

Homomorphic Matrix Operations Under Bicyclic Encoding

Jingwei Chen[✉], Linhan Yang, Wenyuan Wu, Yang Liu, and Yong Feng

Abstract—Homomorphically encrypted matrix operations are extensively used in various privacy-preserving applications. Consequently, reducing the cost of encrypted matrix operations is a crucial topic on which numerous studies have been conducted. In this paper, we introduce a novel matrix encoding method, named *bicyclic encoding*, under which we propose two new algorithms **BMM-I** and **BMM-II** for encrypted matrix multiplication. **BMM-II** outperforms the state-of-the-art algorithms in theory, while **BMM-I**, combined with the *segmented* strategy, performs well in practice, particularly for matrices with high dimensions. Another noteworthy advantage of bicyclic encoding is that it allows for transposing an encrypted matrix entirely free. A comprehensive experimental study based on our proof-of-concept implementation shows that each algorithm introduced in this paper has specific scenarios outperforming existing algorithms, achieving speedups ranging from 2x to 38x.

Index Terms—Secure matrix multiplication, fully homomorphic encryption (FHE), Single Instruction Multiple Data (SIMD), bicyclic encoding.

I. INTRODUCTION

PRVACY-PRESERVING computing or privacy-enhancing technologies, capable of protecting data privacy and fully exploiting data value, is an exceptionally popular area of research. Fully Homomorphic Encryption (FHE) enables computations to be performed on encrypted data without the need for decryption [1], [2], and hence offering a powerful tool for privacy-preserving computation. Among the numerous privacy-preserving applications enabled by homomorphic encryption, matrix operations emerge as a core fundamental. Therefore, the importance of matrix multiplication over encrypted data is self-evident. In this paper, we investigate matrix multiplication over data encrypted using an FHE scheme that supports SIMD (Single Instruction

Multiple Data), such as the BGV [3] and B/FV [4], [5] schemes for integer arithmetic, and the CKKS scheme [6] for approximate arithmetic.

Since Gentry's pioneering work [2], FHE has rapidly developed, giving rise to various schemes such as BGV [3], B/FV [4], [5], CKKS [6], FHEW [7], TFHE [8], and numerous optimizations like data packing for SIMD [9], bootstrapping [10], [11], [12], [13], etc. However, homomorphic matrix multiplication remains a bottleneck in practice. For example, Huang et al. reported in [14] that it takes nearly 6 minutes for an encrypted matrix multiplication with dimensions 2048×8 and 8×2048 .

Almost all existing work about homomorphic matrix multiplication employed schemes that support SIMD, i.e., BGV, B/FV, or CKKS. For these schemes, computational efficiency is primarily influenced by two key factors. The first is *multiplicative depth* (including *ciphertext-ciphertext multiplication* (Mul) and *plaintext-ciphertext multiplication* (CMul)), a metric that directly impacts the computational efficiency of all known FHE schemes. This is because more multiplicative depths imply a larger ciphertext modulus or an increased number of bootstrappings. The second is the number of required *ciphertext rotations* (Rot) on the packed ciphertext. According to our test, for the CKKS scheme implemented in Microsoft SEAL [15], the efficiency ratio between a ciphertext rotation and a ciphertext multiplication can be as large as 8 : 1. Similar observations can be found in [16] as well. Consequently, the primary objective of this paper is to optimize both the multiplicative depth and the required number of ciphertext rotations for encrypted matrix multiplication.

A. Related Work

Secure matrix multiplication has plenty of applications, making it a highly active and impactful research area in privacy-preserving computation. The subject offers diverse perspectives for investigation, including secure multi-party computation [17], [18], [19], [20], [21], information-theoretically privacy [22], etc. Here, we mainly focus on encrypted matrix multiplication methods based on FHE schemes supporting SIMD, although there exist many other algorithms, e.g., [23] and [24], based on a matrix version of GSW [25].

A plaintext of an FHE scheme that supports SIMD is usually an element in a certain polynomial ring, say $R = \mathbb{Z}[X]/\langle X^N + 1 \rangle$. Polynomials in R can be represented in two equivalent ways. One is referred to *coefficients-encoding*, corresponding

Received 18 April 2024; revised 22 September 2024; accepted 22 October 2024. Date of publication 4 November 2024; date of current version 27 January 2025. This work was supported in part by the National Key Research and Development Program of China under Grant 2020YFA0712303; in part by the Natural Science Foundation of Chongqing under Grant 2022yszx-jcx0011cstb, Grant cstb2023yszx-jcx0008, and Grant cstb2023nseq-msx0441; and in part by the Western Young Scholars Program of Chinese Academy of Sciences (CAS). The associate editor coordinating the review of this article and approving it for publication was Dr. Jun Zhao. (Corresponding author: Linhan Yang.)

Jingwei Chen, Wenyuan Wu, and Yong Feng are with the Chongqing Key Laboratory of Secure Computing for Biology, Chongqing Institute of Green and Intelligent Technology, CAS, Chongqing 400714, China, and also with Chongqing College, University of CAS, Chongqing 400714, China (e-mail: chenjingwei@cigit.ac.cn; wuwenyuan@cigit.ac.cn; yongfeng@cigit.ac.cn).

Linhan Yang and Yang Liu are with the School of Information Science and Engineering, Chongqing Jiaotong University, Chongqing 400074, China (e-mail: linhanyang@mails.cqjtu.edu.cn; liuyang13@cqjtu.edu.cn).

Digital Object Identifier 10.1109/TIFS.2024.3490862

TABLE I

THE COST OF SLOTS-ENCODING ALGORITHMS FOR A (n, m, p) MATRIX MULTIPLICATION OVER ENCRYPTED DATA, WHERE THE TWO MATRICES TO BE MULTIPLIED ARE OF DIMENSIONS $n \times m$ AND $m \times p$, RESPECTIVELY, AND ℓ IS THE NUMBER OF PLAINTEXT SLOTS

Method	Max. dim.	#Ctxts	#Mul	#CMul	#Rot	#Mult. depth
[10], [32]*	ℓ	(d, p)	pd	0	$2pd^{\frac{1}{2}}$	1 Mul
[33], [34]	ℓ	(n, m)	mn	mn	$mn \log p$	1 Mul + 1 CMul
[35] [†]	$\sqrt{\ell}$	$(1, 1)$	d	d	$d \log d + d$	1 Mul + 1 CMul
[36] [†]	$\sqrt{\ell}$	$(1, 1)$	d	$5d$	$3d + 5d^{\frac{1}{2}}$	1 Mul + 2 CMul
[37] [†]	$\sqrt{\ell}$	$(1, 1)$	d	$2d$	$2d + 4d^{\frac{1}{2}}$	1 Mul + 1 CMul
[38] [†]	$\sqrt[3]{\ell}$	$(1, 1)$	1	$2d$	$2d + 3 \log d - 2$	1 Mul + 1 CMul
[39], [40] [‡]	$\sqrt{\ell}$	$(1, 1)$	m	m	$m \log d + m$	1 Mul + 1 CMul
[41]*	$O(\sqrt{\ell})$	$(1, 1)$	d	$2d$	$4d$	1 Mul + 1 CMul
[42] [†]	$O(\sqrt{\ell})$	$(1, 1)$	1	d	$4d$	1 Mul + 1 CMul
BMM-I [§]	$\sqrt{\ell/2}$	$(1, 1)$	m	0	$2m + 2$	1 Mul
BMM-II ^{†, **}	$\sqrt[3]{\ell}$	$(1, 1)$	1	0	$3 \log d$	1 Mul

* $d = \max(n, m)$. [†] $d = \max(n, m, p)$. [‡] $d = \max(m, p)$.

[§] (n, m, p) are pairwise coprime and $\max(n, p) < m$.

** (n, m, p) are pairwise coprime and m is a power-of-two integer.

to the coefficient vector of the polynomial. The other is **slots-encoding**, i.e., the polynomial evaluations at certain ℓ points, where ℓ is the number of slots the FHE scheme supports. Roughly speaking, all existing algorithms for encrypted matrix multiplication can be categorized into two types, both of which support SIMD operations.

1) *Coefficients-Encoding Algorithms*: Duong et al. [26] generalize Yasuda et al.'s method for secure inner product [27], [28], and present two algorithms for encrypted matrix multiplication. Suppose that a ciphertext encodes a vector of dimension at most ℓ . Then the two algorithms of Duong et al. support maximal dimensions of $O(\ell^{1/2})$ and $O(\ell^{1/3})$, and require $O(\ell^{1/2})$ and only one Mul, respectively. The drawback of this method is its lack of efficiency in handling consecutive matrix multiplications. Later on, Mishra et al. [29] extended it to support k matrix multiplications successively. However, the maximal dimension is limited only to $O(\ell^{1/(k+1)})$, which makes it impractical. Zheng et al. [30] recently proposed a new framework for homomorphic matrix multiplication under BGV, which supports consecutive matrix multiplication and requires only constant Muls and CMuls, and $O(\log d)$ Rots for square matrix multiplication of dimension d when $d = O(\sqrt[3]{N})$, achieving the best theoretical complexity bound, where N is the ring dimension used in BGV. Zheng et al. also presented a Strassen variant [31] for matrices with large dimensions. However, since their algorithm relies on a hypercube structure of the plaintext space, it does not support B/FV or CKKS. In addition, for all the coefficient encoded methods, to reuse partial entries of the encrypted resulting matrix, we are typically compelled to resort to extra operations for switching between the coefficients and slots encoding, which may slow down the computation.

2) *Slots-Encoding Algorithms*: There are also algorithms that encode matrix data in evaluations of plaintext polynomials, which naturally support consecutive matrix multiplication. We summarize these algorithms in Table I. Halevi and Shoup investigate linear transformation on encrypted vector [10],

[32], i.e., matrix-vector multiplication. Their methods can be directly extended to matrix multiplication. Lu et al. [33] and Wang and Huang [34] extended Halevi and Shoup's method for matrix-matrix multiplication based on the row-order and column-order encoding methods, respectively. Rathee et al. [35] considered an encrypted version of a matrix multiplication algorithm presented in [43]. Jiang et al. presented an algorithm for matrix multiplication over encrypted data in [36]. It uses SIMD operations and the technique for linear transformation [10], [32]. A recent survey [44] identifies Jiang et al.'s algorithm [36] as the state-of-the-art for FHE-based matrix multiplication. Based on Jiang et al.'s algorithm [36], Huang et al. [14] improved the block matrix multiplication for rectangular matrices with special shapes. Chiang [39] and Huang and Zong [40] presented a scheme for non-square matrices, which can be considered as a generalization and optimization of [35]. Zhu et al. [41] have made further improvements based on the method in [40]. Jang et al. [37] presented an adapted CKKS scheme to support data with tensor structure better and improved Jiang et al.'s algorithm [36] in the number of required Rots and CMuls. However, the security of Jang et al.'s variant of CKKS is based on a non-standard hardness assumption called multivariate polynomial learning with errors (m-RLWE). Rizomiliotis and Triakosia [38] introduced a new method for matrix multiplication over encrypted data. This method fully leverages packing techniques, reducing the required number of Muls to just one. However, it only supports matrix dimensions $\ell^{1/3}$.

3) *Other Optimizations*: In addition to optimizing the number of rotations, another optimization direction is to accelerate Rot itself. For example, the hoisting technique proposed in [10] optimizes scenarios involving multiple rotations on the same ciphertext, and the double hoisting introduced by Bossuat et al. [45] makes further progress. These techniques were recently used to accelerate related matrix operations in PCA (Principal Component Analysis) [46]. Gao and Gao [42] construct a fully homomorphic encryption scheme, named GMS, for matrices based on the n -secret LWE assumption [47], which is suitable for matrix multiplication. For encrypted matrices with large dimensions, the Strassen algorithm [31] has been applied recently to this area in [30], [48], and [49].

B. Contribution

Let \mathbf{A} and \mathbf{B} be two matrices with dimensions $n \times m$ and $m \times p$, respectively. Denote by (n, m, p) matrix multiplication the multiplication between \mathbf{A} and \mathbf{B} .

1) *For Small-Dimensional Matrices*: We prove the following

Theorem 1: Let (n, m, p) be pairwise coprime integers and ℓ the number of slots that the FHE scheme supports.

- If $\ell > 2 \cdot \max\{mn, mp, np\}$, there exists an algorithm (BMM-I) that computes a homomorphically encrypted (n, m, p) matrix multiplication within m ciphertext-ciphertext multiplications (Muls), and $2(m + \log[p/m] + \log[n/m] + 1)$ rotations on ciphertexts (Rots), and costs only one Mul multiplicative depth.
- If $\ell > mnp$, there exists an algorithm (BMM-II) that computes a homomorphically encrypted (n, m, p) matrix

multiplication with only one *Mul* and *CMul*, and at most $\log[m] + \log[n] + \log[p]$ *Rots*, and costs one *Mul* and one *CMul* multiplicative depth. In particular, if m is a power-of-two integer, the algorithm can finish the computation with only one *Mul* and at most $\log[m] + \log[n] + \log[p]$ *Rots*, without *CMul*, and hence costs only one *Mul* multiplicative depth.

Both BMM-I and BMM-II support all SIMD-supported FHE schemes (e.g., BGV, B/FV, CKKS) and features the following:

a) *The Required Number of Ciphertext Operations*: As indicated in Table I, BMM-I and BMM-II has their own advantages on the number of ciphertext operations. In particular, BMM-II requires only one *Mul* and $O(\log d)$ *Rots* without *CMul*, better than all other existing algorithms in Table I. Although a similar result can be achieved by Zheng et al.'s algorithm [30], theirs only apply to BGV, not BFV or CKKS.

b) *Multiplicative Depth*: Both BMM-I and BMM-II can achieve optimal multiplicative depth, i.e., only one *Mul* depth. Compared to Halevi-Shoup's method [10], [32], our approach utilizes fewer ciphertexts. For instance, our computation results in a single ciphertext, while theirs produces p ciphertexts.

c) *Transpose for Free*: Our algorithms rely on a novel matrix encoding method given in Section III-A. As a benefit of this encoding, the transpose of an encrypted matrix can be computed for completely free; see Corollary 1. In previous algorithms, e.g., [36], the transpose of an encrypted matrix is reduced to a higher-dimensional linear transformation. This feature is expected to accelerate those applications involving matrix transpose, e.g., computing the covariance matrix in PCA, backpropagation in deep learning, etc.

2) *Matrix With High Dimension*: BMM-I (resp. BMM-II) has a limitation: It requires $\ell > 2 \cdot \max\{mn, mp, np\}$ (resp. $\ell > mnp$). Assuming $n \approx m \approx p$ gives $n \approx \sqrt{\ell/2}$ (resp. $n \approx \sqrt[3]{\ell}$), while Jiang et al.'s algorithm [36] supports encrypted matrix multiplication of dimension $\sqrt{\ell}$. Thus, for handling high-dimensional matrices in a block-wise manner, the number of blocks would be a bit more than that of some existing algorithms, e.g., Jiang et al.'s [36]. This may lead to a lower efficiency of the block-wise algorithms based on BMM-I or BMM-II. To address this problem, we fully exploit the properties of our novel encoding method (on which BMM-I and BMM-II depend) and introduce a segmented version of BMM-I for multiplying high-dimensional matrices (BMM-III), where the utilization rate of slots achieves nearly 100%, similar to, e.g., Jiang et al.'s algorithm. This leads to our second result:

Theorem 2: Assume that (n, m, p) are pairwise coprime integers and ℓ is the number of slots the FHE scheme supports. Then there exists an algorithm (BMM-III) that computes the homomorphically encrypted (n, m, p) matrix multiplication within $m \cdot \lceil \frac{np}{\ell} \rceil$ ciphertext-ciphertext multiplications (*Mul*), and $2m \cdot \lceil \frac{np}{\ell} \rceil$ rotations on ciphertexts (*Rot*), $(4 \cdot \lceil \frac{np}{\ell} \rceil + 2)m + n + p$ plaintext-ciphertext multiplications (*CMul*), and costs only one *Mul* and one *CMul* multiplicative depth.

As a consequence, BMM-III reduces the number of *Rots* by a factor $\frac{1}{3}$ and saves one depth of *CMul* compared with the state-of-the-art algorithm for high-dimensional encrypted rectangular matrices [14]. Furthermore, it depends on a

so-called *segmented* matrix multiplication, which is different from the traditional block matrix multiplication.

3) *Experimental Study*: We implement all BMM-I, BMM-II and BMM-III with CKKS in SEAL [15], including the naïve (textbook) block version and the Strassen [31] version of BMM-I as well. The code is available at https://github.com/hangenba/bicyclic_mat_mul. A comprehensive experimental study demonstrates the performances of these algorithms and identifies how to select different algorithms for different cases to achieve optimal efficiency. In particular, our implementation of BMM-III can compute a $(2048, 8, 2048)$ encrypted matrix multiplication in about 81s, achieving a 2.6x speedup compared with the state-of-the-art algorithm for rectangular matrices [14], and for a task with dimension $(1024, 1024, 1024)$, it costs about 1200s, 5x faster than the block version of Jiang et al.'s algorithm [36]; BMM-II can compute a $(15, 16, 17)$ encrypted matrix multiplication in about 13ms, about 16x faster than the algorithm presented in [38]. SegLKS (obtained from applying the segment strategy to an algorithm presented by Lu et al. [33]) with some optimizations can compute a $(32, 33, 13847)$ encrypted matrix multiplication within 8.24s, about 38x faster the block algorithm based on [36].

C. Technique Overview

Our results mainly rely on two techniques: bicyclic encoding for matrices and segmented matrix multiplication.

1) *Bicyclic Encoding*: We first define a novel encoding map that identifies an $n \times m$ matrix as a vector of dimension mn , provided n and m coprime. We call it the *bicyclic encoding*, which can be roughly viewed as an extension of the diagonal vector of a matrix employed in Halevi-Shoup's algorithm [32]. It follows from the Chinese Remainder Theorem (CRT) that the coprime restriction on n and m guarantees that $\mathbb{Z}/(mn\mathbb{Z}) \cong \mathbb{Z}/(n\mathbb{Z}) \otimes \mathbb{Z}/(m\mathbb{Z})$, which implies that a single vector is enough to traverse all elements of an $n \times m$ matrix. For instance, the bicyclic encoding for

$$A = \begin{pmatrix} 0 & 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 & 9 \end{pmatrix} \text{ and } B = \begin{pmatrix} A & B & C \\ D & E & F \\ G & H & I \\ J & K & L \\ M & N & O \end{pmatrix} \quad (1)$$

are $\mathbf{a} = (0, 6, 2, 8, 4, 5, 1, 7, 3, 9)$ and $\mathbf{b} = (A, E, I, J, N, C, D, H, L, M, B, F, G, K, O)$, respectively. In particular, the k -th entry of \mathbf{a} for $k = 0, 1, \dots, 9$ is the (i, j) -entry of A with $i = k \bmod 2$ and $j = k \bmod 5$.

Let A and B be two matrices to be multiplied, with dimensions (n, m) and (m, p) , respectively, satisfying (n, m, p) pairwise coprime and $m > \max(n, p)$, where the latter condition is just for simplicity and will be removed in our main algorithm. Assume that vectors \mathbf{a} and \mathbf{b} are obtained via the bicyclic encoding of A and B , respectively. Then the bicyclic encoding of the resulting matrix $X = AB$ is exactly $\sum_{0 \leq i < m} \mathbf{a}_i \odot \mathbf{b}_i$, where \odot denotes component-wise multiplication, and \mathbf{a}_i and \mathbf{b}_i are of dimension np obtained by rotating the original \mathbf{a} and \mathbf{b} , respectively. Taking A and B as in Eq. (1), i.e., $(n, m, p) = (2, 5, 3)$, we give in Table II

TABLE II

AN ILLUSTRATIVE EXAMPLE FOR ALGORITHM 1, WHERE s_a (RESP. s_b) IS THE STEP SIZE FOR ROTATING \mathbf{a} (RESP. \mathbf{b}) TO THE LEFT

i	(s_a, s_b)	\mathbf{a}_i	\mathbf{b}_i	$\mathbf{a}_i \odot \mathbf{b}_i$
0	(0, 0)	(0, 6, 2, 8, 4, 5)	(A, E, I, J, N, C)	(0A, 6E, 2I, 8J, 4N, 5C)
1	(8, 3)	(3, 9, 0, 6, 2, 8)	(J, N, C, D, H, L)	(3J, 9N, 0C, 6D, 2H, 8L)
2	(6, 6)	(1, 7, 3, 9, 0, 6)	(D, H, L, M, B, F)	(1D, 7H, 3L, 9M, 0B, 6F)
3	(4, 9)	(4, 5, 1, 7, 3, 9)	(M, B, F, G, K, O)	(4M, 5B, 1F, 7G, 3K, 9O)
4	(2, 12)	(2, 8, 4, 5, 1, 7)	(G, K, O, A, E, I)	(2G, 8K, 4O, 5A, 1E, 7I)

an illustrative example for a matrix multiplication algorithm (Algorithm 1) under bicyclic encoding. We should note that the number of rotation positions for \mathbf{a}_i and \mathbf{b}_i is nontrivial; see Section III-B for details. Clearly, this algorithm requires 5 component-wise vector multiplications and 8 vector rotations.

2) *Segmented Matrix Multiplication*: For matrices with high dimension, the usual approach involves block recursive algorithms such as the Strassen algorithm [31]. However, when applied to encrypted matrix multiplication, using a recursive algorithm leads to a significant memory overhead, while rewriting the recursive algorithm into a loop increases the required multiplication depth. To address this problem, we introduce a technique called *segmented matrix multiplication*.

The basic idea of segmented matrix multiplication is to implement encrypted matrix multiplication strictly following the plaintext algorithm for matrix multiplication under bicyclic encoding. Continuing with the matrices from equation (1) as an example, suppose the number of plaintext slots $\ell = 2$. Then, the bicyclic encoding of \mathbf{A} will be encrypted into five ciphertexts $(\text{ct}.\mathbf{a}_i)_{0 \leq i < 5}$ corresponding to encryptions of (0, 6), (2, 8), (4, 5), (1, 7) and (3, 9). The primary task then becomes how to perform rotations on these ciphertexts. To fulfill this function, we design a subroutine called *Long Rotation* (LongRot). For instance, running a LongRot on $(\text{ct}.\mathbf{a}_i)_{0 \leq i < 5}$ with step size one (i.e., rotating one position towards left) returns ciphertexts of (6, 2), (8, 4), (5, 1), (7, 3) and (9, 0). Based on LongRot, we present BMM-III, which supports high-dimensional encrypted matrix multiplication. For encrypted rectangular matrix multiplication, BMM-III reduces the number of Rots by a factor $\frac{1}{3}$ and saves one CMul depth, compared with the state-of-the-art algorithm [14]. The experiment in Section VII shows that BMM-III is efficient for high-dimensional matrix multiplication. Applying the segmented technique with several optimizations to an algorithm in [33] leads to an algorithm SegLKS requiring the least number of Rots asymptotically among all existing algorithms.

Outline: In Section II, we provide the necessary preliminaries. In Section III, we introduce the bicyclic encoding method and discuss matrix operations under this encoding, including matrix transpose, matrix multiplication, and switching between bicyclic encoding and the commonly used row/column encoding. In Section IV, we present the encrypted version of the aforementioned matrix multiplication algorithms and analyze their cost. In Section V, we present the encrypted segmented matrix multiplication and apply the segmented strategy to Lu et al.'s algorithm [33] in Section VI. We present a comprehensive experimental study in Section VII and conclude with Section VIII.

II. HOMOMORPHIC OPERATIONS

A typical FHE scheme consists of the following algorithms:

- **Setup**(1^λ). Given a security parameter λ , output **parms**.
- **KeyGen**(**parms**). Output a secret key $\text{sk} = s$ and the corresponding public key pk . (For convenience, we also let pk include one or more evaluation keys.)
- **Enc_{pk}**(b). Given a message $b \in \mathcal{M}$, output a ciphertext $c \in \mathcal{C}$, where \mathcal{M} and \mathcal{C} are the plaintext space and ciphertext space, respectively.
- **Dec_{sk}**(c). Given $c \in \mathcal{C}$ as input, output $b \in \mathcal{M}$.
- **Eval_{pk}**($f, (c_1, \dots, c_k)$). Given a function f in k variables, and $(c_i)_{i \leq k}$ with $c_i \leftarrow \text{Enc}_{\text{pk}}(b_i)$, output $c \in \mathcal{C}$ such that $\text{Dec}_{\text{sk}}(c) \neq f(b_1, \dots, b_k)$ holds with negligible probability.

We omit pk or sk for simplicity without ambiguity. An FHE scheme is said to be *secure* if it is IND-CPA secure. The security of almost all existing FHE schemes is based on the assumptions of LWE [50], RLWE [51], or their variants.

Now we recall some homomorphic operations of an SIMD-supported FHE scheme. For convenience, we take CKKS [6] as an example. In CKKS, the plaintext space is $\mathcal{M} = \mathbb{Z}[X]/\langle X^N + 1 \rangle =: R$ while messages are complex vectors in \mathbb{C}^ℓ with $\ell = N/2$, where N is a power-of-two integer. The ciphertext space of CKKS is $\mathcal{C} = R/qR$, where q is the *ciphertext modulus*, a large integer. The restriction of the canonical embedding $\mathbb{R}[X]/\langle X^N + 1 \rangle \rightarrow \mathbb{C}^\ell$ on R maps $m(X) \in R$ into $\mathbf{m} \in \mathbb{C}^\ell$ by evaluating $m(X)$ at the primitive $2N$ -roots of unity $\xi_j = \xi^{Nj}$ for $0 \leq j < \ell$. The inverse of the canonical embedding encodes a message \mathbf{m} as a plaintext $m(X)$. Thus, CKKS naturally supports SIMD operations, i.e., performing an operation on a ciphertext corresponds to performing the same operation on $\ell = N/2$ entries of \mathbf{m} in parallel. Each entry of the message $\mathbf{m} \in \mathbb{C}^\ell$ is called a *plaintext slot*.

For $\mathbf{x} = (x_i)_{0 \leq i < \ell}$ and $\mathbf{y} = (y_i)_{0 \leq i < \ell}$, let $\text{ct}.\mathbf{x}$ and $\text{ct}.\mathbf{y}$ be the ciphertext encrypted by CKKS under the same public key. CKKS supports the following basic operations:

- **Add**($\text{ct}.\mathbf{x}, \text{ct}.\mathbf{y}$): $\text{Dec}(\text{Add}(\text{ct}.\mathbf{x}, \text{ct}.\mathbf{y})) = \mathbf{x} + \mathbf{y}$.
- **Mul**($\text{ct}.\mathbf{x}, \text{ct}.\mathbf{y}$): $\text{Dec}(\text{Mul}(\text{ct}.\mathbf{x}, \text{ct}.\mathbf{y})) = \mathbf{x} \odot \mathbf{y}$, where \odot is for component-wise multiplication.
- **CMul**($\mathbf{m}, \text{ct}.\mathbf{x}$): $\text{Dec}(\text{CMul}(\mathbf{m}, \text{ct}.\mathbf{x})) = \mathbf{m} \odot \mathbf{x}$, where \mathbf{m} is a message in \mathbb{C}^ℓ ; for $m \in \mathbb{C}$, $\text{CMul}(m, \text{ct}.\mathbf{x})$ is a special case of $\text{CMul}(\mathbf{m}, \text{ct}.\mathbf{x})$ with $\mathbf{m} = (m, \dots, m)$.
- **Sl_[i,j]**($\text{ct}.\mathbf{x}$) convert a ciphertext $\text{ct}.\mathbf{x} = \text{Enc}(x_0, \dots, x_{\ell-1})$ into a ciphertext that encrypts $(0, x_i, x_{i+1}, \dots, x_j, 0)$, equivalent to $\text{CMul}(\mathbf{m}, \text{ct}.\mathbf{x})$ for $\mathbf{m} = (\mathbf{0} \in \mathbb{Z}^{i-1}, \mathbf{1} \in \mathbb{Z}^{j-i+1}, \mathbf{0})$.
- **Rot_k**($\text{ct}.\mathbf{x}$) convert $\text{ct}.\mathbf{x} = \text{Enc}(x_0, \dots, x_{\ell-1})$ into a new ciphertext $\text{Enc}(x_k, \dots, x_{\ell-1}, x_0, \dots, x_{k-1})$.

III. BICYCLIC ENCODING FOR MATRICES

In this section, we introduce a novel encoding method for matrices, disclose an intriguing property of this new encoding for matrix transpose, and present two algorithms for matrix multiplication under this new encoding.

To this end, we fix some notations. Let $\mathbf{A} \in \mathcal{R}^{n \times m}$ and $\mathbf{B} \in \mathcal{R}^{m \times p}$ be two matrices over some ring $\mathcal{R} \subseteq \mathbb{C}$. Denote by $\mathbf{X} \in \mathcal{R}^{n \times p}$ the resulting matrix of their multiplication,

i.e., $X = AB$. The transpose of a matrix A is denoted by A^T . Let $[i]_k$ be the non-negative representation of the residue class of i in $\mathbb{Z}/(k\mathbb{Z})$. All indices of vectors and matrices start from 0 unless otherwise specified. For an integer k , we define a *rotation* of a vector $\mathbf{v} = (v_i)_{0 \leq i < n} \in \mathcal{R}^n$ as $\rho_k(\mathbf{v}) = (v_{[k]_n}, v_{[k+1]_n}, \dots, v_{[k+n-1]_n}) \in \mathcal{R}^n$, i.e., $\rho_k(\mathbf{v})$ rotates \mathbf{v} to the left by $[k]_n$ positions.

A. Bicyclic Encoding

Let $A = (a_{i,j}) \in \mathcal{R}^{n \times m}$ be a matrix. Define a map $\varphi_r : \mathcal{R}^{n \times m} \rightarrow \mathcal{R}^r$ with a positive integer $r \leq m \cdot n$ as follows: $\varphi_{n,m,r}(A) = (a_{[0]_n, [0]_m}, a_{[1]_n, [1]_m}, \dots, a_{[r-1]_n, [r-1]_m})$. In particular, if $\gcd(n, m) = 1$ then the elements of $\varphi_{n,m,n \cdot m}(A) \in \mathcal{R}^{n \cdot m}$ exactly traverse each element of A once. The reason is that the indices of the resulting vector are decided by the map $k \mapsto (k \bmod n, k \bmod m)$, which is an isomorphism between $\mathbb{Z}/(mn\mathbb{Z})$ and $\mathbb{Z}/(n\mathbb{Z}) \otimes \mathbb{Z}/(m\mathbb{Z})$ by CRT, provided $\gcd(n, m) = 1$.

We call $\varphi_{n,m,n \cdot m}(A)$ the *bicyclic encoding* of A , denoted by $\mathbf{d}(A)$. For $\ell \geq m \cdot n$, the *bicyclic decoding* map $\psi_{\ell,n,m} : \mathcal{R}^\ell \rightarrow \mathcal{R}^{n \times m}$ that maps a vector $\mathbf{x} \in \mathcal{R}^\ell$ to a matrix $X = (X_{i,j})$ with $X_{i,j} = x_k$ and $k = [i \cdot t \cdot m + j \cdot s \cdot n]_{n \cdot m}$, where (s, t) is a pair of *Bézout coefficients* for (n, m) , i.e., two integers such that $s \cdot n + t \cdot m = 1$. In particular, for the bicyclic encoding $\mathbf{d}(A) = (d_k)_{0 \leq k < n \cdot m}$ of a matrix $A \in \mathcal{R}^{n \times m}$ we have $\psi_{n \cdot m, n, m}(\mathbf{d}(A)) = A$.

The following are some remarks on bicyclic encoding:

- In bicyclic encoding, the traversal of row and column indices of $A \in \mathcal{R}^{n \times m}$ forms two different cycles: one modulo n and the other modulo m . This is where the name comes from. In fact, we can easily generalize the bicyclic encoding to *d-cyclic encoding* for d multidimensional arrays (tensor) $A \in \mathcal{R}^{n_1 \times \dots \times n_d}$ if n_1, \dots, n_d are pairwise coprime, since CRT still holds in this case.
- The bicyclic encoding for matrices can be viewed as an extension of the diagonal vectors [52, Fig. 1-35] employed in Halevi-Shoup's algorithm [32]. The primary difference lies in the fact that diagonal vectors are defined for square matrices, and a d -dimensional square matrix has d diagonal vectors. In contrast, the bicyclic encoding presented in this paper is effective for matrices with coprime dimensions, and the bicyclic encoding of a matrix is exactly a single vector.
- For a matrix $A \in \mathcal{R}^{n \times m}$ with (n, m) coprime, the first component of $\mathbf{d}(A)$ is the $(0, 0)$ -entry of A , which can actually be adjusted to start from any entry of A .

B. Matrix Multiplication Under Bicyclic Encoding

Now we present two algorithms for matrix multiplication under bicyclic encoding.

In Step 4a and 4b, the step size of rotations should be modulo mn and mp , respectively. In Step 4b, the step size of rotation is $i(rm - n)$ with $p \mid (rm - n)$ instead of ip . This nontrivial condition plays a key role in the proof of the correctness of Algorithm 1.

Proposition 1: Algorithm 1 is correct. It requires at most $2(m + \log[p/m] + \log[n/m] + 1)$ vector rotations and m component-wise vector products.

Proof: Since (m, p) are coprime, n can be represented as an integral linear combination of m and p , which guarantees the existence of r in Step 3. Now we assume that (s, t) is a pair of Bézout coefficients for (n, p) . Then according to the definition of bicyclic decoding, the (i, j) -element of X is $x_k = \sum_{0 \leq l < m} \underline{a}_{l,k} \cdot \underline{b}_{l,j}$, where x_k , $\underline{a}_{l,k}$ and $\underline{b}_{l,j}$ are the k -th element of \mathbf{x} , $\underline{\mathbf{a}}$, and $\underline{\mathbf{b}}$, respectively, and $k = [i \cdot t \cdot p + j \cdot s \cdot n]_{np}$. Furthermore, we have x_k is equal to

$$\sum_{0 \leq l < m} \underline{a}_{l,k} \cdot \underline{b}_{l,j} = \sum_{0 \leq l < m} a_{[k-l \cdot n]_n, [k-l \cdot n]_m} \cdot b_{[k+l(r \cdot m - n)]_m, [k+l(r \cdot m - n)]_p} \quad (2)$$

$$= \sum_{0 \leq l < m} a_{[k]_n, [k-l \cdot n]_m} \cdot b_{[k-l \cdot n]_m, [k+l(r \cdot m - n)]_p} \quad (3)$$

$$= \sum_{0 \leq l < m} a_{i, [k-l \cdot n]_m} \cdot b_{[k-l \cdot n]_m, j} \quad (4)$$

$$= \sum_{0 \leq l < m} a_{i,l} \cdot b_{l,j}. \quad (5)$$

Eq. (2) follows from the definitions of bicyclic encoding and the rotation operator, and the construction of \mathbf{a} and \mathbf{b} in Step 2 of Algorithm 1. Eq. (3) easily follows from the modulo arithmetic. Eq. (4) follows from the fact that $[k]_n = i$ and $[k + i(r \cdot m - n)]_p = j$. In fact, according to the definition of k , there exists an integer q such that $k = i \cdot t \cdot p + j \cdot s \cdot n + q \cdot n \cdot p$. So we have $[k]_n = [i \cdot t \cdot p]_n = i$ because of $s \cdot n + t \cdot p = 1$. Similarly, $[k + i(r \cdot m - n)]_p = [j \cdot s \cdot n + i(r \cdot m - n)]_p = j$, where the last equality follows from $s \cdot n + t \cdot p = 1$ and $p \mid (r \cdot m - n)$. To prove (5), we only need to prove that the set $\{[k - l \cdot n]_m : 0 \leq l < m\}$ forms a complete residue system modulo m . Assume that it is not the case, i.e., there exist l and l' such that $0 \leq l' < l < m$ and $[k - l \cdot n]_m \neq [k - l' \cdot n]_m$. This assumption implies that there exists a nonzero integer u such that $k - l \cdot n + m \cdot u = k - l' \cdot n$, so we have $(l - l')n = m \cdot u$. Recalling $\gcd(n, m) = 1$, it gives $n \mid u$, i.e., there exists a nonzero integer v such that $u = n \cdot v$. Hence, $(l - l') = m \cdot v$, which implies $|l - l'| > m$. This contradicts with $0 \leq l' < l < m$, which completes the proof.

Algorithm 1

Input: $A \in \mathcal{R}^{n \times m}$, $B \in \mathcal{R}^{m \times p}$, (n, m, p) pairwise coprime.

Output: Matrix $X = AB$.

- 1) Initialize $\mathbf{a} := \mathbf{d}(A)$, $\mathbf{b} := \mathbf{d}(B)$, and $\mathbf{x} := \mathbf{0} \in \mathcal{R}^{n \cdot p}$.
 $\begin{matrix} \text{[p/m] times} & & \text{[n/m] times} \\ \overbrace{(\mathbf{a}, \dots, \mathbf{a})} & & \overbrace{(\mathbf{b}, \dots, \mathbf{b})} \end{matrix}$
- 2) Update $\mathbf{a} := \overbrace{(\mathbf{a}, \dots, \mathbf{a})}^{\text{[p/m] times}}$ and $\mathbf{b} := \overbrace{(\mathbf{b}, \dots, \mathbf{b})}^{\text{[n/m] times}}$.
- 3) Compute the smallest positive integer r satisfying $r \cdot m - n > 0$ and $p \mid (r \cdot m - n)$.
- 4) For $0 \leq i < m$ do
 - a) Set $\underline{\mathbf{a}}_i := \rho_{-i \cdot n}(\mathbf{a})$, update $\underline{\mathbf{a}}_i$ as its first np entries.
 - b) Set $\underline{\mathbf{b}}_i := \rho_{i \cdot (r \cdot m - n)}(\mathbf{b})$, update $\underline{\mathbf{b}}_i$ as its first np entries.
 - c) Update $\mathbf{x} := \mathbf{x} + \underline{\mathbf{a}}_i \odot \underline{\mathbf{b}}_i$.
- 5) Decode $X := \psi_{np, n, p}(\mathbf{x})$.

The for loop of Algorithm 1 requires $2m - 2$ vector rotations and m vector Hadamard products, and Step 2 requires at most $2(\log[p/m] + \log[n/m] + 2)$ extra vector rotations (see Proposition 4), which completes the proof. \square

Definition 1: For $\mathbf{c} = (c_i)_i \in \mathcal{R}^n$ and an integer k with k dividing n , the segment-sum of \mathbf{c} with length k is defined to be $\mathbf{s} = (s_i)_i \in \mathcal{R}^k$ with $s_i = \sum_{0 \leq j < n/k} c_{j \cdot k + i}$.

Now we propose the second matrix multiplication algorithm (Algorithm 2) under bicyclic encoding.

Proposition 2: Algorithm 2 is correct. It requires at most $\log[n] + \log[m] + \log[p]$ vector rotations and only one component-wise vector product of dimension mnp .

Algorithm 2

Input: $\mathbf{A} \in \mathcal{R}^{n \times m}$, $\mathbf{B} \in \mathcal{R}^{m \times p}$, (n, m, p) pairwise coprime.

Output: Matrix $\mathbf{X} = \mathbf{AB}$.

- 1) Initialize $\mathbf{a} := \mathbf{d}(\mathbf{A})$ and $\mathbf{b} := \mathbf{d}(\mathbf{B})$.
 - 2) Set $\mathbf{a} := (\mathbf{a}, \dots, \mathbf{a}) \in \mathcal{R}^{mnp}$, $\mathbf{b} := (\mathbf{b}, \dots, \mathbf{b}) \in \mathcal{R}^{mnp}$.
 - 3) Compute $\mathbf{x} := \mathbf{a} \odot \mathbf{b}$.
 - 4) Compute the segment-sum of the vector \mathbf{x} with length np .
 - 5) Decode $\mathbf{X} := \psi_{mnp, n, p}(\mathbf{x})$.
-

From the perspective of matrix multiplication in plaintext, both Algorithm 1 and 2 require $O(mnp)$ arithmetic operations. However, their vectorized calculation style may be employed to accelerate matrix multiplication on certain heterogeneous platforms, such as FPGA and GPU.

C. Encrypted Matrix Operations Under Bicyclic Encoding

We now start to discuss some matrix operations on encrypted data under bicyclic encoding, but we defer the encrypted matrix multiplication to the next section.

1) *Switching Between Encrypted Bicyclic Encoding and Row Encoding:* The first operation is how to convert between the bicyclic encoding and the commonly used row encoding or column encoding. Here, we only discuss the row encoding, since the discussion for the column encoding can be obtained similarly. For a matrix $\mathbf{A} = (a_{i,j}) \in \mathcal{R}^{n \times m}$ with (n, m) coprime, the row encoding of \mathbf{A} is defined as a vector $\mathbf{r}(\mathbf{A}) = (a_{[k'/m], [k']_m})_{0 \leq k' < mn}$. Suppose that $\mathbf{d}(\mathbf{A})$ is the bicyclic encoding of \mathbf{A} . Then there exists a linear transformation \mathbf{T} between $\mathbf{r}(\mathbf{A})$ and $\mathbf{d}(\mathbf{A})$. In particular, $\mathbf{d}(\mathbf{A}) = \mathbf{T} \cdot \mathbf{r}(\mathbf{A})$, where $t_{k,k'} = 1$ with k and k' determined as follows. For $0 \leq i < n$ and $0 \leq j < m$, we first decide $k' = i \cdot m + j$, and then compute $k = [i \cdot t \cdot m + j \cdot s \cdot n]_{mn}$, where (s, t) is a pair of Bézout coefficients for (n, m) , as in the bicyclic decoding process. One can decide a matrix \mathbf{T}' satisfying $\mathbf{r}(\mathbf{A}) = \mathbf{T}' \cdot \mathbf{d}(\mathbf{A})$ similarly.

Given an encryption of $\mathbf{r}(\mathbf{A})$, we can use the diagonal encoding introduced by Halevi-Shoup to perform linear transformations on ciphertexts [32], thereby accomplishing the conversion between the two encodings. Assuming the matrix \mathbf{T} has d non-zero diagonal vectors, this transformation can be computed within d CMuls and $2\sqrt{d}$ Rots, and requires only one level of CMul multiplication depth [36].

2) *Encrypted Matrix Transpose Under Bicyclic Encoding:* Under bicyclic encoding, we show that one can transpose a matrix for free, either in plaintext or encrypted form.

Proposition 3: For a matrix $\mathbf{A} \in \mathcal{R}^{n \times m}$ with $\gcd(n, m) = 1$, we have $\mathbf{d}(\mathbf{A}) = \mathbf{d}(\mathbf{A}^T)$.

Corollary 1: Let $(\text{ct}.\mathbf{a}_i)_{i < \lceil \frac{mn}{\ell} \rceil}$ be ciphertexts (under an FHE scheme that supports ℓ slots) of the bicyclic encoding of a matrix $\mathbf{A} \in \mathcal{R}^{n \times m}$ with $\gcd(n, m) = 1$. Then $(\text{ct}.\mathbf{a}_i)_{i < \lceil \frac{mn}{\ell} \rceil}$ are also ciphertexts of the bicyclic encoding of \mathbf{A}^T .

IV. ENCRYPTED MATRIX MULTIPLICATION UNDER BICYCLIC ENCODING

In this section, we always assume that (n, m, p) are coprime, which means bicyclic encoding applies to all matrices \mathbf{A} , \mathbf{B} and $\mathbf{X} = \mathbf{AB}$, denoted by \mathbf{a} , \mathbf{b} , and \mathbf{x} the encoded vectors, respectively. We also assume that all the encoded vectors can be encrypted into a single ciphertext, denoted by $\text{ct}.\mathbf{a}$, $\text{ct}.\mathbf{b}$, and $\text{ct}.\mathbf{x}$, respectively.

A. Building Blocks

For convenience, we first present two building blocks.

1) *The Repeat Operation:* According to Step 2 of Algorithm 1, we need first to convert a ciphertext $\text{ct}.\mathbf{a}$ of \mathbf{a} to a ciphertext of $(\mathbf{a}, \dots, \mathbf{a})$, i.e., repeated with certain times.

Algorithm 3 Repeat

Input: A ciphertext $\text{ct}.\mathbf{a}$ of $(\mathbf{a}, \mathbf{0}) \in \mathcal{R}^\ell$ (where $\mathbf{a} \in \mathcal{R}^d$) and an integer $t = \sum_{i=0}^{\lfloor \log t \rfloor} t_i \cdot 2^i \geq 1$ satisfying $td < \ell$.

Output: A updated ciphertext $\text{ct}.\mathbf{c}$ that encrypts $(\mathbf{a}, \dots, \mathbf{a}, \mathbf{0}) \in \mathcal{R}^\ell$ with \mathbf{a} repeated t times.

- 1) Initialize $\text{ct}.\mathbf{a}_0 \leftarrow \text{ct}.\mathbf{a}$.
 - 2) For $1 \leq i \leq \lfloor \log t \rfloor$ do
 - a) Compute $\text{ct}.\mathbf{a}_i \leftarrow \text{Add}(\text{ct}.\mathbf{a}_{i-1}, \text{Rot}_{-2^{i-1}d}(\text{ct}.\mathbf{a}_{i-1}))$.
 - 3) Set $k = 2^{\lfloor \log t \rfloor} \cdot d$ and $\text{ct}.\mathbf{c} := \text{ct}.\mathbf{a}_{\lfloor \log t \rfloor}$.
 - 4) For $i = \lfloor \log t \rfloor - 1, \lfloor \log t \rfloor - 2, \dots, 1, 0$ do
 - a) If $t_i \neq 0$ then compute $\text{ct}.\mathbf{c} \leftarrow \text{Add}(\text{ct}.\mathbf{c}, \text{Rot}_{-k}(\text{ct}.\mathbf{a}_i))$ and update $k := k + t_i \cdot 2^i \cdot d$.
-

Proposition 4: The Repeat algorithm is correct and requires at most $2 \log t$ Rots and Adds, respectively. In particular, if t is a power-of-two integer, it only requires $\log t$ such operations.

Algorithm 4 SegSum

Input: A ciphertext $\text{ct}.\mathbf{a}$ of $(\mathbf{a}, \mathbf{0}) \in \mathcal{R}^\ell$ (where $\mathbf{a} \in \mathcal{R}^n$) and an integer k satisfying $n = k \cdot m$ for an integer m .

Output: A ciphertext $\text{ct}.\mathbf{c}$ that encrypts $(\mathbf{c}, \mathbf{0}) \in \mathcal{R}^\ell$, where $\mathbf{c} \in \mathcal{R}^k$ is the segment-sum of \mathbf{a} with length k .

- 1) Initialize $\text{ct}.\mathbf{c} \leftarrow \text{ct}.\mathbf{a}$ and $m := n/k$.
 - 2) For $i = \lfloor \log m \rfloor - 1, \lfloor \log m \rfloor - 2, \dots, 1, 0$ do
 - a) Set $\text{ct}.\mathbf{t} := \text{Rot}_{2^i \cdot k}(\text{ct}.\mathbf{c})$
 - b) If $i = \lfloor \log m \rfloor - 1$ then $\text{ct}.\mathbf{t} \leftarrow \text{Sl}_{[0, m-2^i-1]}(\text{ct}.\mathbf{t})$.
// This can be omitted if m is a power-of-two integer.
 - c) Update $\text{ct}.\mathbf{c} \leftarrow \text{Add}(\text{ct}.\mathbf{c}, \text{ct}.\mathbf{t})$.
-

2) *The SegSum Operation:* In Step 4 of Algorithm 2, we need to compute the segment-sum (Definition 1) of a vector $\mathbf{a} = (a_i)_i \in \mathcal{R}^n$ with length k satisfying $n = k \cdot m$ for an integer m . The following algorithm is an encrypted form of this process.

Proposition 5: The SegSum algorithm is correct and requires at most $\lceil \log m \rceil$ Rots and Adds, and one CMul. In particular, if m is a power-of-two integer, it only requires $\log m$ Rots and Adds, without CMul.

B. BMM-I (Encrypted Version of Algorithm 1)

We now propose an encrypted version of Algorithm 1. In comparison to the plaintext algorithm (Algorithm 1), the most significant difference lies in Step 2, where the number of repetitions for vectors \mathbf{a} and \mathbf{b} is doubled. This directly leads to the requirement of $\ell > 2 \cdot \max\{mn, mp, np\}$. The reason is that Rot operations are indeed performed on vectors of dimension ℓ . In contrast, in the plaintext algorithm, the rotation operations for \mathbf{a} (resp. \mathbf{b}) are performed on a vector of dimension mn (resp. mp). Thus, we have to double the number of repetitions of the original vectors to guarantee the correctness of the results.

Algorithm 5 BMM-I (Bicyclic Matrix Multiplication)

Input: $\text{ct}.\mathbf{a}$ and $\text{ct}.\mathbf{b}$, which are ciphertexts of the bicyclic encoding of matrices $\mathbf{A} \in \mathcal{R}^{n \times m}$ and $\mathbf{B} \in \mathcal{R}^{m \times p}$, respectively, where (n, m, p) are coprime and the number of slots $\ell > 2 \cdot \max\{mn, mp, np\}$.

Output: $\text{ct}.\mathbf{x}$, whose first $n \cdot p$ slots correspond to the bicyclic encoding of the resulting matrix $\mathbf{X} \in \mathcal{R}^{n \times p}$.

- 1) Initialize $\text{ct}.\mathbf{x} \leftarrow \text{Enc}(\mathbf{0})$. Compute the smallest positive integer r satisfying $p \mid (r \cdot m - n)$ and $r \cdot m - n > 0$.
 - 2) Update $\text{ct}.\mathbf{a} \leftarrow \text{Repeat}(\text{ct}.\mathbf{a}, 2\lceil p/m \rceil)$ and $\text{ct}.\mathbf{b} \leftarrow \text{Repeat}(\text{ct}.\mathbf{b}, 2\lceil n/m \rceil)$.
 - 3) For $0 \leq i < m$ do
 - a) Compute $\text{ct}.\mathbf{a}_i \leftarrow \text{Rot}_{[-i \cdot n]_{mn}}(\text{ct}.\mathbf{a})$.
 - b) Compute $\text{ct}.\mathbf{b}_i \leftarrow \text{Rot}_{[i \cdot (r \cdot m - n)]_{mp}}(\text{ct}.\mathbf{b})$.
 - c) Update $\text{ct}.\mathbf{x} \leftarrow \text{Add}(\text{ct}.\mathbf{x}, \text{Mul}(\text{ct}.\mathbf{a}_i, \text{ct}.\mathbf{b}_i))$.
-

Another difference lies in the for loop: neither $\text{ct}.\mathbf{a}_i$ nor $\text{ct}.\mathbf{b}_i$ is truncated to keep only the first np elements. Even after the for loop, $\text{ct}.\mathbf{x}$ is still not truncated to contain only the first np elements. In fact, if we denote the ℓ -dimensional vector \mathbf{x} by the decryption of $\text{ct}.\mathbf{x}$ returned by BMM-I, then the first np components of \mathbf{x} exactly correspond to the bicyclic encoding of the result matrix $\mathbf{X} = \mathbf{AB}$. Hence, the truncation can be delayed until decryption. This property implies that BMM-I does not need any plaintext-ciphertext multiplication CMul.

Proof: [Proof of Theorem 1 item 1] Under the assumptions on n, m, p and ℓ , bicyclic encoding is available for each matrix, and each bicyclic encoding vector can be encrypted in one ciphertext. Further, the assumption $\ell > 2 \cdot \max\{mn, mp, np\}$ guarantees the correctness of Rots in Step ca and cb. Thus, BMM-I exactly follows the plaintext Algorithm 1, which implies the correctness. The cost of BMM-I follows from the cost of Repeat and counting. \square

Remark 1: Although BMM-I does not use plaintext-ciphertext multiplication CMul, users should know that the computation results are contained only in the first np slots. Suppose that decrypting the $\text{ct}.\mathbf{x}$ obtains $\mathbf{x} \in \mathcal{R}^\ell$. Then running bicyclic decoding $\psi_{\ell, n, p}(\mathbf{x})$ gives the resulting matrix $\mathbf{X} = \mathbf{AB}$. If necessary, these results can be extracted through a single $S_{\lfloor 0, np-1 \rfloor}(\text{ct}.\mathbf{x})$ operation. In fact, performing this selection

remains unnecessary unless there is a need to utilize the latter $\ell - np$ slots.

Comparison with Existing Algorithms: Compared with algorithms designed specifically for square matrices [35], [36], [37], [38], BMM-I offers greater flexibility in matrix dimensions. Despite the pairwise coprime constraints among (n, m, p) , one can use padding with zeros to meet the requirements. Generally, BMM-I needs fewer padding positions. Compared with algorithms that support encrypted matrix multiplication of arbitrary dimensions [39], [40], as well as those facilitating encrypted approximate number matrix operations [33], [35], [36], [37], [38], BMM-I requires the fewest number of multiplication depths and ciphertext rotations. While the algorithm in [38] requires a smaller number of ciphertext-ciphertext multiplications (Mul), our algorithm does not need plaintext-ciphertext multiplication (CMul); see Table I.

However, the condition $\ell > 2 \cdot \max\{mn, mp, np\}$ in BMM-I constrains the dimensions of matrices. For a set of CKKS parameters with $\ell = 2^{12}$, Jiang et al.'s algorithm [36] supports encrypted matrix multiplication of dimensions $(64, 64, 64)$, whereas BMM-I only supports dimensions of $(43, 45, 44)$, which is about $\sqrt{\ell/2}$. (Note that the algorithm in [38] only supports $(16, 16, 16)$, i.e., $\sqrt[3]{\ell}$, encrypted matrix multiplication under the same setting.) However, experimental results in Section VII show that BMM-I is still practical.

Algorithm 6 BMM-II

Input: $\text{ct}.\mathbf{a}$ and $\text{ct}.\mathbf{b}$, which are ciphertexts of the bicyclic encoding of matrices $\mathbf{A} \in \mathcal{R}^{n \times m}$ and $\mathbf{B} \in \mathcal{R}^{m \times p}$, respectively, where (n, m, p) are pairwise coprime and the number of slots $\ell > n \cdot m \cdot p$.

Output: $\text{ct}.\mathbf{x}$, whose first $n \cdot p$ slots correspond to the bicyclic encoding of the resulting matrix $\mathbf{X} = \mathbf{AB} \in \mathcal{R}^{n \times p}$.

- 1) Initialize $\text{ct}.\mathbf{x} \leftarrow \text{Enc}(\mathbf{0})$.
 - 2) Set $\text{ct}.\mathbf{a} \leftarrow \text{Repeat}(\text{ct}.\mathbf{a}, p)$, $\text{ct}.\mathbf{b} \leftarrow \text{Repeat}(\text{ct}.\mathbf{a}, n)$.
 - 3) Compute $\text{ct}.\mathbf{x} \leftarrow \text{Mul}(\text{ct}.\mathbf{a}, \text{ct}.\mathbf{b})$.
 - 4) $\text{ct}.\mathbf{x} \leftarrow \text{SegSum}(\text{ct}.\mathbf{x}, n \cdot p)$
-

C. BMM-II (Encrypted Version of Algorithm 2)

Similarly, we propose an encrypted version of Algorithm 2 as BMM-II (Algorithm 6). We note that the discussion in Remark 1 also holds for BMM-II.

Proof of Theorem 1 item 1: In comparison to Algorithm 2 in plaintext, all steps of BMM-II are essentially the same. Therefore, the item 1 of Theorem 1 follows from Proposition 2, which also completes the proof of Theorem 1. \square

Comparison with BMM-I: Compared with BMM-I, BMM-II imposes stricter constraints on matrix dimensions, requiring $\ell > nmp$. Given a fixed ℓ , this constraint limits the matrix dimensions that BMM-II can support. However, this stricter requirement ensures that we can always construct vectors of dimension nmp from the diagonal encoding of the two input matrices. The product of these two vectors encapsulates all the information of the matrix multiplication result, leading to an optimized result on the number of ciphertext operations.

Comparison with Existing Algorithms: Similar to BMM-I, both the Rizomiliotis-Triakosia (RT) algorithm [38] and Zheng

et al.'s algorithm [30] have a restriction of $d = O(\sqrt[3]{N})$, where $d = \max(n, m, p)$ and N is the dimension of the ring used in the encryption algorithm. Compared to the RT algorithm, the number of RotS required by BMM-II is reduced from a linear function in d to a logarithmic function. Compared to Zheng et al.'s algorithm, the cost of BMM-II is considerable; theoretically, BMM-II and Zheng et al.'s algorithm are the two most cost-effective algorithms. However, because Zheng et al.'s algorithm relies on the tensor structure of the ring, it currently only supports the BGV scheme, while BMM-II can support any homomorphic encryption scheme that supports SIMD, including BGV, B/FV, CKKS.

However, we note that when computing high-dimensional matrix multiplication in a block-wise manner, these algorithms require more blocks than the other algorithms listed in Table I, which makes them not so practical as shown in Section VII. Next, we will discuss how to circumvent this obstacle to support encrypted matrix multiplication with higher dimensions.

V. ENCRYPTED MATRICES OF HIGH DIMENSIONS

Assume that the matrix multiplication with dimensions (n, m, p) satisfies the condition that $\max(mn, mp, np) > \ell$. This implies that at least one matrix involved in the multiplication (either A , B , or X) requires multiple ciphertexts for storage. Under this setting, traditional methods typically resort to block matrix multiplication. Besides, there exists another natural approach, called *segmented strategy* for handling high-dimensional encrypted matrices, in which any bicyclic encoding vector of matrices with dimension larger than ℓ may be divided into multiple vectors of dimension ℓ .

A. Block Matrix Multiplication

Let $A \in \mathcal{R}^{n \times m}$ and $B \in \mathcal{R}^{m \times p}$ be the two matrices to be multiplied. Assume that the number of slots ℓ supports a (n_0, m_0, p_0) matrix multiplication for packed ciphertexts with $m_0 = \max\{n_0, m_0, p_0\}$. For simplicity, we further assume that $n_1 = n/n_0 = m/m_0 = p/p_0$. Then both A and B can be split into $n_1 \times n_1$ blocks. To compute this (n_1, n_1, n_1) block matrix multiplication, one needs to compute n_1^ω matrix multiplications of dimension (n_0, m_0, p_0) with $2 < \omega < 3$, e.g., for Strassen algorithm [31], $\omega = \log 7 \approx 2.81$. We call the resulting algorithm the *block* version of BMM-I, which requires $m_0 n_1^\omega$ Muls and $2(m_0 + 1)n_1^\omega$ RotS. In fact, the block strategy applies to all matrix multiplication algorithms. In Table III, we summarize the cost of some of them, such as [30], [36], and [38]. Note that $n_0 = m_0 = p_0 = \sqrt{\ell}$ for [36] and $n_0 = m_0 = p_0 = \sqrt[3]{\ell}$ for [38].

It appears that the block version of BMM-II and Zheng et al.'s algorithm [30] are faster than the others in Table III. However, as mentioned previously, these two algorithms have to deal with more blocks. In particular, n_1 for these two algorithms is approximately $d/\ell^{1/3}$, while for the others n_1 is about $d/\ell^{1/2}$, where $d = \max(n, m, p)$.

For the block version of BMM-I, if we assume that $n_0 \approx m_0 \approx p_0 \approx \sqrt{\ell/2}$ and $d = \max\{n, m, p\}$, then $\#Mul \leq (\ell/2)^{\frac{1-\omega}{2}} d^\omega$, and $\#Rot \leq 2(\ell/2)^{(1-\omega)/2} d^\omega$, which seems worse

TABLE III
THE COST OF DIFFERENT BLOCK ALGORITHMS FOR HIGH DIMENSIONAL (n, m, p) ENCRYPTED MATRIX MULTIPLICATION WITH $d = \max\{n, m, p\}$

Method	#Mul	#CMul	#Rot	Mult. depth
Block [36]	$\ell^{\frac{1-\omega}{2}} d^\omega$	$5\ell^{\frac{1-\omega}{2}} d^\omega$	$(3\ell^{\frac{1-\omega}{2}} + 5\ell^{\frac{1}{2}(\frac{1}{2}-\omega)})d^\omega$	1 Mul + 2 CMul
Block [38]	$\ell^{-\frac{\omega}{3}} d^\omega$	$2\ell^{\frac{1-\omega}{3}} d^\omega$	$(2\ell^{\frac{1-\omega}{3}} + \ell^{-\frac{1}{3}\omega} \log \ell) d^\omega$	1 Mul + 1 CMul
Block [30] [†]	n_1^ω	$2n_1^\omega$	$2 \log N \cdot n_1^\omega$	1 Mul + 1 CMul
Block BMM-I	$m_0 n_1^\omega$	0	$2(m_0 + 1)n_1^\omega$	1 Mul
Block BMM-II [‡]	n_1^ω	0	$3 \log m_0 \cdot n_1^\omega$	1 Mul

[†] N is the ring dimension of the BGV scheme.

[‡] m_0 is a power-of-two integer.

than that of algorithms in [36] and [38]. However, the advantages of the block version of BMM-I include: it requires only one multiplicative depth, and it needs no CMul. Experiments in Section VII-D show that the block version of BMM-I performs well in practice.

B. Segmented Matrix Multiplication

The bicyclic encoding introduced in Section III allows us to segment the bicyclic encoding vector of matrices, thereby supporting high-dimensional encrypted matrix multiplication following Algorithm 1 exactly. To facilitate this, it is necessary to introduce a fundamental operation called *LongRot*, used to rotate the segmented vectors of a high-dimensional vector.

1) *The LongRot Algorithm*: Given $\mathbf{a} \in \mathcal{R}^d$ with $d > \ell$, the LongRot operation implements the following functionality:

- Construct $\underline{\mathbf{a}} = (\mathbf{a}, \dots, \mathbf{a}) \in \mathcal{R}^{t \cdot d}$, repeating t times of \mathbf{a} ;
- Rotate the vector $\underline{\mathbf{a}}$ to the left by k positions, resulting in $\underline{\mathbf{a}}' = \rho_k(\underline{\mathbf{a}})$;
- Select the first τ elements of $\underline{\mathbf{a}}'$ and divide them into $\lceil \frac{\tau}{\ell} \rceil$ groups, each containing ℓ elements, possibly zero-padding for the last one.

Clearly, this functionality is designed to construct the ciphertexts of $\underline{\mathbf{a}}_i$ and $\underline{\mathbf{b}}_i$ from the ciphertexts of \mathbf{a} and \mathbf{b} in Step b of Algorithm 1, respectively. We omit the detailed description here since this involves only some tedious control structures. For a complete description, we refer the readers to [53].

Algorithm 7 LongRot

Input: Ciphertexts $(\text{ct}.\mathbf{a}_i)_{0 \leq i < w}$ for $\mathbf{a} = (\mathbf{a}_0, \mathbf{a}_1, \dots, \mathbf{a}_{w-2}, \mathbf{a}_{w-1}) \in \mathcal{R}^d$ with $\mathbf{a}_{w-1} \in \mathcal{R}^z$ and $\mathbf{a}_i \in \mathcal{R}^\ell$ for $i = 0, 1, \dots, w-2$ (i.e., $d = (w-1)\ell + z$ with $0 \leq z < \ell$, where ℓ is the number of slots), the number of repeated times t , the number of positions to be rotated $k \in [0, d)$, the number of selected elements τ .

Output: Ciphertexts $(\text{ct}.\underline{\mathbf{a}}'_i)_{0 \leq i < \lceil \frac{\tau}{\ell} \rceil}$, i.e., $\lceil \frac{\tau}{\ell} \rceil$ ciphertexts corresponding to the first τ elements of $\underline{\mathbf{a}}'$.

Proposition 6: The LongRot algorithm correctly computes the above functionality within $\lceil \frac{\tau}{\ell} \rceil$ RotS, $2\lceil \frac{\tau}{\ell} \rceil + \frac{\tau}{d} + 1$ CMuls, and one CMul multiplicative depth.

2) *BMM-III (Segmented Version of BMM-II)*: Now, we present an algorithm (BMM-III) for high-dimensional encrypted matrix multiplication under bicyclic encoding. It is essentially a direct translation of Algorithm 1 into its encrypted

TABLE IV
THE COST OF BMM-III FOR HIGH-DIMENSIONAL RECTANGULAR
MATRIX MULTIPLICATION

Size	Method	#Mul	#CMul	#Rot	Mult. depth
$(n, n, \ell/n)$	[14]	n	$5n$	$3n + \frac{n^2}{\ell} + \frac{n\sqrt{n}}{\sqrt{\ell}} - 1$	1 Mul + 2 CMul
	BMM-III [†]	n	$7n + \frac{\ell}{n}$	$2n$	1 Mul + 1 CMul
$(n, \ell/n, \ell/n)$	[14]	$\frac{\ell}{n}$	$\frac{5\ell}{n}$	$\frac{3\ell}{n} + \frac{6\sqrt{\ell}}{n\sqrt{n}} + \frac{2\ell}{n^2} - 2$	1 Mul + 2 CMul
	BMM-III [†]	$\frac{\ell}{n}$	$\frac{7\ell}{n} + n$	$\frac{2\ell}{n}$	1 Mul + 1 CMul
$(\ell/n, n, \ell/n)$	[14]	$\frac{\ell}{n}$	$\frac{5\ell}{n}$	$\frac{3\ell}{n} + \frac{6\sqrt{\ell}}{n\sqrt{n}} + \frac{\ell}{n^2} - 1$	1 Mul + 2 CMul
	BMM-III [†]	$\frac{\ell}{n}$	$\frac{6\ell}{n} + 2n$	$\frac{2\ell}{n}$	1 Mul + 1 CMul

[†] For BMM-III, we should make the dimensions pairwise coprime.

version. The only difference lies in Step 3a, where the rotation step size is adjusted from $-jn$ to $(m - j)n$. These two are evidently equivalent and meet the requirement of LongRot.

Proof of Theorem 2: From the structure of BMM-III, it is easy to see that the required multiplicative depth is one Mul and one CMul. It follows from Proposition 6 that Step 3a and 3b requires at most $2\lceil \frac{np}{\ell} \rceil + \frac{p}{m} + 1$ and $2\lceil \frac{np}{\ell} \rceil + \frac{n}{m} + 1$ CMuls, respectively, and both costs at most $\lceil \frac{np}{\ell} \rceil$ Rots. Therefore, totally, it requires at most $2m \cdot \lceil \frac{np}{\ell} \rceil$ Rots and $(4 \cdot \lceil \frac{np}{\ell} \rceil + 2)m + n + p$ CMuls. Step 3c costs $\lceil \frac{np}{\ell} \rceil$ Muls. Thus, the total number of Muls is bounded by $m \cdot \lceil \frac{np}{\ell} \rceil$. \square

Algorithm 8 BMM-III

Input: Ciphertexts $(\text{ct}.a_i)_{i < \lceil mn/\ell \rceil}$ for the bicyclic encoding of $A \in \mathcal{R}^{n \times m}$ and ciphertexts $(\text{ct}.b_i)_{i < \lceil mp/\ell \rceil}$ for the bicyclic encoding of $B \in \mathcal{R}^{m \times p}$, where (n, m, p) are pairwise coprime.

Output: Ciphertexts $(\text{ct}.x_i)_{0 \leq i < \lceil np/\ell \rceil}$ for the bicyclic encoding of the resulting matrix $X \in \mathcal{R}^{n \times p}$.

- 1) For $i = 0, 1, \dots, \lceil np/\ell \rceil - 1$ initialize $\text{ct}.x_i \leftarrow \text{Enc}(\mathbf{0})$.
- 2) Compute the smallest positive integer r such that $p \mid (r \cdot m - n)$ and $r \cdot m - n > 0$.
- 3) For $j = 0, 1, \dots, m - 1$ do the following:
 - a) $(\text{ct}.a'_i)_i \leftarrow \text{LongRot}((\text{ct}.a_i)_{i < \lceil \frac{mn}{\ell} \rceil}, \lceil \frac{p}{m} \rceil, (m - j)n, np)$.
 - b) $(\text{ct}.b'_i)_i \leftarrow \text{LongRot}((\text{ct}.b_i)_{i < \lceil \frac{mp}{\ell} \rceil}, \lceil \frac{n}{m} \rceil, j(rm - n), np)$.
 - c) Compute $\text{ct}.x_i \leftarrow \text{Add}(\text{ct}.x_i, \text{Mul}(\text{ct}.a'_i, \text{ct}.b'_i))$ for $i = 0, 1, \dots, \lceil np/\ell \rceil - 1$.

Comparison: For square encrypted matrix multiplication with dimension d , BMM-III requires $O(d^3)$ ciphertext operations. Therefore, in an asymptotic sense, the number of ciphertext operations required by BMM-I is greater than those required by the block version algorithms in Section V-A. However, experiments in Section VII demonstrate that when $d \leq 1024$, BMM-III has a distinct advantage over those block version algorithms. The reason is that most of the block version algorithms are originally recursive. It is well known that the efficiency of recursive algorithms is not that fast. Usually, one can rewrite a recursive algorithm as a loop algorithm. However, in the case of matrix multiplication over encrypted data, such a rewriting will increase the required multiplicative depth regarding the recursion depth, which is unacceptable.

Additionally, BMM-III supports encrypted matrix multiplication of flexible dimensions. In the literature, Huang et al. [14] investigated encrypted matrix multiplication for high-dimensional rectangular matrices with different shapes, including matrix multiplication of dimensions $(n, n, \ell/n)$, $(n, \ell/n, \ell/n)$, and $(\ell/n, n, \ell/n)$, where ℓ is the number of slots. All of these shapes cost similar ciphertext operations. Compared with theirs in Table IV shows that BMM-III reduces the number of Rots by a factor $\frac{1}{3}$ and saves one CMul depth, at a cost of a bit more CMuls.

VI. ANOTHER APPLICATION OF SEGMENTED STRATEGY

In this section, we consider applying the segmented strategy to the algorithm by Lu et al. [33], which is not under bicyclic encoding. This algorithm is no longer the best for matrix multiplication with smaller dimensions, as shown in Table I. However, combining several optimizations and observations, the segmented version of Lu et al.'s algorithm (SegLKS) requires the fewest number of ciphertext rotations in theory among all currently known encrypted matrix multiplication algorithms for high-dimensional matrices.

A. Lu Et Al.'s Algorithm

The algorithm in [33] for encrypted matrix multiplication relies on the Replicate operation. The function of $\text{Replicate}_i(\text{ct}.a)$ is to transform a ciphertext of $a = (a_0, \dots, a_{\ell-1})$ into a ciphertext of (a_i, a_i, \dots, a_i) . A Replicate can be finished within one CMul plus $\log \ell$ Rots and Adds. Lu et al.'s algorithm supports matrix multiplication of any dimension. For computing an (n, m, p) matrix multiplication $X = AB$ in encrypted form, the algorithm encodes all matrices row by row.

Algorithm 9 SegLKS (Segmented variant of [33])

Input: Ciphertexts $(\text{ct}.a_{i,j})_{0 \leq i < n, 0 \leq j < \lceil m/\ell \rceil}$ for $A \in \mathcal{R}^{n \times m}$ and ciphertexts $(\text{ct}.b_{i,j})_{0 \leq i < m, 0 \leq j < \lceil p/\ell \rceil}$ for $B \in \mathcal{R}^{m \times p}$, where $\text{ct}.a_{i,j}$ is the ciphertext corresponding to the j -th segment of the i -th row of a , similar for $\text{ct}.b_{i,j}$.

Output: Ciphertexts $(\text{ct}.x_{i,j})_{0 \leq i < n, 0 \leq j < \lceil p/\ell \rceil}$ for the resulting matrix $X \in \mathcal{R}^{n \times p}$.

- 1) For $i = 0$ to $n - 1$ do
 - a) For $j = 0$ to $\lceil p/\ell \rceil - 1$ do
 - i) For $k = 0$ to $\lceil m/\ell \rceil$ do
 - A) For $\iota = 0$ to $\ell - 1$ compute $\text{ct}.x_{i,j} \leftarrow \text{Add}(\text{ct}.x_{i,j}, \text{Mul}(\text{Replicate}_\iota(\text{ct}.a_{i,k}), \text{ct}.b_{k\ell+\iota,j}))$.
- 2) Return $(\text{ct}.x_{i,j})_{0 \leq i < n, 0 \leq j < \lceil p/\ell \rceil}$.

When $\max(m, p) < \ell$, the number of ciphertexts in Lu et al.'s algorithm corresponding to A , B and X are n , m and n , respectively. It requires mn Muls and CMuls, and $mn \log p$ Rots. Note that only $\log \ell$ key-switching keys are enough for Lu et al.'s algorithm, which is used to replicate each entry of A . In addition, since the matrices are encoded by rows, it naturally supports the segmented method.

TABLE V

THE COST OF SEGMENTED ALGORITHMS FOR HIGH DIMENSIONAL ENCRYPTED (n, m, p) MATRIX MULTIPLICATION WITH $n \leq p$

Method	#Mul	#CMul	#Rot	#Mult. depth
BMM-III	$m \lceil \frac{np}{\ell} \rceil$	$4m \lceil \frac{np}{\ell} \rceil + 2(n + p + 2m)$	$2m \lceil \frac{np}{\ell} \rceil$	1 Mul + 1 CMul
SegLKS/O3	$mn \lceil \frac{p}{\ell} \rceil$	mn	$mn \log \ell$	1 Mul + 1 CMul
SegLKS/O4	$mn \lceil \frac{p}{\ell} \rceil$	$nm(1 + \frac{1}{\kappa})$	$nm(\frac{\log \frac{\ell}{\kappa}}{\kappa} + \log \kappa)$	1 Mul + 2 CMul

B. SegLKS (Segmented Lu Et Al.'s Algorithm)

Now we consider matrix multiplication of high dimensions, where each row of A or B is encoded and encrypted as multiple ciphertexts in a segmented manner, say, $\min(m, p) > \ell$. The number of ciphertexts corresponding to A , B and X are $n \lceil \frac{m}{\ell} \rceil$, $m \lceil \frac{p}{\ell} \rceil$ and $n \lceil \frac{p}{\ell} \rceil$, respectively. Further, this algorithm requires $mn \lceil \frac{p}{\ell} \rceil$ Muls and CMuls, and $mn \log p$ Rots. However, we can optimize the algorithm further.

1) *On Encoding (O1)*: Since the cost of SegLKS heavily relies on a factor nm , when $nm \gg mp$, it is costly. If this is the case, we can encode the matrix by columns, or, equivalently, consider the matrix multiplication in transpose $X^T = B^T A^T$. So, one can always assume that $n \leq p$.

2) *On the Number of Ciphertexts (O2)*: For matrix $A \in \mathcal{R}^{n \times m}$, since it only involves the Replicate operation, it can be encoded and encrypted into fewer ciphertexts without affecting efficiency. Indeed, it can be encrypted in a row-wise manner into $\lceil \frac{nm}{\ell} \rceil$ ciphertexts. For instance, if $mn < \ell$, this reduces the number of ciphertexts for A from n to 1.

3) *On the Number of Rots (O3)*: For multiplying with the ciphertexts of each row of B , one must replicate each element of A to a vector of dimension p theoretically. However, all segments are the same. So it does not need $mn \log p$ Rots, but only needs $mn \log \ell$ Rots. This observation shows that the required Rots is independent of p .

4) *On Replicate the Same Ciphertext (O4)*: In SegLKS, we need to replicate $a = (a_0, \dots, a_{m-1}) \in \mathcal{R}^m$ to $(a_i, a_i, \dots, a_i) \in \mathcal{R}^\ell$ for $i = 0, \dots, m-1$. This costs $m \log \ell$ Rots and m CMuls. However, to obtain the same results, one may first group a into groups with each group κ elements and repeat each group ℓ/κ times. Then, for each repeated ciphertext, replicate the κ elements. For instance, assume that our goal is to obtain ciphertexts of $(i, \dots, i) \in \mathcal{R}^{16}$ from $a = (1, 2, 3, 4) \in \mathcal{R}^4$ for $i = 1, 2, 3, 4$. First, we select $(1, 2) \in \mathcal{R}^2$ and repeat it 8 times to obtain $(1, 2, \dots, 1, 2)$, which costs 1 CMuls and $\log(\ell/\kappa) = 3$ Rots. Then we can obtain $(1, 0, \dots, 1, 0)$ and $(0, 2, \dots, 0, 2)$ by 2 CMuls. Then we can obtain $(1, \dots, 1)$ and $(2, \dots, 2)$ by 2 Rots. With this optimization, we can reduce the number of Rots from 16 to 10, at a cost of m/κ more CMuls and one more CMul depth. So, to replicate all elements of A , the required Rots and CMuls are bounded by $nm(\frac{\log(\ell/\kappa)}{\kappa} + \log \kappa)$ and $nm(1 + 1/\kappa)$ respectively. The minimum number of Rots achieves when κ satisfies $\kappa e^{\kappa-1} = \ell$.

Summary: We summarize in Table V the segmented algorithms. If $\lceil \frac{np}{\ell} \rceil$ is small, say a constant, the efficiency of BMM-III is relevant since the required number of operations

is only linear in m . In the case of $n, m \ll p$, the advantage of SegLKS is evident, as the required number of Rots and CMuls is independent of p . Furthermore, we also note that the number of Rots required by SegLKS is even less than that of the state-of-the-art algorithm [30], which needs $\ell^{-\frac{2}{3}} d^2 \log d$ Rots with $d = \max\{n, m, p\}$, asymptotically more than $\log \ell \cdot d^2$ (SegLKS/O3) when d tends to infinity.

VII. IMPLEMENTATION AND EVALUATION

In this section, we first introduce our implementation with more details, followed by a comprehensive experimental study of the involved algorithms to evaluate their performance.

A. Implementation

To evaluate the performance of the algorithms introduced in this paper, we implement them all (BMM-I, BMM-II, BMM-III, SegLKS) and the algorithms of [36] and [38] by using the CKKS scheme [6] implemented in Microsoft's open-source SEAL [15]. We also implement the naïve (text-book) and Strassen block version of the algorithms in [36], BMM-I, and BMM-II. Our implementation is open-sourced at https://github.com/hangenba/bicyclic_mat_mul.

In SEAL, the key-switching keys for Rot_k are generated only for $k = 2^i$ by default. When a Rot_k operation is involved for a non-power-of-two k , the task can be accomplished using consecutive rotations, e.g., $\text{Rot}_3(\text{ct}.x) = \text{Rot}_1(\text{Rot}_2(\text{ct}.x))$. This is the main reason for the time-consuming nature of Rot. For efficiency, we pre-generate all the necessary switching keys in our implementation. Although this needs additional time and storage overhead for key generation, according to our test with BMM-III for $N = 8192$, it saves about 15% of computing time. Once these keys are generated, they can be reused in subsequent operations.

There are many fast algorithms for matrix multiplication in plain, e.g., [31] and [54]. Almost all of these algorithms are recursive. It is well known that a recursive program can always be rewritten as an iterative loop. However, in the context of ciphertext computation, such as encrypted Strassen's matrix computation, the multiplication depth of the rewritten version may depend on the recursion depth, which will slow down the algorithm significantly. In contrast, the original recursive algorithm might only require a single Mul depth (possibly plus one or two CMul depths), although it is memory-consuming and hard to optimize further.

B. Setup

All experiments are run on a *single* thread (no parallelization) of a Desktop with an Intel Core i9-12900K at 3.2 GHz and 32 GB memory. We evaluate encrypted matrix multiplication algorithms with different dimensions:

- Small (almost) square: all matrices A , B and X can be packed into one ciphertext. This setting is to evaluate BMM-I, BMM-II, and algorithms in [36] and [38].
- Large (almost) square: all matrices A , B and X have to be packed into several ciphertexts. This setting is to evaluate BMM-III, together with the block version of Jiang et al.'s algorithm [36] and BMM-I.

TABLE VI
PERFORMANCE COMPARISON WITH THE RT ALGORITHM FOR
SMALL-DIMENSIONAL MATRICES. $N = 8192$ AND $\log \Delta = 30$

Method	$\log q$	Dimension	Time (ms)	Speedup
RT [38]	170	(16, 16, 16)	199	1.0x
BMM-I	170	(16, 19, 17)	130	1.5x
BMM-I	140	(16, 19, 17)	82	2.4x
BMM-II	140	(15, 16, 17)	13	16.6x

- Rectangular: Some of A , B and X are rectangular. This setting is to compare the performance of BMM-III, SegLKS and the state-of-the-art algorithm [14].

Recall that the ciphertext of CKKS is $\mathbb{Z}[X]/\langle X^N + 1, q \rangle$, and there is a scale factor Δ in CKKS related to the precision of the computed results. In our tests, all matrix entries are set by $\text{pow}(-1, i + j) * \text{rand}() / \text{pow}(2, 30)$, roughly in the interval $(-2, 2)$. The degree of ciphertext space and ciphertext modulus may vary depending on algorithms. In particular, we always set the bit-size of q as $50 + L \cdot \log \Delta + 60$, where L is the number of multiplicative depth (Mul and CMul) required by the corresponding algorithm (see, e.g., Table I).

Such a setting always achieves at least 128-bit security level according to the latest lattice estimator [55]. We also note that multiple test runs show that in this setup, the maximal absolute errors of the computed results are always less than 10^{-2} . All timings include encryption, computation, and decryption.

C. Matrices With Small Dimension

For CKKS-based encrypted matrix multiplication with small dimension (each of the matrices A , B and $X = AB$ can be encoded and encrypted into a single ciphertext), we compare BMM-I and BMM-II with the algorithms in [36] and [38].

When $N = 8192$, the Rizomiliotis-Triakosia (RT) algorithm [38] supports square matrix multiplication of dimension $\sqrt[3]{N/2} = 16$. As indicated in Table VI, BMM-I achieves a 1.5x speedup under the same parameter settings as the RT algorithm. However, due to the lower multiplicative depth required by BMM-I, it can finish the computation with a smaller ciphertext modulus. Under this setting, compared with the RT algorithm, BMM-I eventually achieves a 2.4x speedup. In the same parameter setting, BMM-II supports a (15, 16, 17) matrix multiplication, which performs the best among these algorithms, achieving a 16.6x speedup. When $N = 16384$ or $N = 32768$, the RT algorithm does not have corresponding matrix multiplication, whereas both BMM-I and BMM-II have. More precisely, the RT algorithm can support matrix multiplication of other dimensions at a cost of more ciphertext operations. The reason is that the rotation operation Rot is designed for vectors of dimension ℓ . If the dimension is not exact ℓ , then one Rot can be finished with two Rots, two CMuls and one Add. The same reason holds for Jiang et al.'s algorithm when $N = 16384$ in Table VII.

We test three different N in Table VII. For the same N , BMM-I demonstrates a clear advantage (achieving a 4.4x speedup at least) compared with Jiang et al.'s algorithm [36], though the matrix dimensions it supports are only about $\sqrt{1/2} \approx 70\%$ of that supported by Jiang et al.'s algorithm.

TABLE VII
PERFORMANCE COMPARISON WITH [36] FOR MATRICES WITH
SMALL DIMENSION, WHERE $\log \Delta = 30$. “—” INDICATES THAT THE
ALGORITHM DOES NOT SUPPORT MATRIX MULTIPLICATION
FOR THE DIMENSIONS

Method	$\log q$	$N = 8192$		$N = 16384$		$N = 32768$	
		Dim.	Time (ms)	Dim.	Time (ms)	Dim.	Time (ms)
[36]	200	(64, 64, 64)	1453	—	—	(128, 128, 128)	13526
BMM-I	140	(43, 45, 44)	284	(61, 64, 63)	1003	(89, 91, 90)	3059
BMM-II	140	(15, 16, 17)	13	(21, 16, 23)	31	(31, 16, 33)	48

TABLE VIII
PERFORMANCE COMPARISON FOR (128, 128, 128) MATRIX
MULTIPLICATION. $\log \Delta = 30$ EXCEPT FOR BMM-III
WITH $\log \Delta = 40$

Method	$\log q$	$N = 8192$		$N = 32768$	
		Basic block	Time (s)	Basic block	Time (s)
Naïve block [36]	200	(64, 64, 64)	11.34	(128, 128, 128)	14.17
Strassen + [36]	200	(64, 64, 64)	10.34	(128, 128, 128)	14.17
Naïve block BMM-I	140	(43, 45, 44)	7.59	(86, 89, 87)	13.32
Strassen + BMM-I	140	(32, 35, 33)	8.31	(64, 67, 65)	11.99
Naïve block BMM-II	140	(15, 16, 17)	4.40	(21, 32, 23)	8.32
Strassen + BMM-II	140	(11, 8, 9)	84.89	(17, 16, 19)	127.69
BMM-III	190	(128, 131, 129)	11.05	(128, 131, 129)	41.60

In fact, the matrix dimension supported by all these algorithms are constrained by the number of plaintext slots ℓ . For instance, for $N = 8192$ ($\ell = 4096$), the algorithm in [36] can support square matrix multiplication of dimension $\sqrt{\ell} = 64$, the algorithm in [38] is limited to $\sqrt[3]{\ell} = 16$, and BMM-I (resp. BMM-II) can support matrix multiplication of dimensions (43, 45, 44) (resp. (15, 16, 17)). However, the dimensions supported by BMM-I and BMM-II are quite flexible. For example, when $N = 16384$, BMM-I can support matrix multiplication of dimensions (61, 64, 63) or others, say (29, 128, 31) with a runtime of 1940 ms when $\log q = 140$, while the other two algorithms do not support square matrix multiplication under this parameter setting.

Comparison with the Algorithm in [30]: Since we did not implement our algorithms for the BGV scheme, we cannot directly compare the performance with Zheng et al.'s algorithm [30]. As analyzed previously, the performance of BMM-II should be a bit better than that of theirs. In particular, according to [30, Table 2], it is reported that their algorithm costs 1 Mul, 2 CMuls and 14 Rots for a (16, 16, 16) matrix multiplication (their implementation costs about 119ms on our desktop for the first example given in [30, Tab. 3]), while BMM-II requires only 1 Mul and 13 Rots for a (15, 16, 17) matrix multiplication in Table VII.

D. Large Square Matrix Multiplication

To evaluate the performance of algorithms for square matrices of high dimensions, we assume a scenario involving matrix multiplications with dimensions of 128, 256, 512 and 1024, respectively. For these tasks, we test different algorithms, including the naïve and Strassen version of block matrix multiplication, and the segmented BMM-III; see Tables VIII–X for details. Since it follows from Table VI that the algorithm in [38] is not comparable with BMM-I, we here only consider

TABLE IX
PERFORMANCE COMPARISON FOR (256, 256, 256) MATRIX
MULTIPLICATION. $\log \Delta = 30$ EXCEPT FOR ALGO. 8
WITH $\log \Delta = 40$

Method	$\log q$	$N = 8192$		$N = 32768$	
		Basic block	Time (s)	Basic block	Time (s)
Naïve block [36]	200	(64, 64, 64)	89.33	(128, 128, 128)	112.01
Strassen + [36]	200	(64, 64, 64)	71.54	(128, 128, 128)	97.90
Naïve block BMM-I	140	(43, 45, 44)	59.57	(86, 89, 87)	73.35
Strassen + BMM-I	140	(32, 35, 33)	56.98	(64, 67, 65)	80.99
Naïve block BMM-II	140	(15, 16, 17)	76.60	(21, 32, 23)	76.19
BMM-III	190	(256, 259, 257)	42.90 [†]	(256, 259, 257)	110.99

[†] We pre-generate the key-switching keys. Otherwise, it takes about 50s.

the Jiang et al.'s algorithm [36], BMM-I, and BMM-II as the base algorithms for the block version.

1) *Memory and the Number of Ciphertexts*: In Table VIII–X, we only list the dimensions of the involved basic block, from which, together with the input dimensions, the number of blocks and hence the number of ciphertexts can be determined. For example, each matrix is split into 2×2 blocks for Naïve block [36] in Table VIII with $N = 8192$, and hence each matrix is encrypted as 4 ciphertexts; while for the naïve block BMM-II with the same N , the second matrix is split into 8×8 blocks, and hence encrypted as 64 ciphertexts. Although using BMM-I or BMM-II for high-dimensional matrix multiplication usually generates more ciphertexts than, e.g., Jiang et al.'s [36], which makes these algorithms require more memory than that of [36], those algorithms based on BMM-III do not have this disadvantage. For the last row in Table VIII as another example, each matrix is encrypted as 5 ciphertexts, only 1 more than that of [36]. Once again, we note that the timings include all for encryption, computation, and decryption.

2) *Timings*: From Table VIII, it is evident that for (128, 128, 128) matrix multiplication, the most efficient is BMM-II using the naïve block method, while surprisingly, the Strassen version of BMM-II is the worst. In fact, this is consistent with the previous analysis, as although the basic version of BMM-II is theoretically the best among these algorithms, it supports the smallest matrix dimensions, resulting in more blocks and thus affecting the performance. Therefore, we will not consider the Strassen version of BMM-II for further tests.

As indicated in Table IX, compared with the naïve block version, the Strassen block version does provide a speedup. In addition, the Strassen block version of BMM-I is more efficient than that of Jiang et al.'s. However, it should be noted that the Strassen block version of BMM-I requires more ciphertexts. For instance, when $N = 8192$, the Strassen block version of BMM-I needs 64 ciphertexts to store a matrix, while the Strassen block version of Jiang et al.'s algorithm requires only 16. Table IX also shows that the segmented algorithm BMM-III with $N = 8192$ outperforms all block algorithms, each matrix stored in 17 ciphertexts due to the coprime limitation of the dimensions. For example, it can finish a (256, 259, 257) encrypted matrix multiplication within 42.90 seconds (2x faster than naïve block algorithm in [36]), of

TABLE X
PERFORMANCE COMPARISON FOR MATRIX MULTIPLICATION OF
DIMENSION 512 AND 1024 WITH $N = 8192$

Method	$\log q$	(512, 512, 512)		(1024, 1024, 1024)	
		Basic block	Time (s)	Basic block	Time (s)
Naïve block [36]	200	(64, 64, 64)	728	(64, 64, 64)	6028
Strassen + [36]	200	(64, 64, 64)	479	(64, 64, 64)	3514
Naïve block BMM-I	140	(43, 45, 44)	490	(43, 45, 44)	4240
Strassen + BMM-I	140	(32, 35, 33)	390	(32, 35, 33)	2757 [†]
Naïve block BMM-II	140	(15, 16, 17)	1766	(15, 16, 17)	—
BMM-III	190	(512, 515, 513)	181 [†]	(1024, 1027, 1025)	1200 [†]

[†] By default, $\log \Delta = 30$ except for these with $\log \Delta = 40$.

TABLE XI
PERFORMANCE FOR RECTANGULAR MATRIX MULTIPLICATION (I)

Huang et al.'s algorithm [14] [†]		BMM-III [‡]			
Dimension	Time (s)	Dimension	KeyGen (s)	Total (s)	Speedup
(256, 256, 16)	6.19	(256, 257, 17)	16.15	23.60	
(256, 16, 256)	6.23	(256, 17, 257)	14.69	19.37	
(1024, 1024, 16)	108.22	(1024, 1025, 17)	14.68	55.79	1.9x
(1024, 16, 1024)	108.31	(1024, 17, 1025)	14.47	43.64	2.4x
(2048, 2048, 8)	218.09	(2048, 2049, 11)	14.49	98.01	2.2x
(2048, 8, 2048)	218.09	(2049, 8, 2051)	14.63	81.09	2.6x

[†] For the algorithm in [14], N is set as the same as theirs and $\log \Delta = 30$.

[‡] For BMM-III, $N = 8192$ and $\log \Delta = 40$.

which 15.80 seconds are spent generating the switching keys required for Rot. Indeed, once these keys are generated, they can be reused, hence possibly further improving the efficiency for subsequent computations in practical applications.

Based on observations derived from Table VIII and IX, the setting $N = 8192$ performs better than $N = 32768$ for all tested algorithms. Thus, in tests with dimensions of 512 and 1024, we fix $N = 8192$. From Table X, the naïve block BMM-II perform the worst, since, again, more blocks slow down the speed. In addition, we have a similar observation to that from Table IX, indicating that BMM-III has the best performance among these algorithms. Specifically, for the task of (1024, 1024, 1024) encrypted matrix multiplication, BMM-III is 5x faster than the naïve block version Jiang et al.'s algorithm [36]. By the way, We also test the performance of the naïve block BMM-III, which costs about 260s for the (512, 512, 512) case with the basic block as (256, 259, 257).

E. Rectangular Matrix Multiplication

Huang et al. in [14] investigated high-dimensional encrypted rectangular matrix multiplication for different shapes. As indicated in Table XI, for the different dimensions in [14], as the dimensions increase, the efficiency of BMM-III gradually outperforms that of Huang et al.'s algorithm, with a 2.6x speedup at most. Furthermore, BMM-III can be used to cases of even larger dimension. For instance, it can finish a (8191, 8192, 15) encrypted matrix multiplication within 1452.16 seconds, while Huang et al.'s algorithm in its current setup does not support this instance.

Tables XII and XIII include two additional cases of rectangular encrypted matrix multiplication not discussed in [14]. In the first case (Table XII), matrix A is short and wide, while matrix B is tall and narrow, i.e., $n, p \ll m$. The naïve block version of BMM-I exhibits the best performance, about

TABLE XII
PERFORMANCE (IN SEC.) FOR RECTANGULAR MATRIX
MULTIPLICATION (II). $N = 8192$, $\log \Delta = 30$

Dimension	Naïve block [36]	Naïve block BMM-I	Speedup
(4, 1636, 5)	36.92	9.76	3.7x
(8, 3405, 9)	78.58	19.92	3.9x
(16, 6903, 17)	157.00	38.57	4.0x
(32, 13847, 33)	317.31	76.73	4.1x

TABLE XIII
PERFORMANCE (IN SEC.) FOR RECTANGULAR MATRIX
MULTIPLICATION (III). $N = 8192$, $\log \Delta = 30$

Dimension	Naïve block [36]	Naïve block BMM-I	SegLKS+O4
(4, 5, 1636)	36.31	0.10	0.25
(8, 9, 3405)	77.56	0.65	0.63
(16, 17, 6903)	157.32	4.96	2.23
(32, 33, 13847)	316.76	38.82	8.24

4x faster than the corresponding version of Jiang et al.'s algorithm. Note that BMM-III is a bit faster than the naïve block version of Jiang et al.'s, but slower than that of BMM-I. So, we omit its performance here. Note that the naïve block BMM-II applies to the tests in Table XII. It costs 0.20, 1.52, 23.02, and 513.07 seconds, respectively.

In Table XIII, we report experimental results for another case, i.e., $n \approx m \ll p$ and $mn < \ell$ (for this case the LongRot algorithm and hence the BMM-III algorithm do not work), for which SegLKS demonstrates the best performance. This is because the dimensions (n, m) of the matrix A are relatively small, and the number of Rots and CMuls required by SegLKS are independent of the number of columns p of the matrix B (as already indicated in Section VI-B). We implement all optimizations in Section VI-B. Experiments show that SegLKS with O4 is about 3x faster than that without O4. Thus, we only list the performance of SegLKS with O4 in Table XIII. It follows from Table XIII that the naïve block version of BMM-I is significantly faster than the naïve block version of Jiang et al.'s, but slower than SegLKS with O4 (in Table V $\kappa = 8$ is fixed for $\ell = 4096$). In particular, for the case of (32, 33, 13847), SegLKS with O4 is about 38x faster than the naïve block version of [36].

F. Comparison With Most-Recent Algorithms

In a recent work [41], Zhu et al. presented four different algorithms for different matrix dimensions (square and rectangular). In a very recent work [42], Gao and Gao presented a homomorphic encryption scheme named GMS for encrypted matrix multiplication. The two works considered similar cases and designed and implemented algorithms in HELIB. The code of the two works does not seem to be open-sourced. So we take the time from their paper directly. For a (128, 128, 128) encrypted matrix multiplication, Zhu et al.'s algorithm costs about 200 seconds (with a 2.5GHz CPU), while Gao and Gao's algorithm takes about 102 seconds (with a 2.6 GHz CPU). Note that the naïve block version of BMM-II costs only 4.4 seconds (see Table VIII). Similar observations can be made for other square or rectangular cases as well. While such significant performance differences may be attributed to the use

of different software libraries, different encryption schemes, or different parameters, the theoretical results in Table I show that BMM-I and BMM-II have their own advantages.

VIII. CONCLUSION

In this paper, we introduce bicyclic encoding, a novel method for matrix encoding. We design several new algorithms for encrypted matrix multiplication under bicyclic encoding, and investigate the block and segmented methods for handling high-dimensional matrices. In the context of CKKS, the following conclusions can be drawn from our theoretical analysis and comprehensive experimental study:

- For encrypted matrix multiplication of small dimensions, BMM-I or BMM-II is the optimal choice;
- When dealing with larger-scale encrypted square matrix multiplication, although theoretically, variants based on the Strassen block-wise strategy are faster, the practical performance is better with the segmented strategy (BMM-III);
- For those types of rectangular encrypted matrix multiplication discussed in [14], BMM-III shows better performance;
- For rectangular encrypted matrix multiplication with $n \approx p \ll m$, the naïve block BMM-I is effective;
- For rectangular encrypted matrix multiplication with $n \approx m \ll p$, SegLKS that combines Lu et al.'s algorithm [33] with the segmented strategy demonstrates exceptional performance.

From the perspective of practical application, the implementations in this paper can still be further optimized. For example, employing a multi-thread parallelization can yield additional acceleration; for block algorithms, the intermediate results after rotations can be reused, and hence further acceleration may be achieved; the hoisting and double-hoisting technique (e.g., [45]) may be used to accelerate our implementation further as well.

Furthermore, implementing these algorithms using BGV or B/FV schemes would significantly enhance integer matrix multiplication on encrypted data. Additionally, exploring other matrix operations beyond transpose and multiplication under bicyclic encoding are worth further investigation.

Finally, the error analysis for matrix multiplication on encrypted data is another intriguing direction, which is closely related to choosing parameters, e.g., the scale factor Δ for CKKS, or the plaintext modulus for BGV and B/FV.

ACKNOWLEDGMENT

The authors thank Zhicong Huang for sharing with them, together with Cheng Hong, Chenkai Weng, and Wenjie Lu, their codes for [14], and Xiaopeng Zheng for discussion with them on the multiplicative depth of the encrypted Strassen algorithm. They also thank an anonymous reviewer from ACM CCS '24 whose suggestions inspired them to design BMM-II, and other anonymous reviewers for their suggestions that helped present the paper more clearly. The work of Linhan Yang was done when he was visiting with CIGIT, CAS.

REFERENCES

- [1] R. Rivest, L. Adleman, and M. Dertouzos, "On data banks and privacy homomorphisms," in *Foundations of Security Computation*. New York, NY, USA: Academic, 1978, pp. 165–179.
- [2] C. Gentry, "Fully homomorphic encryption using ideal lattices," in *Proc. 41st Annu. ACM Symp. Theory Comput. (STOC)*, 2009, pp. 169–178.
- [3] Z. Brakerski, C. Gentry, and V. Vaikuntanathan, "(Leveled) fully homomorphic encryption without bootstrapping," *ACM Trans. Comput. Theory*, vol. 6, no. 3, pp. 1–36, Jul. 2014.
- [4] Z. Brakerski, "Fully homomorphic encryption without modulus switching from classical GapSVP," in *Proc. Annu. Cryptol. Conf.*, 2012, pp. 868–886.
- [5] J. Fan and F. Vercauteren. (2012). *Somewhat Practical Fully Homomorphic Encryption*. [Online]. Available: <https://ia.cr/2012/144>
- [6] J. H. Cheon, A. Kim, M. Kim, and Y. Song, "Homomorphic encryption for arithmetic of approximate numbers," in *Proc. Int. Conf. Theory Appl. Cryptol. Inf. Secur.*, 2017, pp. 409–437.
- [7] L. Ducas and D. Micciancio, "FHEW: Bootstrapping homomorphic encryption in less than a second," in *Proc. Annu. Int. Conf. Theory Appl. Cryptograph. Techn.*, Apr. 2015, pp. 617–640.
- [8] I. Chillotti, N. Gama, M. Georgieva, and M. Izabachène, "TFHE: Fast fully homomorphic encryption over the torus," *J. Cryptol.*, vol. 33, no. 1, pp. 34–91, Jan. 2020.
- [9] N. P. Smart and F. Vercauteren, "Fully homomorphic SIMD operations," *Designs, Codes Cryptography*, vol. 71, no. 1, pp. 57–81, Apr. 2014.
- [10] S. Halevi and V. Shoup, "Bootstrapping for HELib," *J. Cryptol.*, vol. 34, no. 1, pp. 641–670, Jan. 2021.
- [11] A. Kim et al., "General bootstrapping approach for RLWE-based homomorphic encryption," *IEEE Trans. Comput.*, vol. 73, no. 1, pp. 86–96, Jan. 2024.
- [12] F.-H. Liu and H. Wang, "Batch bootstrapping I: A new framework for SIMD bootstrapping in polynomial modulus," in *Proc. EUROCRYPT*, 2023, pp. 321–352.
- [13] B. Xiang, J. Zhang, Y. Deng, Y. Dai, and D. Feng, "Fast blind rotation for bootstrapping FHEs," in *Proc. CRYPTO*, 2023, pp. 3–36.
- [14] Z. Huang, C. Hong, C. Weng, W.-J. Lu, and H. Qu, "More efficient secure matrix multiplication for unbalanced recommender systems," *IEEE Trans. Dependable Secure Comput.*, vol. 20, no. 1, pp. 551–562, Jan. 2023.
- [15] (2023). *Microsoft SEAL (Release 4.1.1)*. Accessed: Microsoft. [Online]. Available: <https://github.com/microsoft/SEAL>
- [16] F. Boemer, R. Cammarota, D. Demmler, T. Schneider, and H. Yalame, "MP2ML: A mixed-protocol machine learning framework for private inference," in *Proc. ARES '20*, 2020, p. 14.
- [17] S. Fu, Y. Yu, and M. Xu, "A secure algorithm for outsourcing matrix multiplication computation in the cloud," in *Proc. 5th ACM Int. Workshop Secur. Cloud Comput.*, Apr. 2017, pp. 27–33.
- [18] J. Dumas et al., "Secure multiparty matrix multiplication based on Strassen-Winograd algorithm," in *Proc. IWSEC*, 2019, pp. 67–88.
- [19] L. Zhao and L. Chen, "Sparse matrix masking-based non-interactive verifiable (Outsourced) computation, revisited," *IEEE Trans. Dependable Secur. Comput.*, vol. 17, no. 6, pp. 1188–1206, Nov. 2020.
- [20] H. Chen, M. Kim, I. Razenshteyn, D. Rotaru, Y. S. Song, and S. Wagh, "Maliciously secure matrix multiplication with applications to private deep learning," in *Proc. Int. Conf. Theory Appl. Cryptol. Inf. Secur.*, 2020, pp. 31–59.
- [21] S. Balla and F. Koushanfar, "HELiKs: HE linear algebra kernels for secure inference," in *Proc. ACM SIGSAC Conf. Comput. Commun. Secur.*, Nov. 2023, pp. 2306–2320.
- [22] J. Zhu, S. Li, and J. Li, "Information-theoretically private matrix multiplication from MDS-coded storage," *IEEE Trans. Inf. Forensics Security*, vol. 18, pp. 1680–1695, 2023.
- [23] R. Hiromasa, M. Abe, and T. Okamoto, "Packing messages and optimizing bootstrapping in GSW-FHE," *IEICE Trans. Fundamentals Electron., Commun. Comput. Sci.*, vol. 99, no. 1, pp. 73–82, 2016.
- [24] Y. Bai, X. Shi, W. Wu, J. Chen, and Y. Feng, "SeIMC: A GSW-based secure and efficient integer matrix computation scheme with implementation," *IEEE Access*, vol. 8, pp. 98383–98394, 2020.
- [25] C. Gentry, A. Sahai, and B. Waters, "Homomorphic encryption from learning with errors: Conceptually-simpler, asymptotically-faster, attribute-based," in *Proc. CRYPTO*, 2013, pp. 75–92.
- [26] D. H. Duong, P. K. Mishra, and M. Yasuda, "Efficient secure matrix multiplication over LWE-based homomorphic encryption," *Tatra Mountains Math. Publications*, vol. 67, no. 1, pp. 69–83, Sep. 2016.
- [27] M. Yasuda, T. Shimoyama, J. Kogure, K. Yokoyama, and T. Koshiba, "Secure statistical analysis using RLWE-based homomorphic encryption," in *Proc. ACISP*, 2015, pp. 471–487.
- [28] M. Yasuda, T. Shimoyama, J. Kogure, K. Yokoyama, and T. Koshiba, "New packing method in somewhat homomorphic encryption and its applications," *Secur. Commun. Netw.*, vol. 8, no. 13, pp. 2194–2213, Sep. 2015.
- [29] P. K. Mishra, D. H. Duong, and M. Yasuda, "Enhancement for secure multiple matrix multiplications over ring-LWE homomorphic encryption," in *Proc. ISPEC*, 2017, pp. 320–330.
- [30] X. Zheng, H. Li, and D. Wang, "A new framework for fast homomorphic matrix multiplication," *Cryptology ePrint Archive*, vol. 1, p. 1649, Apr. 2023.
- [31] V. Strassen, "Gaussian elimination is not optimal," *Numerische Math.*, vol. 13, no. 4, pp. 354–356, Aug. 1969.
- [32] S. Halevi and V. Shoup, "Algorithms in HELib," in *Proc. Annu. Cryptol. Conf.*, 2014, pp. 554–571.
- [33] W.-J. Lu, S. Kawasaki, and J. Sakuma, "Using fully homomorphic encryption for statistical analysis of categorical, ordinal and numerical data," in *Proc. Netw. Distrib. Syst. Secur. Symp.*, 2017, pp. 1–12.
- [34] S. Wang and H. Huang, "Secure outsourced computation of multiple matrix multiplication based on fully homomorphic encryption," *KSII Trans. Internet Info Syst.*, vol. 13, no. 11, pp. 5616–5630, 2019.
- [35] D. Rathee, P. K. Mishra, and M. Yasuda, "Faster PCA and linear regression through hypercubes in HELib," in *Proc. Workshop Privacy Electron. Soc.*, Jan. 2018, pp. 42–53.
- [36] X. Jiang, M. Kim, K. Lauter, and Y. Song, "Secure outsourced matrix computation and application to neural networks," in *Proc. ACM SIGSAC Conf. Comput. Commun. Secur.*, Oct. 2018, pp. 1209–1222.
- [37] J. Jang et al., "Privacy-preserving deep sequential model with matrix homomorphic encryption," in *Proc. ACM Asia Conf. Comput. Commun. Security*, 2022, pp. 377–391.
- [38] P. Rizomiliotis and A. Triakosia, "On matrix multiplication with homomorphic encryption," in *Proc. Cloud Comput. Secur. Workshop*, Nov. 2022, pp. 53–61.
- [39] J. Chiang, "Volley revolver: A novel matrix-encoding method for privacy-preserving neural networks (inference)," 2022, *arXiv:2201.12577*.
- [40] H. Huang and H. Zong, "Secure matrix multiplication based on fully homomorphic encryption," *J. Supercomput.*, vol. 79, no. 5, pp. 5064–5085, Mar. 2023.
- [41] L. Zhu, Q. Hua, Y. Chen, and H. Jin, "Secure outsourced matrix multiplication with fully homomorphic encryption," in *Proc. ESORICS*, 2023, pp. 249–269.
- [42] J. Gao and Y. Gao, "GMS: An efficient fully homomorphic encryption scheme for secure outsourced matrix multiplication," *J. Supercomput.*, vol. 80, no. 18, pp. 26435–26461, Dec. 2024.
- [43] G. C. Fox, S. W. Otto, and A. J. G. Hey, "Matrix algorithms on a hypercube I: Matrix multiplication," *Parallel Comput.*, vol. 4, no. 1, pp. 17–31, Feb. 1987.
- [44] M. Babenko et al., "A comparative study of secure outsourced matrix multiplication based on homomorphic encryption," *Big Data Cognit. Comput.*, vol. 7, no. 2, p. 84, Apr. 2023.
- [45] J.-P. Bossuat, C. Mouchet, J. Troncoso-Pastoriza, and J.-P. Hubaux, "Efficient bootstrapping for approximate homomorphic encryption with non-sparse keys," in *Proc. Annu. Int. Conf. Theory Appl. Cryptograph. Techn.*, 2021, pp. 587–617.
- [46] X. Ma, C. Ma, Y. Jiang, and C. Ge, "Improved privacy-preserving PCA using optimized homomorphic matrix multiplication," *Comput. Secur.*, vol. 138, Mar. 2024, Art. no. 103658.
- [47] N. Genise, C. Gentry, S. Halevi, B. Li, and D. Micciancio, "Homomorphic encryption for finite automata," in *Proc. ASIACRYPT*, 2019, pp. 473–502.
- [48] D. Ö. Simsek and M. Cenk, "Faster secure matrix multiplication with the BGV algorithm," in *Proc. 16th Int. Conf. Inf. Secur. Cryptol.*, Oct. 2023, pp. 1–5.
- [49] J. Hu, F. Wang, and K. Chen, "Faster matrix approximate homomorphic encryption," *Comput. Standards Interface*, vol. 87, Jan. 2024, Art. no. 103775.
- [50] O. Regev, "On lattices, learning with errors, random linear codes, and cryptography," *J. ACM*, vol. 56, no. 6, p. 34, 2009.
- [51] V. Lyubashevsky, C. Peikert, and O. Regev, "On ideal lattices and learning with errors over rings," *J. ACM*, vol. 60, no. 6, pp. 1–35, Nov. 2013.

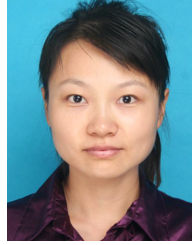
- [52] F. T. Leighton, *Introduction to Parallel Algorithms and Architectures*. San Mateo, CA, USA: Morgan Kaufmann, 1992.
- [53] J. Chen, L. Yang, W. Wu, Y. Liu, and Y. Feng. (2024). *Homomorphic Matrix Operations Under Bicyclic Encoding*. [Online]. Available: <https://ia.cr/2024/1762>
- [54] J. Alman and V. V. Williams, "A refined laser method and faster matrix multiplication," in *Proc. SIAM Symp. Discrete Algorithms*, 2021, pp. 522–539.
- [55] M. R. Albrecht, R. Player, and S. Scott, "On the concrete hardness of learning with errors," *J. Math. Cryptol.*, vol. 9, no. 3, pp. 169–203, Oct. 2015.



Wenyuan Wu received the Ph.D. degree from the Department of Applied Mathematics, University of Western Ontario, Canada, in 2007. He is currently a Professor, a Ph.D. Supervisor, and the Director of the Center for Automated Reasoning and Cognition, Chongqing Institute of Green and Intelligent Technology, CAS. His main research areas include artificial intelligence, automated reasoning, and privacy-preserving computing.



Jingwei Chen received the Ph.D. degree in computer science from the University of Chinese Academy of Sciences in 2013. He is currently an Associate Professor with Chongqing Institute of Green and Intelligent Technology, CAS. His research focus on symbolic computation, lattice-based cryptography, and privacy-preserving computing.



Yang Liu is currently an Associate Professor with Chongqing Jiaotong University. Her research areas include formal verification for software systems, artificial intelligence, and information security.



Linhan Yang is currently pursuing the master's degree with Chongqing Jiaotong University. His main research areas are homomorphic encryption and privacy computing.



Yong Feng is currently a Professor, a Ph.D. Supervisor, and the Director of the Institute of Electronic Information Technology, Chongqing Institute of Green and Intelligent Technology, CAS. His main research areas include automated reasoning and privacy-preserving computing.