialize $\mathbf{w}_A, \mathbf{w}_B$, generate and exchange shares of model weights. More specifically, P_A generates shares $\langle \mathbf{w}_A \rangle_1, \langle \mathbf{w}_A \rangle_2$ ($\langle \mathbf{w}_A \rangle_1 + \langle \mathbf{w}_A \rangle_2 = \mathbf{w}_A$), holds $\langle \mathbf{w}_A \rangle_1$, and sends $\langle \mathbf{w}_A \rangle_2$ to P_B , where $\langle \cdot \rangle$ stands for a secret share. Next, P_A conducts $X_A \llbracket \langle \mathbf{w}_A \rangle_2 \rrbracket$ at line 11 in Algorithm 1 of CAESAR [16].

As shown in Fig. 5b and Alg. 2, at first, P_B divides $\langle \mathbf{w}_A \rangle_2$ into $\frac{m}{N'}$ blocks, encrypts, and sends them to P_A . P_A conducts input partitioning on matrix X_A and obtains blocks. For each vector and matrix block, P_A performs the MatMult method without the final RaS operations. Next, P_A conducts addition over the results of all blocks, randomly generates vector \mathbf{r}_A and subtracts it from the resulting ciphertext. Finally, P_A sends the outcome to P_B . P_B decrypts it and applies lazy RaS. P_A conducts cleartext RaS over \mathbf{r}_A .



(a) The VFL-LinR algorithm. We set $m = n = 4, N' = 8$. The input packing and lazy RaS techniques are also adopted.

(b) The CAESAR algorithm. We set $m = 8, n = 2, N' = 4$. The input partitioning and lazy RaS techniques are also adopted.

(c) The VFL-NN algorithm. We set $m = n = N' = 4$. The transposed matrices' diagonal conversion is also adopted.

Figure 5: Illustration of adapting PackVFL' MatMult method to three SOTA VFL algorithms.

Algorithm 2: Modification to CAESAR Algorithm

Input: P_A receives partitioned $[[w_A]_2]$ from P_B (line 1 in Protocol 1 of CAESAR [16]);

Output: P_A, P_B get $\langle\langle z_A \rangle_2\rangle_1, \langle\langle z_A \rangle_2\rangle_2$, respectively (output in Protocol 1 of CAESAR);

- 1 P_A conducts input partitioning on matrix X_A and obtains $\frac{m}{N'}$ matrix blocks;
- 2 **while** $j = 0, 1, \dots, \frac{m}{N'}$ **do**
- 3 P_A performs PackVFL-MatMult in block j without the final RaS operations;
- 4 **end**
- 5 P_A adds up the intermediate results of $\frac{m}{N'}$ blocks;
- 6 P_A randomly generates vector r_A ;
- 7 P_A subtracts the sum of intermediate results with r_A , following Protocol 2 in CAESAR;
- 8 P_A sends the subtraction result to P_B ;
- 9 P_B decrypts the received ciphertext and applies lazy RaS to obtain $\langle\langle z_A \rangle_2\rangle_2$;
- 10 P_A conducts cleartext RaS over r_A and gets $\langle\langle z_A \rangle_2\rangle_1$.

5.2.1 Multiplication Level Reduction. The efficiency of PackedHE is highly related to its parameters. Parameters allowing less multiplication refer to more efficient PackedHE operations. Therefore, we modify the computation process of CAESAR to reduce the multiplication level. More detailedly, in line 21 of the Algorithm 1 in CAESAR, $[[e]]^T = [[\hat{y}]] - y = q_0 + q_1[z] + q_2[z^3] - y$ is computed, which consumes one multiplication level, i.e., $q_1[z]$ and $q_2[z^3]$. Then, $[[g_B]] = [[e]]^T X_B$ is computed based on the resulted $[[e]]^T$, which needs one more multiplication level. We integrate them to:

$$\begin{aligned}
 [[g_B]] &= [[e]]^T X_B = ([[y]] - y)^T X_B = (q_0 + q_1[z] + q_2[z^3] - y)^T X_B \\
 &= (q_0 \vec{1} - y)^T X_B + (q_1 X_B)[z] + (q_2 X_B)[z^3],
 \end{aligned} \tag{7}$$

which decreases the multiplication level by one since the computation of $(q_1 X_B)[z]$ and $(q_2 X_B)[z^3]$ only need one multiplication level.

5.2.2 Cleartext Inverse RaS. However, Eq. 7 introduces extra cleartext addition term $(q_0 \vec{1} - y) \cdot X_B$ before secretly sharing $[[g_B]]$. We cannot directly perform lazy RaS over the shares because the addition term needs no inner sum operation. One option is to separately share the above cleartext addition term and the remaining ciphertext term, only conducting lazy RaS over the share of the ciphertext term. But this approach increases communication cost. Therefore, we design the cleartext inverse RaS method, which iteratively conducts split and concatenation over the cleartext plaintext term before addition, which properly aligns the number of required RaS between cleartext and ciphertext addition terms. For example, after one inverse RaS, cleartext vector $[M_0, M_1]$ turns to $[M_{0,0}, M_{1,0}, M_{0,1}, M_{1,1}]$, where $M_0 = M_{0,0} + M_{0,1}, M_1 = M_{1,0} + M_{1,1}$.

5.3 Adaption to VFL-NN

VFL-NN [77] implements a distributed neural network [2] via the SplitNN [64] architecture. It divides NN layers into P_A, P_B , separately owning X_A, X_B, y_B . VFL-NN, capable of handling more complex networks, addresses "harder" tasks. We also refer to the VFL-NN implementation in FATE [44]. In VFL-NN, P_A feeds data into the bottom layer and obtains embedding $\alpha_A \in \mathbb{R}^{m \times n}$, where m is the batch size, and n is the number of neurons in the bottom layer. Then, P_A encrypts α_A into $[[\alpha_A]]$ and sends them to P_B who

contains the interactive and top layers. Next, P_B conducts two separate MatMult $\llbracket \alpha_A \rrbracket \cdot W_A$ and $\delta_{\text{act}} \cdot \llbracket \alpha_A \rrbracket = (\llbracket \alpha_A^T \rrbracket \cdot \delta_{\text{act}}^T)^T$ in the forward and backward propagation calculation. We have two observations: 1) the MatMult is between matrices; 2) $\llbracket \alpha_A \rrbracket$ locates at opposite positions in two MatMult. For the first observation, we choose to apply diagonal encoding on the ciphertext matrix instead of the cleartext matrix to eliminate all the existing **O4 (HstRot)**, which is unprecedented in earlier works.

Algorithm 3: Modification to the Forward Propagation Process of VFL-NN Algorithm

Input: P_B receives diagonally encoded $\llbracket \alpha_A \rrbracket$ from P_A (line 5 in Alg. 1 of VFL-NN [77]);

Output: P_A obtains $z_A + \epsilon_B$ (line 8 in Alg. 1 of VFL-NN);

- 1 P_B randomly generates $\epsilon_B \in \mathbb{Z}^{m \times n'}$;
 - 2 **while** $k = 0, 1, \dots, n'$ **do**
 - 3 P_B conducts PackVFL-MatMult $\llbracket z_A[:, k] \rrbracket = \llbracket \alpha_A \rrbracket W_A[:, k]$;
 - 4 P_B computes $\llbracket z_A[:, k] \rrbracket + \epsilon_B[:, k]$;
 - 5 **end**
 - 6 P_B obtains $\llbracket z_A + \epsilon_B \rrbracket = \{\llbracket z_A[:, k] \rrbracket + \epsilon_B[:, k] \}_{k=0,1,\dots,n'}$ and sends it to P_A ;
 - 7 P_A decrypts the received ciphertexts and applies lazy RaS to obtain $z_A + \epsilon_B$.
-

5.3.1 *Transposed Matrices' Diagonal Conversion.* For the second observation, we designed the **transposed matrices' diagonal conversion** mechanism, which further reduces half of the communication cost. Therefore, instead of transmitting both $\llbracket \alpha_A \rrbracket$ and $\llbracket \alpha_A^T \rrbracket$ from P_A to P_B , we could only transmit $\llbracket \alpha_A \rrbracket$. The mechanism allows conversion between diagonal encodings \tilde{e} and \tilde{e}' of transposed matrices:

$$\tilde{e}'_i = \text{RotR}(\tilde{e}[(\min(m, n) - i) \bmod \min(m, n)], \max(m, n) - \min(m, n) + i). \quad (8)$$

Shown in Alg. 3 and Fig. 5c, our modification to forward propagation process starts from line 5 in Alg. 1 of VFL-NN [77]. First, P_A diagonally encodes, encrypts $\llbracket \alpha_A \rrbracket$, and sends it to P_B . P_B generates random mask $\epsilon_B \in \mathbb{R}^{m \times n'}$, where n' is the number of neurons in the interactive layer $W_A \in \mathbb{R}^{n \times n'}$. For each column k of W_A in parallel, P_B conducts MatMult between $\llbracket \alpha_A \rrbracket$ and W_A to obtain $\llbracket z_A[:, k] \rrbracket$, where $[:, k]$ represents the k column of matrix. The MatMult is efficient since we eliminate all **O4 (HstRot)**. Besides, P_B masks each resulted ciphertext. Next, P_B sends them to P_A for decryption and lazy RaS. In the idle period of waiting for P_A 's response, P_B conducts **transposed matrices' diagonal conversion** according to Eq. 8 and obtains $\llbracket \alpha_A^T \rrbracket$.

Shown in Alg. 4 and Fig. 5c, our modification to the backward propagation process starts from line 3 in Alg. 2 of VFL-NN [77]. First, P_B randomly generates mask $\epsilon'_B \in \mathbb{Z}^{n' \times n}$. Then, for each column k of δ_{act}^T , P_B conducts MatMult between $\llbracket \alpha_A^T \rrbracket$ and $\delta_{\text{act}}^T[:, k]$, where $\delta_{\text{act}} \in \mathbb{R}^{n' \times m}$ is the error back-propagated from the top layer. Besides, P_B conducts addition between $\llbracket \alpha_A^T \rrbracket + \delta_{\text{act}}^T[:, k]$ and $\epsilon_B'^T[:, k]$. Finally, P_B collects $\llbracket \alpha_A^T \delta_{\text{act}}^T + \epsilon_B'^T \rrbracket$ and sends it to P_A . P_A decrypts them, conducts lazy RaS and transpose to obtain $\delta_{\text{act}} \alpha_A + \epsilon_B$.

Algorithm 4: Modification to the Backward Propagation Process of VFL-NN Algorithm

Input: P_B obtains $\llbracket \alpha_A^T \rrbracket$ via Eq. 8 (line 3 in Alg. 2 of VFL-NN [77]);

Output: P_A obtains $\delta_{\text{act}} \alpha_A + \epsilon_B$ (line 7 in Alg. 2 of VFL-NN);

- 1 P_B randomly generates $\epsilon'_B \in \mathbb{Z}^{n' \times n}$;
 - 2 **while** $k = 0, 1, \dots, n'$ **do**
 - 3 P_B conducts PackVFL-MatMult
 $\llbracket \alpha_A^T \delta_{\text{act}}^T[:, k] \rrbracket = \llbracket \alpha_A^T \rrbracket \delta_{\text{act}}^T[:, k]$;
 - 4 P_B computes $\llbracket \alpha_A^T \delta_{\text{act}}^T[:, k] \rrbracket + \epsilon_B'^T[:, k]$;
 - 5 **end**
 - 6 P_B obtains
 $\llbracket \alpha_A^T \delta_{\text{act}}^T + \epsilon_B'^T \rrbracket = \{\llbracket \alpha_A^T \delta_{\text{act}}^T[:, k] \rrbracket + \epsilon_B'^T[:, k] \}_{k=0,1,\dots,n'}$ and sends it to P_A ;
 - 7 P_A decrypts the received ciphertexts, applies lazy RaS, and conducts transpose operation on them to obtain
 $\delta_{\text{act}} \alpha_A + \epsilon_B = (\alpha_A^T \delta_{\text{act}}^T + \epsilon_B'^T)^T$.
-

6 Correctness and Security Analysis

PackVFL innovates on the top of but didn't change basic PackedHE operations. Therefore, ciphertext correctness/security is fully protected by PackedHE, which is well recognized [35, 76]. Besides, PackVFL substitutes Paillier with PackedHE in HE-based VFL for cryptographic computation, which enhances the ciphertext security since PackedHE is quantum-resilient [55] and Paillier is vulnerable to post-quantum attacks [61].

We empirically show that our end-to-end training accuracy loss is small enough in §7.1. In addition, security of VFL protocols varies based on their originally designed exposure of cleartexts during training. For instance, VFL-LinR and VFL-NN are more susceptible to attacks than CAESAR due to greater information exposure. CAESAR exposes no intermediate cleartext results during training. As a plug-in method, PackVFL also didn't modify the original federated protocols [16, 70, 77]. Thus, their security analysis of VFL protocols, under various threat models, *e.g.*, semi-honest/malicious, remains valid and can be our reference.

7 Experiment

Setup and Implementation. The experiments are conducted on three x86 servers, each with 128GB memory and 40-core Intel Xeon Glod 5115 CPU. The bandwidth and latency among servers are 50MB/s and 20ms, respectively. We build our framework based on FATE [44] and Lattigo [1]. We compile the CKKS [19] interfaces of Lattigo with CGO⁵, implement our MatMult method over the compiled interfaces, and integrate it into VFL-LinR, CAESAR, and VFL-NN of FATE. The implementation of GALA refers to the open-source codes⁶.

Dataset. We utilize the SUSY dataset [5] to conduct the efficiency experiments, focusing on both computation and communication aspects. To meet different requirements of data size, we adjust the sample count and feature dimensions through sampling. Additionally, we use the NUS-WIDE dataset [20] for the accuracy experiments, since the NUS-WIDE dataset contains heterogeneous

⁵<https://golang.org/cmd/cgo/>

⁶<https://github.com/mc2-project/delphi>

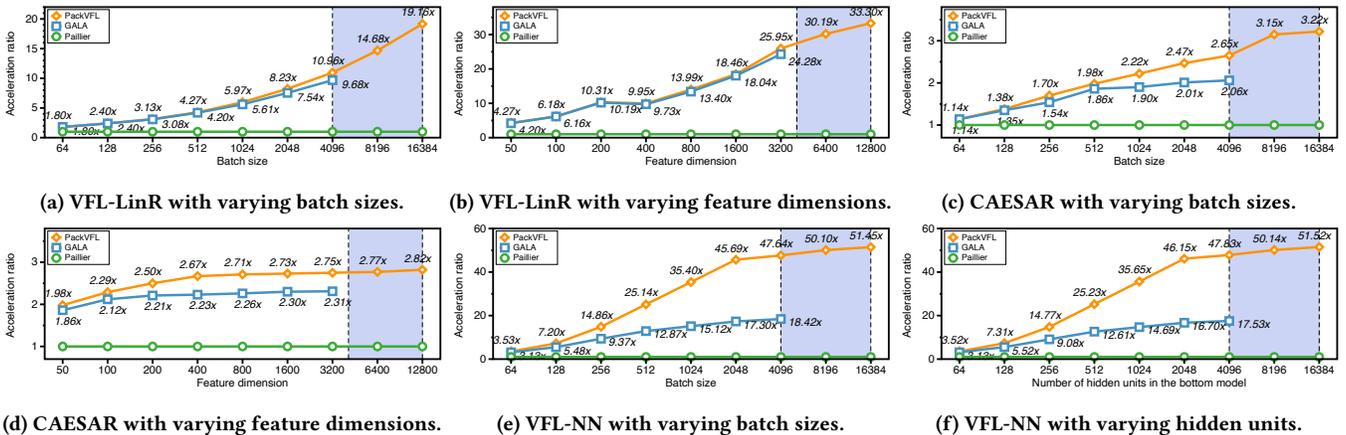


Figure 6: End-to-end acceleration of PackVFL and GALA for original Paillier-based VFL-LinR, CAESAR and VFL-NN algorithms. The ratio is computed using the time of each training epoch.

features, *i.e.*, text and image features, and is naturally apt for the VFL scenario ⁷. We separate its text and image features into two parties and create the vertically partitioned data situation. From the 81 groundtruth concepts, we select the most balanced one to serve as the label for our binary classification task.

Baseline. The baselines of the end-to-end experiment are utilizing Paillier and GALA [76] for MatMult computation in VFL-LinR, CAESAR, and VFL-NN, respectively. When batch size and feature dimension exceed $N' = N/2 = 4096$, we do not show the results of GALA because GALA have no specific design of input partitioning.

Selection of PackedHE Parameters. The parameters of PackedHE, *i.e.*, N , q , are pretty relative to the efficiency and security. q also decides the multiplication level. We set $N = 8192$ for all the experiments. We set $\log q = 156$ for the naive method and GALA in CAESAR, which allows two multiplication operations, and sets $\log q = 122$ for the other methods, which allows one multiplication operation. According to [4], the involved PackedHE MatMult methods offer a 192-bit security level. Therefore, we set the security level of the adopted Paillier as 128-bit [15] for a fair comparison.

7.1 End-to-End VFL

We show the end-to-end acceleration of GALA and PackVFL over Paillier for the VFL training process in Fig. 6. We could conclude that PackVFL always has the most significant acceleration ratio, which reflects speed gains over Paillier-based VFL algorithms, in different situations. Additionally, we demonstrate that PackVFL does not compromise the accuracy of the VFL training process in Tab. 4.

7.1.1 Accelerating VFL-LinR. In Fig. 6a and Fig. 6b, we vary the batch size and feature dimension, fixing the feature dimension and batch size as 50 and 512, respectively. As batch size and feature dimension increase, the speedups of GALA and PackVFL over Paillier also increase. PackVFL is more efficient than GALA from the beginning, with fewer rotation operations **O3 (Rot)** and **O4 (HstRot)**.

⁷We didn't use other datasets with homogeneous features, such as MNIST [38], CIFAR [37], and LEAF [14].

PackVFL achieves the largest acceleration ratio up to 33.30 \times over Paillier-based VFL-LinR when feature dimension extends to 12800 in Fig. 6b. As explained in §4.2, scenarios involving large feature dimensions are quite typical in VFL.

7.1.2 Accelerating CAESAR. In Fig. 6c and Fig. 6d, we also fix feature dimension and batch size as 50 and 512, separately. The speedups of GALA and PackVFL over Paillier in CAESAR are all smaller than those in VFL-LinR because CAESAR also contains many SS operations, which PackVFL cannot accelerate. Besides, the performance of GALA degrades more than PackVFL because they have no design for reducing multiplication level, thus leading to a less efficient PackedHE parameter $\log q = 156$. PackVFL still has the best performance from the beginning, with the largest speedup over Paillier up to 3.22 \times .

7.1.3 Accelerating VFL-NN. In Fig. 6e and Fig. 6f, we vary the batch size and number of hidden units in the bottom model, fixing the number of hidden units in the bottom model and batch size as 32 and 32, respectively. Besides, the feature dimension is 50, the number of hidden units in the interactive layer is 32, and the top model is linear. We show that PackVFL's acceleration for VFL-NN is much more significant than acceleration for both VFL-LinR and CAESAR. The reason is that PackVFL chooses to diagonally encode the ciphertext matrix, which eliminate all **O3 (Rot)** and **O4 (HstRot)**. Therefore, PackVFL has the most considerable speedup over Paillier up to 51.52 \times for VFL-NN. Besides, GALA is not explicitly designed for this scenario containing MatMult between matrices. They are much slower than PackVFL.

7.1.4 Communication Analysis. We also show the communication costs of different methods in Fig. 7, where the batch size is set to 512 for CAESAR and 32 for VFL-NN. The feature dimension is set to 800 for VFL-LinR and 50 for VFL-NN. The number of hidden units in the interactive layer is 32 for VFL-NN ⁸. The communication cost is computed as the sum of ciphertexts that a single party receives and sends in one training iteration.

⁸A smaller batch size of 32 was selected to avoid exceeding server memory limits.

Squared loss	VFL-LinR		Log loss	CAESAR		Log loss	VFL-NN	
	Paillier	PackVFL		Paillier	PackVFL		Paillier	PackVFL
epoch = 1	1.4375	1.4283 (-0.0092)	epoch = 1	0.5563	0.5526 (-0.0037)	epoch = 1	0.6442	0.6515 (+0.0073)
epoch = 5	0.6991	0.6918 (-0.0073)	epoch = 5	0.3287	0.3329 (+0.0042)	epoch = 5	0.3797	0.3708 (-0.0088)
epoch = 10	0.4643	0.4661 (+0.0018)	epoch = 10	0.2748	0.2691 (-0.0057)	epoch = 10	0.2144	0.2173 (+0.0029)
AUC	0.9025	0.9038 (+0.0013)	AUC	0.9652	0.9717 (+0.0065)	AUC	0.9815	0.9783 (-0.0032)

Table 4: Accuracy comparison between Paillier-based VFL algorithms and PackVFL. We configured identical hyperparameters for both, such as setting the learning rate to 0.001. We present an analysis that includes both the intermediate loss and the final Area Under the Curve (AUC) score on the training samples, with differences noted in parentheses.

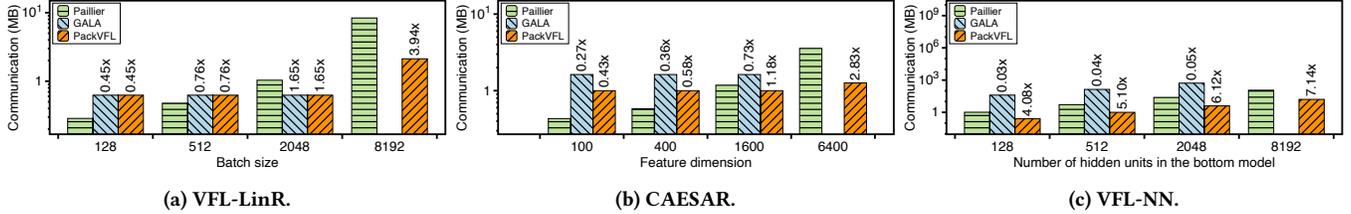


Figure 7: Communication comparison for applying PackVFL and GALA in VFL-LinR, CAESAR, and VFL-NN. The ratio is computed between PackedHE methods and Paillier over the three algorithms.

PackVFL generally has less communication cost than Paillier and GALA. More concretely, in Fig. 7a and Fig. 7b, with small batch sizes, the Paillier has less communication cost than PackVFL since Paillier implementation in FATE [44] has extra engineering optimizations. As batch size increases, PackVFL outperforms Paillier, with the largest improvement of 3.94 \times and 2.83 \times over VFL-LinR and CAESAR, respectively. In Fig. 7c, PackVFL has less communication cost than Paillier for VFL-NN, even when batch size is small. It has the largest improvement of 7.14 \times . The reason is the adoption of vertical input packing and transposed matrices’ diagonal conversion.

Besides, for VFL-LinR, PackVFL has the same communication cost as GALA for VFL-LinR, as shown in Fig. 7a, because their transmitted ciphertext amount and PackedHE parameters are both equal. For CAESAR, shown in Fig 7b, the communication costs of GALA are larger than PackVFL since their parameter $q = 156$ is larger, resulting in a larger ciphertext size. In VFL-NN, PackVFL has a larger advantage over GALA, shown in Fig. 7c. They suffer from extensive communication overhead since they still diagonally encode the cleartext matrix. The matrix to be encrypted can only be encrypted by each row/column with no input packing.

7.1.5 Model Accuracy Discussion. We demonstrate that PackVFL maintains consistent accuracy in Tab. 4. Our comparison of loss and AUC between Paillier-based VFL algorithms and PackVFL reveals minimal discrepancies, with differences not surpassing the third decimal place. This remarkable precision is attributable to the negligible noise inherent in the CKKS scheme, which could be effectively managed through rescaling operations [21] and aligns well with the generalization capabilities of machine learning models [54].

7.2 Computation Analysis for MatMult

In this part, we compare the computation efficiency of PackVFL, naive and GALA on the MatMult task between ciphertext vector and

Matrix size	Naive	GALA	PackVFL
512*64	14.3 (545 \times)	0.0304 (1.15 \times)	0.0265
512*256	17.3 (161 \times)	0.124 (1.16 \times)	0.107
512*1024	21.0 (46.6 \times)	0.541 (1.20 \times)	0.452
512*4096	24.7 (13.0 \times)	2.21 (1.16 \times)	1.90
64*512	2.33 (89.2 \times)	0.0299 (1.14 \times)	0.0261
256*512	9.91 (94.3 \times)	0.125 (1.19 \times)	0.105
1024*512	38.7 (84.8 \times)	0.530 (1.16 \times)	0.456
4096*512	155 (85.8 \times)	2.17 (1.20 \times)	1.81
64*64	1.73 (846 \times)	0.00204 (1.00 \times)	0.00204
256*256	8.78 (38.0 \times)	0.27 (1.17 \times)	0.231
1024*1024	41.9 (47.4 \times)	1.09 (1.24 \times)	0.883
4096*4096	198 (13.3 \times)	17.9 (1.20 \times)	14.9

Table 5: Computation cost comparison with SOTA methods for single MatMult operation. We show each method’s consuming time (s) and the speedup by PackVFL.

cleartext matrix. As shown in Tab. 5, we conduct experiments over different matrix sizes and provide the consuming time and speedup of PackVFL compared to the baselines. We obtain the following conclusions: 1) PackVFL and GALA are all more efficient than the naive method; 2) PackVFL is the most efficient over every matrix size; 3) with bigger m and n , PackVFL has a larger speedup over GALA, up to 1.24 \times . This experiment indicates that our MatMult method is also more efficient than GALA in communication-insensitive cases. We can extend it to other areas, e.g., secure model inference [48].

8 Conclusion

In this paper, we accelerate the main-stream Paillier-based VFL algorithms with PackedHE, which provides a counter-intuitive speedup. The intuition of VFL researchers is that PackedHE tends to be less efficient than Paillier since PackedHE is more complicated and supports more cryptographic operations. However, this paper shows that PackedHE is more suitable for VFL tasks considering both efficiency and security.

References

- [1] 2022. Lattigo v4. Online: <https://github.com/tuneinsight/lattigo>. EPFL-LDS, Tune Insight SA.
- [2] Oludare Isaac Abiodun, Aman Jantan, Abiodun Esther Omolara, Kemi Victoria Dada, Nachaat Abdelatif Mohamed, and Humaira Arshad. 2018. State-of-the-art in artificial neural network applications: A survey. *Heliyon* 4, 11 (2018), e00938.
- [3] Abbas Acar, Hidayet Aksu, A Selcuk Uluagac, and Mauro Conti. 2018. A survey on homomorphic encryption schemes: Theory and implementation. *ACM Computing Surveys (Csur)* 51, 4 (2018), 1–35.
- [4] Martin Albrecht, Melissa Chase, Hao Chen, Jintai Ding, Shafi Goldwasser, Sergey Gorbunov, Shai Halevi, Jeffrey Hoffstein, Kim Laine, Kristin Lauter, et al. 2021. Homomorphic encryption standard. In *Protecting Privacy through Homomorphic Encryption*. Springer, 31–62.
- [5] Pierre Baldi, Peter Sadowski, and Daniel Whiteson. 2014. Searching for exotic particles in high-energy physics with deep learning. *Nature communications* 5, 1 (2014), 1–9.
- [6] Amos Beimel. 2011. Secret-sharing schemes: A survey. In *International conference on coding and cryptography*. Springer, 11–46.
- [7] Ayoub Benaissa, Bilal Retiat, Bogdan Ceber, and Alaa Eddine Belfedhal. 2021. TenSEAL: A library for encrypted tensor operations using homomorphic encryption. *arXiv preprint arXiv:2104.03152* (2021).
- [8] Daniel Bernau, Jonas Robl, Philip W Grassal, Steffen Schneider, and Florian Kerschbaum. 2021. Comparing local and central differential privacy using membership inference attacks. In *IFIP Annual Conference on Data and Applications Security and Privacy*. Springer, 22–42.
- [9] Zvika Brakerski, Craig Gentry, and Shai Halevi. 2013. Packed ciphertexts in LWE-based homomorphic encryption. In *International Workshop on Public Key Cryptography*. Springer, 1–13.
- [10] Zvika Brakerski, Craig Gentry, and Vinod Vaikuntanathan. 2014. (Leveled) fully homomorphic encryption without bootstrapping. *ACM Transactions on Computation Theory (TOCT)* (2014).
- [11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [12] Dongqi Cai, Yan Kang, Lixin Fan, Mengwei Xu, Shanguang Wang, and Yang Qiang. 2022. Accelerating Vertical Federated Learning. *arXiv preprint arXiv:2207.11456* (2022).
- [13] Shuwei Cai, Di Chai, Liu Yang, Junxue Zhang, Yilun Jin, Leye Wang, Kun Guo, and Kai Chen. 2022. Secure Forward Aggregation for Vertical Federated Neural Networks. *arXiv preprint arXiv:2207.00165* (2022).
- [14] Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMahan, Virginia Smith, and Ameet Talwalkar. 2018. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097* (2018).
- [15] Dario Catalano, Rosario Gennaro, and Nick Howgrave-Graham. 2001. The bit security of Paillier’s encryption scheme and its applications. In *International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 229–243.
- [16] Chaochao Chen, Jun Zhou, Li Wang, Xibin Wu, Wenjing Fang, Jin Tan, Lei Wang, Alex X Liu, Hao Wang, and Cheng Hong. 2021. When homomorphic encryption marries secret sharing: Secure large-scale sparse logistic regression and applications in risk control. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*.
- [17] Hao Chen, Wei Dai, Miran Kim, and Yongsoo Song. 2021. Efficient homomorphic conversion between (ring) LWE ciphertexts. In *International Conference on Applied Cryptography and Network Security*. Springer, 460–479.
- [18] Kewei Cheng, Tao Fan, Yilun Jin, Yang Liu, Tianjian Chen, Dimitrios Papadopoulos, and Qiang Yang. 2021. Secureboost: A lossless federated learning framework. *IEEE Intelligent Systems* (2021).
- [19] Jung Hee Cheon, Andrey Kim, Miran Kim, and Yongsoo Song. 2017. Homomorphic encryption for arithmetic of approximate numbers. In *International conference on the theory and application of cryptology and information security*.
- [20] Tat-Seng Chua, Jinhui Tang, Richang Hong, Haojie Li, Zhiping Luo, and Yan-Tao Zheng. July 8–10, 2009. NUS-WIDE: A Real-World Web Image Database from National University of Singapore. In *Proc. of ACM Conf. on Image and Video Retrieval (CIVR’09)*. Santorini, Greece.
- [21] Roshan Dathathri, Blagovesta Kostova, Olli Saarikivi, Wei Dai, Kim Laine, and Madan Musuvathi. 2020. EVA: An encrypted vector arithmetic language and compiler for efficient homomorphic computation. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation*. 546–561.
- [22] Cynthia Dwork, Aaron Roth, et al. 2014. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science* 9, 3–4 (2014), 211–407.
- [23] Fatima Zahra Errounda and Yan Liu. 2022. A Mobility Forecasting Framework with Vertical Federated Learning. In *2022 IEEE 46th Annual Computers, Software, and Applications Conference (COMPSAC)*. IEEE, 301–310.
- [24] Junfeng Fan and Frederik Vercauteren. 2012. Somewhat practical fully homomorphic encryption. *Cryptology ePrint Archive* (2012).
- [25] Olga Fink, Torbjørn Netland, and Stefan Feuerriegel. 2021. Artificial intelligence across company borders. *Commun. ACM* 65, 1 (2021), 34–36.
- [26] 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*. 1316–1330.
- [27] Shai Halevi and Victor Shoup. 2014. Algorithms in helib. In *Annual Cryptology Conference*. Springer, 554–571.
- [28] Meng Hao, Hongwei Li, Hanxiao Chen, Pengzhi Xing, Guowen Xu, and Tianwei Zhang. 2022. Iron: Private inference on transformers. *Advances in Neural Information Processing Systems* 35 (2022), 15718–15731.
- [29] Elad Hoffer, Itay Hubara, and Daniel Soudry. 2017. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. *Advances in neural information processing systems* 30 (2017).
- [30] Yimin Huang, Xinyu Feng, Wanwan Wang, Hao He, Yukun Wang, and Ming Yao. 2022. EFMVFL: an efficient and flexible multi-party vertical federated learning without a third party. *arXiv preprint arXiv:2201.06244* (2022).
- [31] Zhicong Huang, Wen-jie Lu, Cheng Hong, and Jiansheng Ding. 2022. Cheetah: Lean and fast secure {two-party} deep neural network inference. In *31st USENIX Security Symposium (USENIX Security 22)*. 809–826.
- [32] Cornelia Ilin, Sébastien Annan-Phan, Xiao Hui Tai, Shikhar Mehra, Solomon Hsiang, and Joshua E Blumenstock. 2021. Public mobility data enables COVID-19 forecasting and management at local and global scales. *Scientific reports* 11, 1 (2021), 13531.
- [33] Zhifeng Jiang, Wei Wang, and Yang Liu. 2021. Flashe: Additively symmetric homomorphic encryption for cross-silo federated learning. *arXiv preprint arXiv:2109.00675* (2021).
- [34] Chao Jin, Jun Wang, Sin G Teo, Le Zhang, CS Chan, Qibin Hou, and Khin Mi Mi Aung. 2022. Towards End-to-End Secure and Efficient Federated Learning for XGBoost. (2022).
- [35] Chiraag Juvekar, Vinod Vaikuntanathan, and Anantha Chandrakasan. 2018. {GAZELLE}: A low latency framework for secure neural network inference. In *27th USENIX Security Symposium (USENIX Security 18)*.
- [36] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. 2016. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836* (2016).
- [37] Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images. (2009).
- [38] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324.
- [39] Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. 2019. Certified robustness to adversarial examples with differential privacy. In *2019 IEEE symposium on security and privacy (SP)*. IEEE, 656–672.
- [40] Ming Li, Yiwei Chen, Yiqin Wang, and Yu Pan. 2020. Efficient asynchronous vertical federated learning via gradient prediction and double-end sparse compression. In *2020 16th international conference on control, automation, robotics and vision (ICARCV)*. IEEE, 291–296.
- [41] Gang Liang and Sudarshan S Chawathe. 2004. Privacy-preserving inter-database operations. In *International Conference on Intelligence and Security Informatics*. Springer, 66–82.
- [42] Tao Lin, Lingjing Kong, Sebastian Stich, and Martin Jaggi. 2020. Extrapolation for large-batch training in deep learning. In *International Conference on Machine Learning*. PMLR, 6094–6104.
- [43] Bo Liu, Ying Wei, Yu Zhang, and Qiang Yang. 2017. Deep Neural Networks for High Dimension, Low Sample Size Data. In *IJCAI*. 2287–2293.
- [44] Yang Liu, Tao Fan, Tianjian Chen, Qian Xu, and Qiang Yang. 2021. Fate: An industrial grade platform for collaborative learning with data protection. *The Journal of Machine Learning Research* 22, 1 (2021), 10320–10325.
- [45] Yang Liu, Yan Kang, Xinwei Zhang, Liping Li, Yong Cheng, Tianjian Chen, Mingyi Hong, and Qiang Yang. 2019. A communication efficient collaborative learning framework for distributed features. *arXiv preprint arXiv:1912.11187* (2019).
- [46] Wen-jie Lu, Zhicong Huang, Qizhi Zhang, Yuchen Wang, and Cheng Hong. 2023. Squirrel: A Scalable Secure Two-Party Computation Framework for Training Gradient Boosting Decision Tree. *Cryptology ePrint Archive* (2023).
- [47] Vadim Lyubashevsky, Chris Peikert, and Oded Regev. 2013. On ideal lattices and learning with errors over rings. *Journal of the ACM (JACM)* 60, 6 (2013), 1–35.
- [48] Zoltán Ádám Mann, Christian Weinert, Daphnee Chabal, and Joppe W Bos. 2022. Towards Practical Secure Neural Network Inference: The Journey So Far and the Road Ahead. *Cryptology ePrint Archive* (2022).
- [49] H Brendan McMahan et al. 2021. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning* (2021).
- [50] Pratyush Mishra, Ryan Lehmkuhl, Akshayaram Srinivasan, Wenting Zheng, and Raluca Ada Popa. 2020. Delphi: A cryptographic inference service for neural networks. In *29th USENIX Security Symposium (USENIX Security 20)*. 2505–2522.

- [51] Payman Mohassel and Peter Rindal. 2018. ABY3: A mixed protocol framework for machine learning. In *Proceedings of the 2018 ACM SIGSAC conference on computer and communications security*. 35–52.
- [52] Payman Mohassel and Yupeng Zhang. 2017. Secureml: A system for scalable privacy-preserving machine learning. In *2017 IEEE symposium on security and privacy (SP)*. IEEE, 19–38.
- [53] Christian Mouchet, Juan Ramón Troncoso-Pastoriza, and Jean-Pierre Hubaux. 2020. Multiparty Homomorphic Encryption: From Theory to Practice. *IACR Cryptol. ePrint Arch.* 2020 (2020), 304.
- [54] Arvind Neelakantan, Luke Vilnis, Quoc V Le, Ilya Sutskever, Lukasz Kaiser, Karol Kurach, and James Martens. 2015. Adding gradient noise improves learning for very deep networks. *arXiv preprint arXiv:1511.06807* (2015).
- [55] Hamid Nejatollahi, Nikil Dutt, Sandip Ray, Francesco Regazzoni, Indranil Banerjee, and Rosario Cammarota. 2019. Post-quantum lattice-based cryptography implementations: A survey. *ACM Computing Surveys (CSUR)* 51, 6 (2019), 1–41.
- [56] Wei Ou, Jianhuan Zeng, Zijun Guo, Wanqin Yan, Dingwan Liu, and Stelios Fuentes. 2020. A homomorphic-encryption-based vertical federated learning scheme for rick management. *Computer Science and Information Systems* (2020).
- [57] Pascal Paillier. 1999. Public-key cryptosystems based on composite degree residuosity classes. In *International conference on the theory and applications of cryptographic techniques*. Springer, 223–238.
- [58] Bjarne Pfitzner, Nico Steckhan, and Bert Amrich. 2021. Federated Learning in a Medical Context: A Systematic Literature Review. *ACM Transactions on Internet Technology (TOIT)* (2021).
- [59] Nikola Samardzic, Axel Feldmann, Aleksandar Krastev, Srinivas Devadas, Ronald Dreslinski, Christopher Peikert, and Daniel Sanchez. 2021. F1: A fast and programmable accelerator for fully homomorphic encryption. In *MICRO-54: 54th Annual IEEE/ACM International Symposium on Microarchitecture*. 238–252.
- [60] Sinem Sav, Apostolos Pyrgelis, Juan Ramón Troncoso-Pastoriza, David Froelicher, Jean-Philippe Bossuat, Joao Sa Sousa, and Jean-Pierre Hubaux. 2021. POSEIDON: Privacy-Preserving Federated Neural Network Learning. In *28th Annual Network and Distributed System Security Symposium, NDSS 2021, virtually, February 21-25, 2021*.
- [61] Peter W Shor. 1999. Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer. *SIAM review* 41, 2 (1999), 303–332.
- [62] Yang Su, Bailong Yang, Chen Yang, and Luogeng Tian. 2020. FPGA-based hardware accelerator for leveled ring-lwe fully homomorphic encryption. *IEEE Access* 8 (2020), 168008–168025.
- [63] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).
- [64] Praneeth Vepakomma, Otkrist Gupta, Tristan Swedish, and Ramesh Raskar. 2018. Split learning for health: Distributed deep learning without sharing raw patient data. *arXiv preprint arXiv:1812.00564* (2018).
- [65] Hong Wang, Yuanzhi Zhou, Chi Zhang, Chen Peng, Mingxia Huang, Yi Liu, and Lintao Zhang. 2023. XFL: A High Performance, Lightweight Federated Learning Framework. *arXiv preprint arXiv:2302.05076* (2023).
- [66] Depeng Xu, Shuhan Yuan, and Xintao Wu. 2019. Achieving Differential Privacy in Vertically Partitioned Multiparty Learning. *arXiv preprint arXiv:1911.04587* (2019).
- [67] Guowen Xu, Xingshuo Han, Shengmin Xu, Tianwei Zhang, Hongwei Li, Xinyi Huang, and Robert H Deng. 2022. Hercules: Boosting the Performance of Privacy-preserving Federated Learning. *IEEE Transactions on Dependable and Secure Computing* (2022).
- [68] Kai Yang, Tao Fan, Tianjian Chen, Yuanming Shi, and Qiang Yang. 2019. A quasi-newton method based vertical federated learning framework for logistic regression. *arXiv preprint arXiv:1912.00513* (2019).
- [69] Liu Yang, Di Chai, Junxue Zhang, Yilun Jin, Leye Wang, Hao Liu, Han Tian, Qian Xu, and Kai Chen. 2023. A Survey on Vertical Federated Learning: From a Layered Perspective. *arXiv preprint arXiv:2304.01829* (2023).
- [70] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. 2019. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)* (2019).
- [71] Shengwen Yang, Bing Ren, Xuhui Zhou, and Liping Liu. 2019. Parallel distributed logistic regression for vertical federated learning without third-party coordinator. *arXiv preprint arXiv:1911.09824* (2019).
- [72] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414* (2022).
- [73] Yiteng Zhai, Yew-Soon Ong, and Ivor W Tsang. 2014. The emerging “big dimensionality”. *IEEE Computational Intelligence Magazine* 9, 3 (2014), 14–26.
- [74] Chengliang Zhang, Suyi Li, Junzhe Xia, Wei Wang, Feng Yan, and Yang Liu. 2020. {BatchCrypt}: Efficient homomorphic encryption for {Cross-Silo} federated learning. In *2020 USENIX annual technical conference (USENIX ATC 20)*. 493–506.
- [75] Junxue Zhang, Xiaodan Cheng, Liu Yang, Jinbin Hu, Ximeng Liu, and Kai Chen. 2022. Sok: Fully homomorphic encryption accelerators. *arXiv preprint arXiv:2212.01713* (2022).
- [76] Qiao Zhang, Chunsheng Xin, and Hongyi Wu. 2021. GALA: Greedy Computation for Linear Algebra in Privacy-Preserved Neural Networks. In *28th Annual Network and Distributed System Security Symposium, NDSS 2021, virtually, February 21-25, 2021*.
- [77] Yifei Zhang and Hao Zhu. 2020. Additively homomorphic encryption based deep neural network for asymmetrically collaborative machine learning. *arXiv preprint arXiv:2007.06849* (2020).
- [78] Ligeng Zhu, Zhijian Liu, and Song Han. 2019. Deep leakage from gradients. *Advances in neural information processing systems* (2019).