

RESEARCH ARTICLE

Homomorphic Encryption on GPU

ALI ŞAH ÖZCAN¹, CAN AYDUMAN, ENES RECEP TÜRKOĞLU¹,
AND ERKAY SAVAŞ¹, (Member, IEEE)

Faculty of Engineering and Natural Sciences, Sabancı University, 34956 Istanbul, Turkey

Corresponding authors: Ali Şah Özcan (alisah@sabanciuniv.edu), Can Ayduman (canayduman@sabanciuniv.edu),

Enes Recep Türkoğlu (eturkoglu@sabanciuniv.edu), and Erkay Savaş (erkays@sabanciuniv.edu)

This work was supported in part by the Scientific and Technological Research Council of Turkey under Grant 118E725.

ABSTRACT Homomorphic encryption (HE) is a cryptosystem that allows the secure processing of encrypted data. One of the most popular HE schemes is the Brakerski-Fan-Vercauteren (BFV), which supports somewhat (SWHE) and fully homomorphic encryption (FHE). Since overly involved arithmetic operations of HE schemes are amenable to concurrent computation, GPU devices can be instrumental in facilitating the practical use of HE in real world applications thanks to their superior parallel processing capacity. This paper presents an optimized and highly parallelized GPU library to accelerate the BFV scheme. This library includes state-of-the-art implementations of Number Theoretic Transform (NTT) and inverse NTT that minimize the GPU kernel function calls. It makes efficient use of the GPU memory hierarchy and computes 128 NTT operations for ring dimension of 2^{14} only in 176.1 μs on RTX 3060Ti GPU. To the best of our knowledge, this is the fastest implementation in the literature. The library also improves the performance of the homomorphic operations of the BFV scheme. Although the library can be independently used, it is also fully integrated with the Microsoft SEAL library, which is a well-known HE library that also implements the BFV scheme. For one ciphertext multiplication, for the ring dimension 2^{14} and the modulus bit size of 438, our GPU implementation offers **63.4** times speedup over the SEAL library running on a high-end CPU. The library compares favorably with other state-of-the-art GPU implementations of NTT and BFV operations. Finally, we implement a privacy-preserving application that classifies encrypted genome data for tumor types and achieves speedups of 42.98 and 5.7 over CPU implementations using single and 16 threads, respectively. Our results indicate that GPU implementations can facilitate the deployment of homomorphic cryptographic libraries in real-world privacy-preserving applications.

INDEX TERMS Lattice based cryptography, homomorphic encryption, number theoretic transform (NTT), GPU, parallel processing, secure computation.

I. INTRODUCTION

Fully Homomorphic Encryption (FHE) enables computation over encrypted data, which had been considered as the most sought-after cryptographic primitive for many years. In [1], Gentry proposed the first functional FHE scheme, which is described over ideal lattices and permits the homomorphic evaluation of arbitrary circuits. Later, more practicable schemes based on learning with errors problem over rings (RLWE) [2] were proposed, where plaintext and ciphertext messages are represented as polynomials and ciphertext contains “noise,” which, increases

as homomorphic operations are applied. Thus, the scheme has a noise budget sufficient only for a certain number of homomorphic operations; and if noise reaches a certain limit, the homomorphic property will not hold and the ciphertext message does not decrypt due to excessive noise. This scheme is, thus, aptly called somewhat homomorphic encryption (SHE). To continue with the homomorphic operations, a technique referred as bootstrapping was proposed originally by Gentry [1], whereby the ciphertext is homomorphically decrypted to obtain a ciphertext with a replenished noise budget. This process can be applied repeatedly to obtain a fully homomorphic scheme, but bootstrapping is generally deemed to be a prohibitively expensive operation.

The associate editor coordinating the review of this manuscript and approving it for publication was Thomas Canhao Xu¹.

The first implementation of an FHE scheme was realized by Gentry and Halevi, as explained in [3]. Then, several FHE realizations were introduced such as those in [4] and [5]. One of the most promising approaches is the Brakerski-Fan-Vercauteren scheme [6], and there are several practical implementations of this and other similar schemes such as those provided by well-known software libraries SEAL [7], PALISADE [8], and HELib [9].

However, due to their compute-intensive operations in involved mathematical structures, current FHE implementations are far from being easily deployable in practice such as in large-scale practical cloud applications. Besides algorithmic optimizations and theoretical advances, using hardware accelerators is also the most viable option for bridging the gap between FHE performance and the requirements of real-world applications. GPU, FPGA and ASIC architectures can be profitably utilized as accelerators [10], [11], [12], to push the boundaries of FHE performance. A recently announced software library [13] provides support for hardware acceleration integration to software implementations of HE schemes using a standard Hardware Abstraction Layer (HAL).

In this paper, we present algorithms and implementation techniques to accelerate the BFV scheme of the SEAL library via NVIDIA GPUs. Our implementation developed in computing Unified Device Architecture (CUDA) program model [14] accelerates all homomorphic operations in the BFV scheme utilizing various parallelization strategies that can be applied on GPU architectures. To the best of our knowledge, ours is the first work, in which the entire SEAL BFV scheme (including addition, multiplication, relinearization, and rotation operations) can be offloaded onto GPU. We provide a GPU library for the BFV homomorphic encryption scheme, which can be used as a standalone application or integrated with the SEAL library to accelerate chosen homomorphic operations. The fully GPU-operated version of the library is publicly available on Github.¹ Our implementations used in library achieves a very high level of parallelization on GPU, targeting the compute-intensive nature of FHE operations.

It is established that multiplication in polynomial rings (which requires the multiplication of high-degree polynomials and their division by high-degree cyclotomic polynomials) is the most time and resource-critical operation in all FHE and SHE implementations.

Fortunately, there is a good deal of room for algorithmic research to accelerate the polynomial multiplication by utilizing the inherent parallelism in the operation and the hardware infrastructures (FPGA, ASIC, GPU) to exploit it in different ways. GPU architectures support many concurrent threads, which can be employed to perform the multiplication of very high-degree polynomials. Therefore, the noise budget can be made sufficiently large to homomorphically evaluate relatively complex circuits without having to use the bootstrapping method.

Here, we present a full GPU implementation for homomorphic operations of the BFV scheme and show that it can be used to accelerate real-world applications significantly. Our work introduces a new set of NTT implementation improvements for polynomial multiplication adapted to GPU architecture and proves to be the fastest in comparison to those reported in the literature to the best of our knowledge. We can summarize our contribution in this paper as follows:

- Our implementation uses only two GPU kernels, which minimizes the prohibitively slow global memory accesses independent of polynomial degrees used in this paper
- While the second GPU kernel is implemented using the conventional method as in many other works in the literature, we propose a novel algorithm for the implementation of the first GPU kernel (see Algorithm 6).
- In both kernels, global memory is accessed twice, one in the beginning and the other at the end of the kernel computations. In the first kernel, the input array is kept in the registers during the computations without any additional memory accesses (see Figure 9). As the registers are the fastest storage type, performing NTT computation via registers leads to significant acceleration in all homomorphic operations.

Our fast GPU implementation of polynomial multiplication can be used to accelerate existing implementations of other SHE and FHE schemes such as CKKS and BGV and other more involved homomorphic operations such as bootstrapping and scheme switching [15] as it is the most time consuming operation thereof. The latency of a single NTT operation of our implementation is also superior to those of other implementations in the literature; but, in particular, our NTT implementation is optimized for concurrent execution of many NTT operations, which is a typical use case scenario of all homomorphic operations and applications.

The rest of the paper is organized as follows. In Section II, we briefly explain the notation that we use and the mathematical background of NTT, Barrett reduction, Residue Number System (RNS), and homomorphic operations of the SEAL BFV scheme. In Section III, we present the essential working principle of the GPU architecture. In Section IV, we explain our GPU implementation and the algorithms that we used. In Section V, we discuss our implementation results and compare the results with state-of-the-art.

II. BACKGROUND

This section presents the notation used throughout the paper and explains the Barrett Reduction, Residue Number System, Number Theoretic Transform, and FHE operations of the SEAL BFV scheme.

A. NOTATION

The SHE scheme used in this work is BFV, one of the most efficient and widely used cryptographic schemes in the literature. The scheme is based on the ring learning with errors (RLWE) problem, whose difficulty serves as the secu-

¹https://github.com/Alisah-Ozcan/HE_GPU

urity assumption for some post-quantum cryptography and homomorphic encryption algorithms. The RLWE problem, or more precisely, learning with errors problem over rings, is a more efficient and practicable version of the learning with errors (LWE) problem, which is specialized to work with polynomial rings over finite fields, whose details are given below

The BFV scheme makes use of the polynomial ring $\mathbf{R}_q = \mathbb{Z}_q/\Phi(x)$, where \mathbb{Z}_q represents the finite ring $\{0, 1, \dots, q - 1\}$, in which the arithmetic is performed modulo q . Here, n is the degree of the cyclotomic polynomial $\Phi(x)$, and when its degree is selected as a power of two, we obtain $\Phi(x) = x^n + 1$. Then, the arithmetic in the ring \mathbf{R}_q is optimized as the polynomial division is performed with $x^n + 1$. Abusing the terminology we sometimes refer n as the dimension of \mathbf{R}_q and use the notation $\mathbf{R}_{q,n}$ to indicate its dimension.

Symbols and operations used in the subsequent parts of the paper are as follows: $\lceil \cdot \rceil$, $\lfloor \cdot \rfloor$, $\lceil \cdot \rceil$ represent round up, round down and round to nearest integer, respectively. The notation, $[a]_t$, indicates that the integer a lies in $[-t/2, t/2]$ while $|a|_t$ reduces a to the interval $[0, t - 1]$. A polynomial $a(x) \in \mathbf{R}_q$ can be treated as a vector of n integers in \mathbb{Z}_q , which is composed of its coefficients. When the number theoretic transformation (NTT), which is a form of discrete Fourier transformation over rings \mathbb{Z}_m (section III), is applied to the vector of $a(x)$, a vector of the same dimension is obtained, which is shown as $\bar{a}(x)$ (or just \bar{a}). While the symbols $+$, $-$ and \times (or just \cdot) represent addition, subtraction and multiplication, respectively in either \mathbb{Z}_q or \mathbf{R}_q the symbol \odot represents modular pointwise multiplication for vector representation of the elements of \mathbf{R}_q in the NTT domain. Namely, an element in a vector is multiplied by the elements of another vector with the same index value, where multiplications are in \mathbb{Z}_q (i.e., modulo q multiplication). λ is the security parameter denoted in unary notation. $a \leftarrow \mathbb{S}$ stands for the uniform sampling of a from the set \mathbb{S} . χ_{err} , a truncated zero-mean discrete Gaussian distribution, is used to sample the coefficients of error polynomials. The distribution is parameterized by the error bound β_{err} and standard deviation σ .

Now, we can give the most general and simplified definition of the RLWE problem. Suppose $a \leftarrow \mathbf{R}_q$ and the secret s and the error e are the elements of \mathbf{R} , whose coefficients are sampled from χ_{err} . Also suppose we have $b = as + e$. Then, the ‘‘search’’ RLWE problem can be defined as follows: Given a and b , it is hard to find s . In an HE scheme, s is the secret key whereas the (b, a) are the public key.

B. BARRETT REDUCTION

In the RNS variant of homomorphic cryptographic schemes such as [16], there is a multitude of modular multiplication operations that dominate the execution times of all homomorphic operations.

The Barrett reduction [17] and the Montgomery reduction [18] algorithms are two popular algorithms that perform the modular reduction operation efficiently. Since the Montgomery reduction needs the extra step for the transformation of integers to the Montgomery domain, the Barrett reduction algorithm is selected here for its simplicity.

The Barret reduction is described in Algorithm 1. Here, μ is the precomputed value, $\lfloor \frac{2^{2k}}{q} \rfloor$, where q is the modulus and k is the bit length of the modulus. The Barrett reduction algorithm includes multiplication, shift, and subtraction operations instead of an expensive division operation, which is needed in the computation of $C \bmod q$ by conventional modular multiplication algorithms.

Algorithm 1 Barrett Reduction

Input: $C = a \times b$, where $a, b < q$; $k = \lceil \log_2(q) \rceil$; $\mu = \lfloor \frac{2^{2k}}{q} \rfloor$

Output: $C_{out} (C \bmod q)$

- 1: $r \leftarrow C \gg (k - 2)$
 - 2: $r \leftarrow r \cdot \mu$
 - 3: $r \leftarrow r \gg (k + 2)$
 - 4: $r \leftarrow q \cdot r$
 - 5: $C_{out} \leftarrow (C - r)$
 - 6: **if** $C_{out} \geq 2q$ **then** $C_{out} \leftarrow C_{out} - 2q$
 - 7: **else if** $C_{out} \geq q$ **then** $C_{out} \leftarrow C_{out} - q$
 - 8: **else** $C_{out} \leftarrow C_{out}$
 - 9: **end if**
-

C. RESIDUE NUMBER SYSTEM (RNS)

An integer $X < M$, can be represented using residues x_i , where $x_i = X \bmod m_i$ for $i = 1, \dots, r$, if $M = \prod_{i=1}^r m_i$. Here, m_i s forms a set of pair-wise relatively prime integers that are known as moduli or ‘‘base’’ and a common notation is that $[X]_{m_i} = X \bmod m_i$. Due to the Chinese Remainder Theorem (CRT) we have

$$|X|_M = \left| \sum_{i=1}^r |x_i \cdot M_i^{-1}|_{m_i} \cdot M_i \right|_M,$$

where $M_i = \frac{M}{m_i}$. The RNS is preferred in cryptographic applications as it allows concurrent arithmetic with a set of small moduli in place of a big modulus; this is useful especially when the small moduli fit the word length of the underlying computing platform [19]. It is also showed [20], that RNS proves to be useful in accelerating the R-LWE based lattice-base somewhat homomorphic encryption schemes [6], [21]. Furthermore, RNS-variants of such schemes are proposed [16] and their implementations achieved good speedups on platforms where the concurrency of RNS is exploited [22].

D. NUMBER THEORETIC TRANSFORM (NTT)

The number theoretic transform (NTT) is a version of Discrete Fourier Transform (DFT) over the ring \mathbb{Z}_q . Any vector $a = [a_0, a_1, \dots, a_{n-1}]$ which has n elements in the polynomial domain can be transformed to another vector

$\bar{a} = [\bar{a}_0, \bar{a}_1, \dots, \bar{a}_{n-1}]$ which also has n elements in the NTT domain. The forward and inverse NTTs are defined as in Eqns 1 and 2:

$$\bar{a}_i = \sum_{j=0}^{n-1} a_j \omega^{i \times j} \pmod q \quad \text{for } i = 0, 1, \dots, n-1, \quad (1)$$

$$a_i = \frac{1}{n} \sum_{j=0}^{n-1} \bar{a}_j \omega^{-i \times j} \pmod q \quad \text{for } i = 0, 1, \dots, n-1. \quad (2)$$

The NTT (Eqn 1) and INTT (Eqn 2) calculations require the powers of a constant value $\omega \in \mathbb{Z}_q$ referred as the twiddle factor. Two types of twiddle factors are used:

- $\omega \in \mathbb{Z}_q$, which is the n -th root of unity in \mathbb{Z}_q and satisfies the conditions $\omega^n \equiv 1 \pmod q$ and $\omega^i \neq 1 \pmod q \forall i < n$, where $q \equiv 1 \pmod n$.
- ψ , where $\psi \in \mathbb{Z}_q$ is the $2n$ -th root of unity and it satisfies the conditions $\psi^{2n} \equiv 1 \pmod q$ and $\psi^i \neq 1 \pmod q \forall i < 2n$, where $q \equiv 1 \pmod n$. Note that ω and ψ are related with $\omega = \psi^2 \pmod q$ and $\psi^n \pmod q = -1$.

As the formulas in Eqns 1 and 2 result in quadratic complexity, for efficient computation of NTT and its inverse, Algorithms 2 and 3 are utilized [23].

Algorithm 2 Merge Forward NTT

Input: $a(x) \in \mathbb{Z}_q[x]/(x^n + 1)$ polynomial standard-order
Input: $\Psi_{br}[k] = \Psi^{br(k)}$ (Powers of Ψ stored in bit-reversed order)
Input: $n = 2^l, q (q \equiv 1 \pmod{2n})$
Output: $\bar{a} \in \mathbb{Z}_q^n$ in bit-reversed order

- 1: $t = n; m = 1$
- 2: **do**
- 3: $t = t/2$
- 4: **for** i from 0 by 1 to m **do**
- 5: $j_1 = 2it$
- 6: $j_2 = j_1 + t - 1$
- 7: **for** j from j_1 by 1 to $j_2 + 1$ **do**
- 8: $U = a_j$
- 9: $V = a_{j+t} \cdot \Psi_{br}[m+i] \pmod q$
- 10: $a_j = U + V \pmod q$
- 11: $a_{j+t} = U - V \pmod q$
- 12: **end for**
- 13: **end for**
- 14: $m = 2 \times m$
- 15: **while** $m < n$
- 16: **for** i from 0 by 1 to n **do**
- 17: $\bar{a}_i = a_i \pmod q$
- 18: **end for**

Both algorithms are based on the factorization of the cyclotomic polynomial $x^n + 1$ into n degree-1 polynomials as follows:

$$x^n + 1 \equiv \prod_{i=0}^{n-1} (x - \psi^{2i+1}) \pmod q \quad (3)$$

Algorithm 3 Merge Inverse NTT

Input: $\bar{a} \in \mathbb{Z}_q^n$ in bit-reversed order
Input: $\Psi_{rev}[k]$ (power of Ψ^{-1} stored in bit-reverse order ($\Psi_{rev}[k] = \Psi^{-br(k)} \pmod q$))
Input: $n = 2^l, q (q \equiv 1 \pmod{2n})$
Output: $a(x) \in \mathbb{Z}_q[x]/(x^n + 1)$ standard-order

- 1: $t = 1; m = n$
- 2: **do**
- 3: $j_1 = 0; h = m/2$
- 4: **for** i from 0 by 1 to h **do**
- 5: $j_2 = j_1 + t - 1$
- 6: **for** j from j_1 by 1 to $j_2 + 1$ **do**
- 7: $U = \bar{a}_j; V = \bar{a}_{j+t}$
- 8: $\bar{a}_j = U + V \pmod q$
- 9: $\bar{a}_{j+t} = (U - V) \cdot \Psi_{rev}[h+i] \pmod q$
- 10: **end for**
- 11: $j_1 = j_1 + 2 \times t$
- 12: **end for**
- 13: $t = 2 \times t$
- 14: $m = m/2$
- 15: **while** $m < n$
- 16: **for** i from 0 by 1 to n **do**
- 17: $a_i = (\bar{a}_i \cdot n^{-1}) \pmod q$
- 18: **end for**

By reducing a given polynomial $a(x)$ by these degree-1 polynomials, we obtain n integers, which are, in fact, the coefficients of $\bar{a}(x)$. This computation can be performed recursively. We first use the following factorization

$$(x^n + 1) \equiv (x^n - \psi^n) \equiv (x^{n/2} - \psi^{n/2})(x^{n/2} + \psi^{n/2}) \pmod q \quad (4)$$

and reduce $a(x)$ by polynomials $(x^{n/2} - \psi^{n/2})$ and $(x^{n/2} + \psi^{n/2})$. Reducing $a(x)$ by the first and second factors can be realized by employing the equations $x^{n/2} = \psi^{n/2}$ and $x^{n/2} = -\psi^{n/2}$, respectively. This accounts for the addition and subtraction operations in Steps 11 and 12 of Algorithm 2.

Factorization is further utilized as follows:

$$(x^{n/2} - \psi^{n/2}) \equiv (x^{n/4} + \psi^{n/4})(x^{n/4} + \psi^{n/4}) \pmod q$$

and

$$(x^{n/2} + \psi^{n/2}) \equiv (x^{n/2} - \psi^{n/2+n})(x^{n/4} - \psi^{n/4+n/2}) \times (x^{n/4} + \psi^{n/4+n/2}) \pmod q$$

The factorization is repeated until degree-1 polynomials are obtained.

As can be observed from Algorithm 2, different powers of ψ are stored in the bit-reverse order in the table Ψ_{br} , which simply means that the i th power of ψ is stored in the $(br(i) - 1)$ th element of Ψ_{br} . For instance, for $n = 8$ the first element of Ψ_{br} holds ψ^4 as the bit-reversed order of $4 = 100$ is 001 .

The inverse NTT operation, whose steps are given in Algorithm 3, is performed following the recursive factorization of $(x^n + 1)$ in the reverse order of that applied during the NTT computation.

To illustrate the inverse NTT algorithm, its last iteration is demonstrated, which yields the final result. The vector \bar{a} before the last iteration is as follows:

$$a = (a_0 + a_{n/2}\psi^{n/2}), \dots, (a_{n/2-1} + a_{n-1}\psi^{n/2}) \\ + (a_0 - a_{n/2}\psi^{n/2}), \dots, (a_{n/2-1} - a_{n-1}\psi^{n/2}) \quad (5)$$

If the first half is added to the second half, the first half of the resulting vector multiplied by 2 is obtained, $2(a_0, \dots, a_{n/2-1})$

Furthermore, if the second half of \bar{a} is subtracted from its first half,

$$2\psi^{n/2}(a_{n/2}, \dots, a_{n-1}). \quad (6)$$

is obtained. Thus, the result in Eqn 6 needs to be multiplied by $\psi^{-n/2}$. This elaborates the core butterfly operation in Step 12 of Algorithm 3.

As there are $\log_2 n$ iterations in the outermost loop of Algorithm 3 and the vector elements are effectively multiplied by 2 in every iteration, the result needs to be divided by n in \mathbb{Z}_q .

The schoolbook multiplication of polynomials $c(x) = a(x) \times b(x)$, where $a(x), b(x) \in \mathbf{R}_q$, can be performed using the method given in Eqn 7.

$$c(x) = \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} a_i \times b_j \times x^{i+j} \pmod q \quad (7)$$

Due to the quadratic complexity of the schoolbook method, the multiplication in \mathbf{R}_q is slow and inefficient. Moreover, the degree of the resulting polynomial is $2n - 2$ as a result of the multiplication, and thus division with $\phi(x)$ must be applied in order to obtain the final result, which is in \mathbf{R}_q ; i.e., a polynomial with degree at most $n - 1$.

The NTT-based polynomial multiplication operation, on the other hand, has logarithmic complexity. Recall that the vector \bar{a} is made up of scalar integers reduced by degree-1 polynomials that are factors of $(x^n + 1)$. This means that if two such vectors \bar{a} and \bar{b} , where $\bar{a} = NTT(a(x))$ and $\bar{b} = NTT(b(x))$, are multiplied element-wise in \mathbb{Z}_q , the result is \bar{c} in NTT, where $c(x) = a(x)b(x)$. When the inverse NTT is applied on \bar{c} , $c(x)$ is obtained.

Consequently, an NTT multiplication algorithm can be defined for an efficient multiplication in \mathbf{R}_q as described in Eqns. 8 and 9.

$$\bar{c}(x) = NTT_n(a(x)) \odot NTT_n(b(x)) \quad (8)$$

$$c(x) = INTT_n(\bar{c}(x)) \pmod q \quad (9)$$

Note that NTT operations are n -point and no extra polynomial reduction step by $(x^n + 1)$ is needed as $(x^n + 1)$ is factorized into degree-1 polynomials. This makes element-wise multiplications in Eqn. 8 isomorphic to the multiplication in \mathbf{R}_q .

E. SEAL BFV SCHEME

In this section, we first briefly explain four main operations of the BFV homomorphic encryption scheme; namely key generation, evaluation key generation, encryption, and decryption. Then, we give a more detailed explanation of three fundamental homomorphic operations that are common in many homomorphic cryptographic applications: addition, multiplication, and rotation over encrypted ciphertexts.

1) KEY GENERATION, ENCRYPTION AND DECRYPTION

For some integer $t > 1$, where $t \ll q$, the ciphertext and plaintext spaces are taken as $\mathbf{R}_{q,n}$ and $\mathbf{R}_{t,n}$, respectively. Also, we note that neither q nor t has to be a prime integer. The key generation, evaluation key generation, encryption, and decryption operations of the BFV scheme is shown below, where $\Delta = \lfloor q/t \rfloor$ and χ , ℓ , and w represent a discrete Gaussian distribution, the number of evaluation keys, and, the decomposition base, respectively.

- **Key Generation:** $a \leftarrow \mathbf{R}_{q,n}$, $s \leftarrow \mathbf{R}_{2,n}$ and $e \leftarrow \chi$, $sk = s$, $pk = (p_0, p_1) = (\lfloor -(as + e) \rfloor_q, a)$
- **Evaluation Key Generation:** $a_i \leftarrow \mathbf{R}_{n,q}$ and $e_i \leftarrow \chi$ for $j = 0, \dots, r - 1$; $i = 0, \dots, r - 2$ where $f^i = q_{r-1} \pmod{q_i}$ $(evk_i^j[0], evk_i^j[1]) = (\lfloor -(a_j s_j + e_j) + f^i s_j^2 \rfloor_{q_i}, a_j)$
- **Encryption:** $m \in \mathbf{R}_{t,n}$, $u \leftarrow \mathbf{R}_{2,n}$ and $e_1, e_2 \leftarrow \chi$, $ct = (c[0], c[1]) = (\lfloor m \cdot \Delta + p_0 u + e_1 \rfloor_q, \lfloor p_1 u + e_2 \rfloor_q)$
- **Decryption:** $ct = (c[0], c[1]) \in \mathbf{R}_{q,n}$ and $Sk \in \mathbf{R}_{2,n}$, $m = \lfloor \frac{t}{q} [c[0] + c[1]s]_q \rfloor_t$

Note that the evaluation keys, evk , are needed to remove the “nonlinear” parts $c[2]$ of the ciphertext $(c[0], c[1], c[2])$ that occur after homomorphic multiplication operations; a process often referred as relinearization. The number of evaluation keys are $2r(r - 1)$ in total. Note also that in the RNS variant of the BFV scheme, all operations have to be repeated for each prime base q_i .

2) ADDITION

In the BFV scheme, the most straightforward operations are addition and subtraction. It just consists of modular addition and subtraction of the coefficients of ciphertext polynomials that are in $\mathbf{R}_{q,n}$. As shown in Algorithm 4, two pairs of ciphertext polynomials in the same bases are added or subtracted coefficient-wise, where the moduli are q_i for $i = 0, \dots, r - 1$. Here, ct_i stands for the ciphertext pair in the modulus q_i for $i = 0, \dots, r - 1$; namely $ct_i = [ct]_{q_i}$ for ease of notation.

3) MULTIPLICATION

In this section, we explain the homomorphic multiplication operation as illustrated in Figure 1. As pointed out earlier in the RNS variant of the BFV scheme, a set of smaller moduli q_i is used instead of one large coefficient modulus q for the ring arithmetic; a technique known as residue number system (hence, the abbreviation RNS). Using RNS arithmetic allows to perform operations in parallel and removes the need for arbitrary-precision arithmetic.

Algorithm 4 BFV Addition

Input: $ct_i, \tilde{ct}_i \in \mathbf{R}_{q_i}$ for $0 \leq i < r - 1$
Output: $ct_i + \tilde{ct}_i \in \mathbf{R}_{q_i}$ for $0 \leq i < r - 1$

- 1: **for** i from 0 by 1 to $(r - 1)$ **do**
- 2: **for** k from 0 by 1 to 2 **do**
- 3: $\tilde{ct}_i[k] = [ct_i[k] + \tilde{ct}_i[k]]_{q_i}$
- 4: **end for**
- 5: **end for**
- 6: **return** $\tilde{ct} = \tilde{ct}_0, \dots, \tilde{ct}_{r-1}$

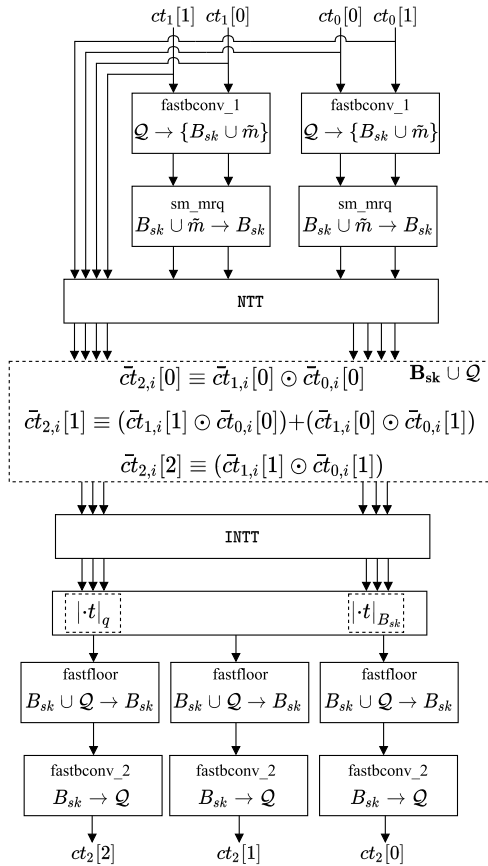


FIGURE 1. Homomorphic multiplication operation in the BFV scheme.

The homomorphic multiplication operation takes two ciphertexts as inputs, each of which consists of two polynomials in $\mathbf{R}_{q,n}$ and performs a tensor product that produces three polynomials as output in each RNS base.

Due to complications of using RNS arithmetic in homomorphic multiplication (see [16] for more details), the SEAL library uses the base extension technique and introduces additional auxiliary base (\mathcal{B} and m_{sk}) in addition to the RNS *thcalQ* base $\{q_0, q_1, \dots, q_{r-1}\}$. The auxiliary base \mathcal{B} consists of $\{B_0, B_1, \dots, B_{\rho-1}\}$, which are pairwise co-prime while m_{sk} is a prime integer. Generally, the auxiliary base \mathcal{B} and the prime m_{sk} are joined to form the base $\mathcal{B}_{sk} (= \mathcal{B} \cup m_{sk})$.

Thus, the homomorphic multiplication operation in BFV requires conversion between the \mathcal{Q} base and the auxiliary base \mathcal{B}_{sk} . The conversion is implemented using a technique

known as “fast base conversion,” which can introduce extra multiples of q in the computations that can lead to error in the ciphertext. To remedy this, a reduction operation through another modulus \tilde{m} is required after the fast base conversion operation is applied.

As shown in Figure 1, the BFV multiplication operation starts by performing the fast base conversion operation `fastbconv_1`, which convert the inputs in \mathcal{Q} to the base $\{\mathcal{B}_{sk} \cup \tilde{m}\}$. The `fastbconv_1` operation is followed by the reduction operation, for which the additional base \tilde{m} is used; this operation is known as small Montgomery reduction modulo q , `sm_mrq`. It limits the impact of the error and converts the inputs in the $\{\mathcal{B}_{sk} \cup \tilde{m}\}$ base to the \mathcal{B}_{sk} base. After the `sm_mrq` operation, the NTT operation is applied to all ciphertext components (both in \mathcal{B}_{sk} and \mathcal{Q} bases) and ciphertext multiplication operation is performed coefficient-wise to all vectors in all bases. Then, the inverse NTT operation is performed to convert the result to the polynomial domain. After the inverse NTT operation, ciphertexts are multiplied with plaintext modulus t . Then, the floor operation is used instead of the rounding operation; via a method is called “`fastfloor` function,” and convert the ciphertext in the base $\{q \cup \mathcal{B}_{sk}\}$ bases to the base \mathcal{B}_{sk} as it involves division by q . Finally, the `fastbconv_2` function is used to perform conversion from the \mathcal{B}_{sk} base back to the original RNS base \mathcal{Q} . The reader is referred to [16] for more detail.

4) RELINEARIZATION

The SEAL BFV uses the switchkey technique (Figure 2), which consists of the mix of three different methods for relinearization operation [24], [25], [26]. The most current method of these techniques is the special modulus method, which improves relinearization in terms of noise performance. The switch-key method shown in Figure 2 is the main building block of the relinearization and the rotation operations.

As shown in Figure 2, after the homomorphic multiplication, in addition to $ct[0]$ and $ct[1]$, the third ciphertext component $ct[2]$ is obtained. Recall that a ciphertext component is written in $r - 1$ moduli excluding q_{r-1} after encryption; ct_i for $i = 0, \dots, r - 1$. Firstly, all $ct_i[2]$ are transformed to the NTT domain using all moduli q_i in the RNS base to be multiplied with the evaluation keys that are already in the NTT domain.

The number of NTT operations is, therefore, $r(r - 1)$. After the NTT operations, the ciphertexts are multiplied with the evaluation keys in the NTT domain, where the multiplication is component-wise modulo multiplication. The modulus used in the multiplication is written next to the box that represents component-wise multiplication in Figure 2. Then, all results from the multiplication using the same modulus q_i in the RNS base are summed switchkey resulting vectors are transformed back to the polynomial domain using inverse NTT operation. Finally, as shown in Figure 2, necessary operations are applied to accommodate $ct[2]$ in $ct[0]$ and $ct[1]$. In the figure

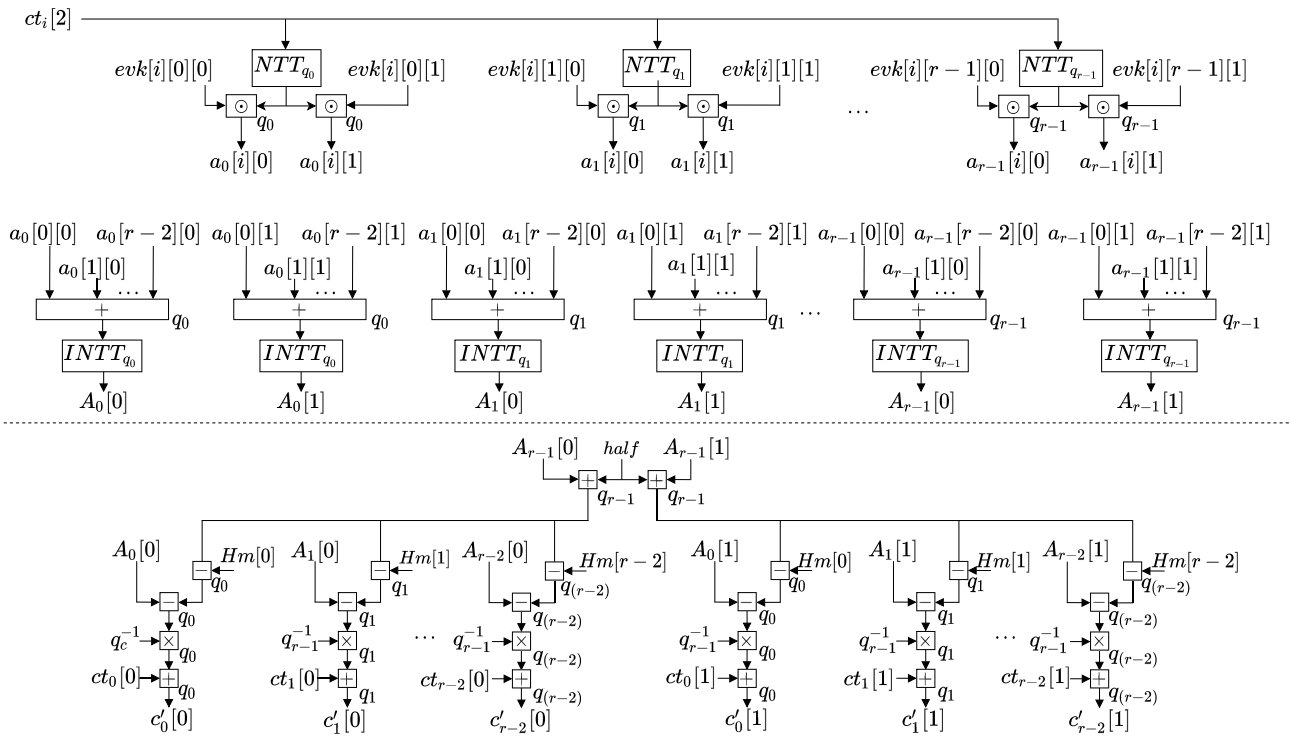


FIGURE 2. Switch key operation in the BFV scheme. The symbols $+$, $-$ and \times represent addition, subtraction and multiplication, respectively in either \mathbb{Z}_q or \mathbf{R}_q while the symbol \odot represents modular pointwise multiplication for vector representation of the elements of \mathbf{R}_q in the NTT domain.

the half mode $Hm[i] = \lfloor \lfloor q_{r-1}/2 \rfloor \rfloor_{q_i}$. See Algorithm 9 for the details.

5) ROTATION

The rotation operation also uses the switch-key operation as in the case of relinearization. However, the operation is based on Galois automorphism [27], and therefore, Galois keys are used for the switch-key operation. For each power of 2, there is a different set of Galois keys and if the rotation amount is a power of 2, the switch-key operation is executed using the corresponding Galois key. On the other hand, if the rotation amount is not a power of two, the amount is written as the combination of powers of two, and the switch-key operation is applied multiple times with different Galois keys. For instance, if the rotation amount is 10, it can be implemented using two switch-key operations; the former uses the Galois keys for 8, the latter for 2.

III. GPU ARCHITECTURE

A graphics processing unit (GPU) is a computing platform, which consists of many cores that can operate on many tasks concurrently, which makes it more suitable for parallel computations. On the other hand, GPU cores are much simpler than CPU cores and run at lower clock frequencies (*cf.* AMD Ryzen7 3800X’s cores working at up to 4.2 GHz and NVIDIA RTX 3060Ti’s cores working only at 1.66 GHz). Thus, GPUs become much more favorable for performing many simple tasks simultaneously. We will show in Section IV that many

time-critical BFV operations can be arranged as independent `for` loops, executed in GPU threads simultaneously.

One of the essential parts of GPUs is “streaming multiprocessors” (SMs); a unit of computing cores running the same GPU kernel. At the highest level of SMs, threads are combined as a 3-dimensional structure called blocks. Also, a grid is a group of blocks launched per GPU kernels. Using kernel launch parameters, one can determine the dimension of blocks and the number of threads per block as needed.

As shown in Figure 3, GPUs have different types of logical memory spaces; namely, shared memory (SMEM), registers, local memory (LMEM), constant memory (CMEM), texture memory (TMEM), global memory (GMEM). They have different sizes and different usages. For instance, GMEM has a large capacity; however, it has high access latency, especially in case of the low locality of access. As shown in Table 1, registers and SMEM are the fastest types of memory, and their read and write data speeds are similar to a typical L1 cache of a CPU. CMEM is a read-only data memory and since it is accessible from all threads, it performs well when multiple threads access the same data. Table 1 compares several aspects of GPU memory types [28].

IV. GPU IMPLEMENTATION

This section presents our implementation techniques, methods, and algorithms for four different homomorphic operations of the SEAL BFV scheme: addition, multiplication, relinearization, and rotation. Moreover, we present the implementation of our NTT and INTT algorithms on GPUs.

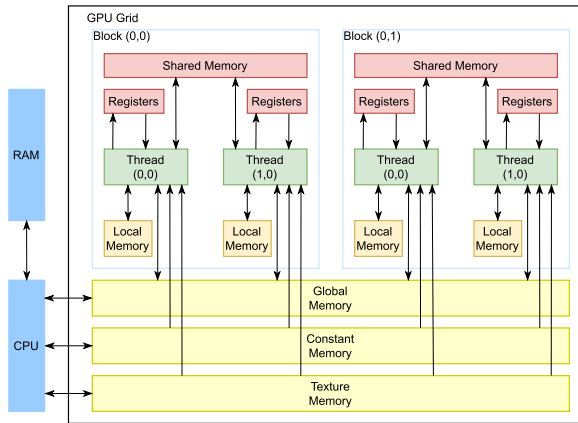


FIGURE 3. CUDA memory model.

TABLE 1. Variables and access penalties on modern GPUs memory architecture.

Variable Declaration	Memory	Life Time	Perform. Penalty
<code>int localVar;</code>	Register	Thread	1×
<code>int LocalArr[10];</code>	Local	Thread	100×
<code>__device__ int GVar;</code>	Global	Application	100×
<code>__constant__ int CVar;</code>	Constant	Application	1×
<code>__shared__ int SVar;</code>	Shared	Block	1×

In all GPU implementations that are performed in this section, we minimize the number of kernels to two. Here, our concern is due to the fact that transferring data from one kernel to another is only possible by using the global memory, and accessing the global memory is prohibitively expensive as explained in Section III. Note that the global memory is accessed only at the beginning and at the end of the kernel. Also, we designed the operations within the kernel in such a way that all required data sharing or exchange among the threads are facilitated only through shared memory, which is much faster than the global memory. The optimal use of global and shared memory decreased the total number of clock cycles spent in memory accesses and boosted memory throughput.

The implementation of the first kernel is involved as threads in different blocks may need to share data that requires inter-block synchronization (which calls for global memory access). In our implementation, we use only two blocks of threads and the number of threads in a block is optimized for the ring dimension. While this technique does not fully utilize the full concurrency in one loop of the NTT algorithm, it eliminates inter-block synchronization as the threads in the blocks always access the same memory locations. In addition, each thread in the blocks uses an access pattern to the global memory in such a way that, after global memory access, it keeps all its input values in its registers needed during the first kernel operations. Details of the particular algorithms and techniques are given in the subsequent sections.

Algorithm 5 Merge in Place Forward NTT on GPU (With `syncthreads`)

Input: $A[n], PsiTable[n], q$

Output: $A[n]$

```

1:  $Idx = blockIdx.x \times blockDim.x + threadIdx.x$ 
2: for loop from 0 by 1 to  $\log_2(n)$  do
3:    $t = (n/2) \gg loop$ 
4:    $m = 2 \ll loop$ 
5:    $address = int(Idx/t) * t + Idx$ 
6:    $U = A[address]$ 
7:    $V = A[address + t]$ 
8:    $Psi = PsiTable[int(Idx/t) + m]$ 
9:    $A[address] = (U + V) \% q$ 
10:   $V = (V \times Psi) \% q$ 
11:   $A[address + t] = (U - V) \% q$ 
12:  __syncthreads()
13: end for

```

A. NTT

The Cooley-Tukey NTT and Gentleman-Sande INTT algorithms described in the preliminary section is implemented on the GPU. NTT and INTT perform as the opposite of each other in terms of algorithm. Therefore, in this section, explanations of NTT and INTT separately are not needed. This section explains the challenges for fast and efficient implementation of NTT and presents our solutions to overcome them. Algorithm 5 shows the GPU pseudo-code for the NTT algorithm, which is essentially the same as the one given in Algorithm 2. One important adaptation to GPU is the synchronization operation in Step 12, whose effect on the correctness of the computations will be explained later in this section.

The NTT operation consists of $\log_2 n$ back-to-back loops, each of which contains $n/2$ butterfly operations independent of each other, which can be performed simultaneously using $n/2$ threads on the GPU.

Each GPU can run a certain number of streaming multiprocessors (SM), the number of which depends on the GPU model and computational capability (version) of the GPU. Each SM consists of 4 warps of 32 threads for all GPU models, so the total number of physical threads equals $(\#SM) \times 4 \times 32$.

When a GPU code is executed, the tasks are performed by warp groups. For example, if the code uses a number of threads in the range of [96-128], a total of four warps is needed in both cases. Also, even if the warps perform the same task, they may not finish their share of tasks simultaneously. Therefore, shared data usage among threads can lead to synchronization problems.

For instance, as the ring dimension n or the number of simultaneous (I)NTT operations increases, synchronization problems can occur if proper synchronization operations are not employed during the execution of Algorithm 5 (suppose Step 12 of Algorithm 5 is not present). This can be explained



FIGURE 4. Execution of Alg. 5 without synchronization in ideal circumstances where $n = 32$.

with a simple example. Suppose that we have a hypothetical GPU with a total thread count of 16 and a maximum block size of 4 threads. Let one warp of this GPU consist of two threads and let Algorithm 5 without synchronization be executed for $n = 32$ on this GPU. Figure 4 portrays a visualization of the execution of Algorithm 5 without synchronization on the hypothetical GPU.

The figure shows a total of $\log_2 32 = 5$ iterations. 16 threads are used, whose indexes are between 0 and 15 (T_0, \dots, T_{15}). We use a different color for the oval rectangle that encircles a block of four threads. Each thread T_i accesses two different memory locations using the address and $address+t$ values in Algorithm 5, performs the butterfly operation, writes two pieces of the results again in the same two memory locations. When a thread finishes its own task, it moves to the iteration of the algorithm and performs the same operation, only with a different value of t this time. Figure 4 represents the execution of the algorithm in the ideal circumstances as the threads are assumed to finish their tasks simultaneously.

However, Algorithm 5 does not always execute in ideal circumstances. Suppose that the number of threads is only 12 for the scenario in Figure 4. Even when the number of threads is less than the number of tasks, incidentally the execution can still be correct. One such scenario is depicted in Figure 5, where we only show the first two iterations. As the scenario requires 16 threads, but the hypothetical GPU has 12 threads, the number of threads is not sufficient, and the code runs sequentially after a point. In the figure, we use primed letters to distinguish the multiple assignments of the same thread to different tasks. For instance, T_0 and T'_0 show that performs thread executes two different butterfly operations in the same iteration sequentially. Incidentally again, this does not necessarily lead to incorrect execution as shown in Figure 5.

Nevertheless, since thread synchronization is not implemented, an error in calculations can occur as visualized in Figure 6. For instance, suppose four threads in the dashed red line, namely T'_2, T'_3, T'_6, T'_7 , are assumed to be scheduled simultaneously. And, since they operate on the same memory locations in two consecutive iterations, there is a data dependency between the first two and the last two threads. This will definitely lead to a race condition, resulting in incorrect results.

It is impossible to put a barrier between the warps to solve the aforementioned synchronization problem. Therefore, only block-level barriers can be used as shown in Step 12 of Algorithm 5, which resolves all synchronization problems as long as the ring dimension n is less than or equal to the block size. Since a block in GPU has a maximum of 1024 threads for all GPU models, the barrier `__syncthreads()` in Step 12 of Algorithm 5 cannot resolve the synchronization issue for higher values of n or when performing many NTT operations in batches.²

The latter issue can be explained over another execution scenario of Algorithm 5 on the hypothetical GPU, depicted in Figure 7. The eight threads enclosed in the dashed red line belong to two different blocks as the block size of the hypothetical GPU is just four. Here, the thread block in the 2nd iteration run on data that has not yet been completed, leading to incorrect results.

For n values much higher than the block size and the high number of multiple NTT operations running simultaneously, an obvious solution to resolve all synchronization issues is simply using more than one kernel depending on the size of n or the number of NTT computations. For example, for $n = 32$ on our hypothetical GPU, to resolve the synchronization

²When multiple and independent NTT operations are executed, the threads are scheduled as if those independent NTT calculations are combined into a single big NTT operation.

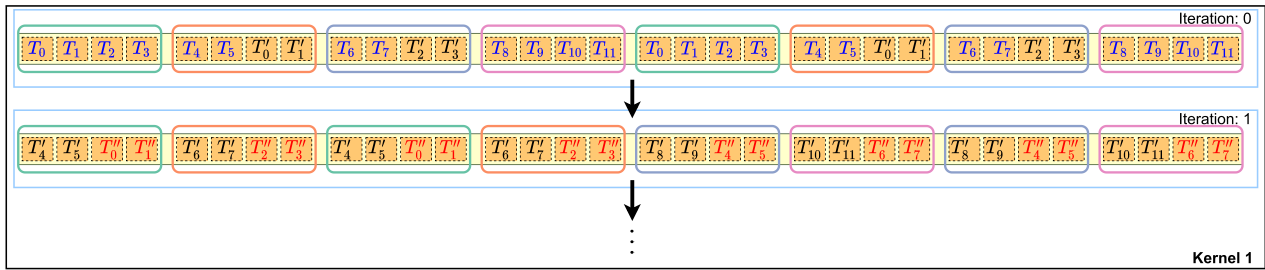


FIGURE 5. One good scenario for Alg. 5 without synchronization, where $n = 32$ and maximum block size = 4.

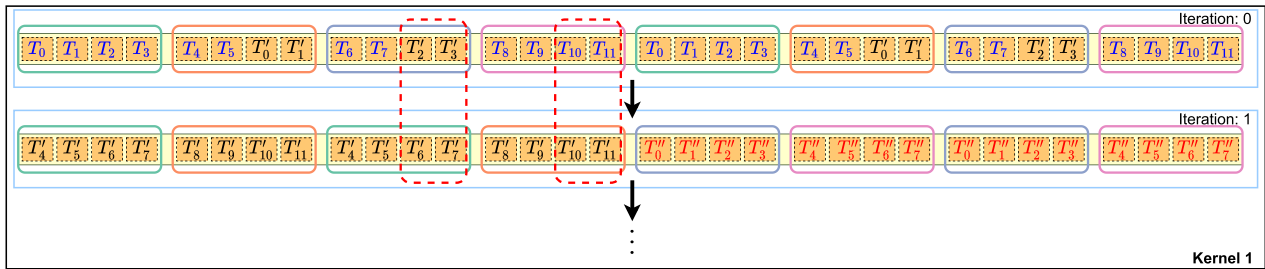


FIGURE 6. One problematic scenario for Alg. 5 without synchronization, where $n = 32$.

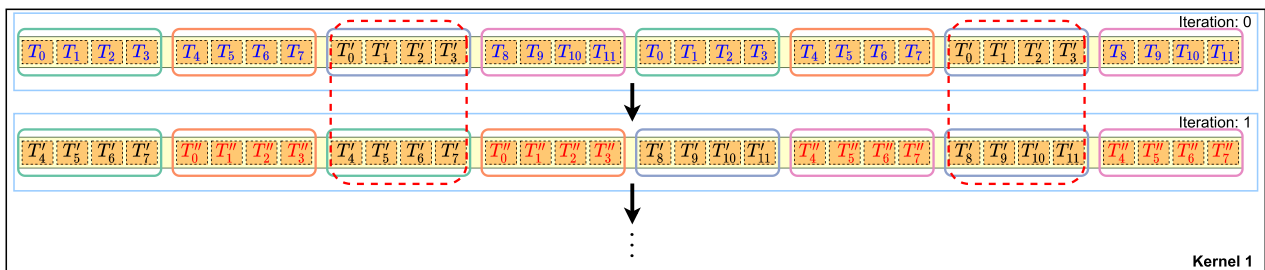


FIGURE 7. Another problematic scenario for Alg. 5, where $n = 32$.

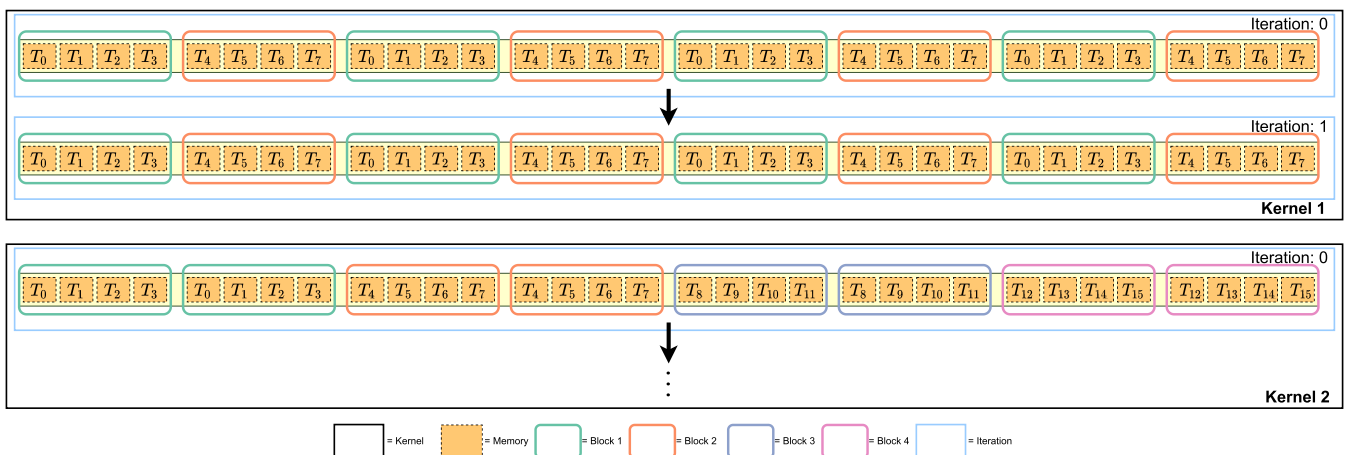


FIGURE 8. An example of our NTT algorithm where, $n = 32$ and maximum block size = 4.

issue in Figure 7, we can use two consecutively executing kernels for the first two iterations of the NTT computation.

After the first two iterations are completed, the threads in one block will never need or use the data processed by another block and no synchronization problem occurs. Therefore,

after the first two iterations, execution can continue within the third kernel using shared memory and block synchronization.

However, when more than one kernel is used, the only way to share data across kernels is to use global memory, which is the slowest of all GPU memory types (see Section III). As the

value of n increases, this method becomes prohibitively inefficient as the number of kernels required for NTT will also increase. This approach is used in [29], and as demonstrated in the subsequent sections, our new approach in this paper scales better as n increases.

The approaches described above have either synchronization issues or inefficient memory usage, both of which are efficiently addressed by our new NTT implementation. The new implementation consists of always two kernels for all values of n .

An example with $n = 32$ is visualized in the hypothetical GPU in Figure 8, where the first two iterations are performed in the first kernel. Operations of these iterations in the first kernel are performed sequentially on purpose. In the example in Figure 8, two blocks and eight threads (recall 2 blocks = 8 threads on the hypothetical GPU) are scheduled twice in the first two iterations. Each thread in the two blocks writes the addresses of interest in the global memory to its registers, as illustrated in Figure 9. Then, each thread performs butterfly operations using its register memory. When a thread finishes its task in one iteration, it writes the data in its register memory to the corresponding global memory. Although the first kernel seems to be slower because it uses fewer number of threads than the above mentioned examples, the acceleration here comes not from the number of threads, but from the more efficient usage of memory as we minimize the global memory access. On the other hand, in the approach employed in [29], the number of kernels along with global memory access increases as n becomes larger. The pseudo-code for the algorithm used to implement the operations in the first kernel is given in Algorithm 6.

In the second kernel in Figure 8, the number of threads in a block suffices to complete the remaining NTT operations. Since data sharing among threads within the block is required, each block has its shared memory consisting of $2 \times \text{blocksize}$. This poses no problem as the shared memory is the fastest type of GPU memory (in fact, as fast as the register memory). Here, each thread accesses its own part of the memory using its respective indexes for each iteration, performs a butterfly operation, and writes the result to the shared memory of the block it is connected to until the last iteration. After all threads finishes their executions, the result, which is in the shared memory, are written to these global memory, and the NTT operation is terminated. The new NTT implementation is fast and free of synchronization issues for all values of n and multiple concurrent NTT computations.

B. SEAL GPU IMPLEMENTATION

This section explains our GPU implementations of homomorphic addition, multiplication, relinearization and rotation operations of the BFV homomorphic encryption scheme. Algorithms for all these homomorphic operation are given as pseudo-codes as implemented in the Microsoft SEAL library. All of these algorithms are implemented so that they use our

Algorithm 6 Kernel 1 in Figure 8

Input: $A[n]$, $PsiTable[n]$, q
Input: bc : no. of blocks ($bc = 2$)
Output: $A[n]$

```

1:  $idx = blockIdx.x \times blockDim.x + threadIdx.x$ 
2:  $m = 1$ 
3:  $k = n / (2 \times blockDim.x \times bc)$ 
4:  $t = n$ 
5: for  $i$  from 0 to  $n / (2 \times blockDim.x)$  do
6:    $reg[i] = A[idx + (i \times (2 \times blockDim.x))]$ 
7: end for
8: for  $i$  from 0 to  $\log_2(n / (2 \times blockDim.x \times bc)) + 1$  do
9:   for  $j$  from 0 to  $n / (2 \times blockDim.x \times bc) - 1$  do
10:     $location = \lfloor \frac{j}{k} \rfloor \times k + j$ 
11:     $U = reg[location]$ 
12:     $V = reg[location + k]$ 
13:     $address = \lfloor \frac{idx}{t} \rfloor + m$ 
14:     $V = (V \times PsiTable[address]) \bmod q$ 
15:     $reg[location] = (U + V) \bmod q$ 
16:     $reg[location + k] = (U - V) \bmod q$ 
17:   end for
18:    $m = m \times 2$ 
19:    $k = k / 2$ 
20:    $t = t / 2$ 
21: end for
22: for  $i$  from 0 by 1 to  $n / (blockDim.x \times 2)$  do
23:    $A[idx + (i \times (blockDim.x \times 2))] = reg[i]$ 
24: end for

```

GPU implementation of the NTT algorithm as described in Section IV-A.

1) HOMOMORPHIC ADDITION/SUBTRACTION

As explained in Section II-E2, addition/subtraction operations of the BFV scheme are simple and inexpensive and their implementation consists of only one kernel. In this kernel, each ring element are represented as a vector over Z_{q_i} for each modulus in the RNS base, and modulo addition/subtraction is performed over the elements of the vectors.

2) HOMOMORPHIC MULTIPLICATION

In addition to kernel functions to implement NTT and INTT operations, ten different CUDA kernel functions are implemented for the multiplication operation (see Figure 1 for these operations). Each of the kernel functions use a one-dimensional block and thread indexing. Before the GPU computation, all necessary parameters are generated on CPU of the host computer, then sent to GPU. In what follows, we briefly mention all of them, but provide pseudo-codes for some important ones in case they are more involved.

The first two CUDA kernel functions are employed to implement base conversion operation from the RNS base \mathcal{Q} to \mathcal{B}_{sk} . The pseudo-code of the base conversion operation is given in Algorithm 7, as it is implemented in the Microsoft

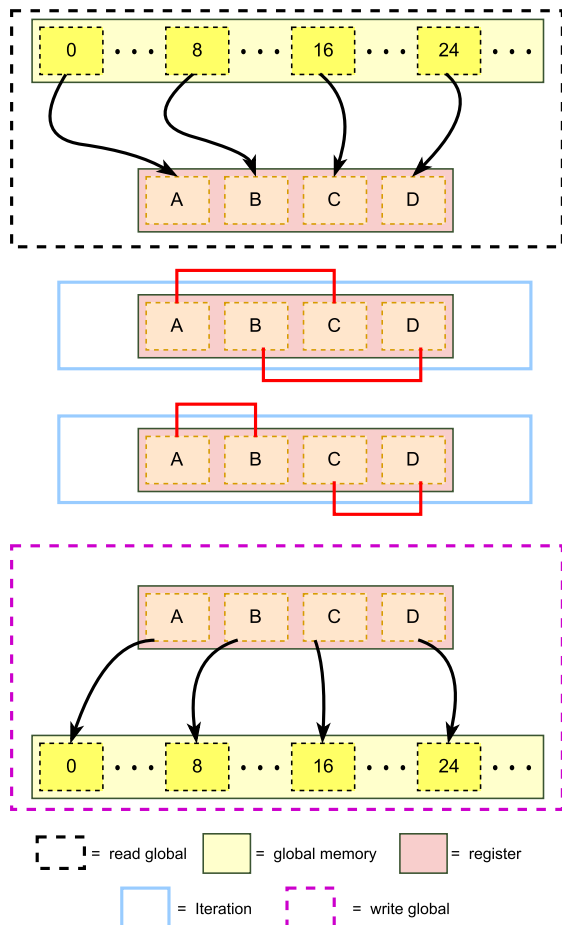


FIGURE 9. Register memory usage in our NTT algorithm for $n = 32$.

SEAL library. In particular, the first kernel implements the for loop in lines 1-3 of Algorithm 7. As the result of the for loop is needed in the subsequent operations in lines 4-12 of Algorithm 7 a second kernel is used.

The third CUDA kernel function is implemented to perform the small Montgomery reduction operation (i.e., `sm_mrq` Figure 1), which is employed to eliminate errors due to the base conversion operation in the previous step. After the NTT operations are applied to all vectors both in the RNS and extension bases, the fourth and fifth CUDA kernel functions are used to perform multiplication of the ciphertexts; the former in the RNS base \mathcal{Q} , `mult_q` and the latter in the extension base \mathcal{B}_{sk} , `mult_BSK` (see the middle block in Figure 1). Then, the INTT operation follows the multiplication operation to convert the ciphertexts back to the polynomial domain. The sixth CUDA kernel function is used to implement the multiplication of ciphertexts with the plaintext modulus t , `multip_t`. The seventh CUDA kernel function, named `first_fast_floor`, implements the first step of the `fast_floor` function (see Algorithm 8 for the pseudo-code): the results of the `multip_t` kernel function in \mathcal{Q} and \mathcal{B}_{sk} bases are converted to the \mathcal{B}_{sk} base. The eighth CUDA kernel function, named `second_fast_floor`, eliminates errors with the

Algorithm 7 Fast Convert Array

```

Input:  $r_i^n \in \mathbf{R}_{q_i,n}$ 
Input:  $P_i = [\frac{q_i}{q}]_{q_i}$ ,  $base\_c_i^j = [\frac{q_i}{q_i}]_{Bsk_j}$ 
Output:  $ct_i^n \in \mathbf{R}_{Bsk,n}$ 
1: for  $i$  from 0 by 1 to  $r-2$  do
2:    $mult_i^n = [r_i^n \times P_i]_{q_i}$ 
3: end for
4: for  $i$  from 0 by 1 to  $Bsk\_len - 1$  do
5:    $sum_i^n = 0$ 
6:   for  $j$  from 0 by 1 to  $r - 2$  do
7:     for  $k$  from 0 by 1 to  $n - 1$  do
8:        $sum_i^k = [sum_i^k + (mul_i^k \times base\_c_i^j)]_{Bsk_j}$ 
9:     end for
10:  end for
11:   $ct_i^n = sum_i^n$ 
12: end for
    
```

Algorithm 8 Fast Floor

```

Input:  $r_i^n \in \mathbf{R}_{q_i,n}$ ,  $r_j^n \in \mathbf{R}_{Bsk_j,n}$ ,  $P_i = [\frac{q_i}{q}]_{q_i}$ 
Input:  $base\_c_i^j = [\frac{q_i}{q_i}]_{Bsk_j}$ 
Output:  $c_i^n \in \mathbf{R}_{q_i,n}$ 
1:  $ct_i^n \leftarrow fast\_conv\_array(r_i^n, P_i, base\_c_i^j)$ 
2: for  $i$  from 0 by 1 to  $Bsk\_len - 1$  do
3:   for  $k$  from 0 by 1 to  $n - 1$  do
4:      $ct_i^k = [r_j^k - ct_i^k]_{Bsk_i}$ 
5:      $c_i^k = [ct_i^k \times [q]_{Bsk_i}^{-1}]_{Bsk_i}$ 
6:   end for
7: end for
    
```

flooring method instead of the rounding method. After the `fast_floor` kernel function, the fast base conversion function is performed in the ninth and tenth CUDA kernel functions. The ninth kernel functions performs the conversion from the extension base \mathcal{B}_{sk} to the RNS base \mathcal{Q} . Finally, the tenth kernel function, `second_fastbdconv_sk` is used to eliminate the rounding errors.

3) RELINEARIZATION

The BFV relinearization operation uses `switchkey` operation as explained in Figure 2, a pseudo-code of which is given in Algorithm 9 as it is implemented in the Microsoft SEAL library. The relinearization operation usually follows a homomorphic multiplication of ciphertexts, which are given in polynomial domain in BFV. The third component of the ciphertext, $c[2]$, which are to be multiplied with evaluations keys are first converted to the NTT domain using our NTT implementation (see line 5 of Algorithm 9). Then, the multiplication with evaluation keys are performed in the lines 2-9 of Algorithm 9, which are implemented in a single kernel function.

Due to the fact that no polynomial multiplication is needed after line 9, the results are converted back to the polynomial domain (see lines 10-13 of Algorithm 9). The lines 14 and

Algorithm 9 Switch Key

Input: $c_i[0], c_i[1], c_i[2] \in \mathbf{R}_{q_i, n}$
Input: $evk_i^k[k] \in \mathbf{R}_{q_i, n}$, where $k \in \{0, 1\}$, $0 \leq j < r$, and $0 \leq i < (r-1)$
Output: $ct_i[0], ct_i[1] \in \mathbf{R}_{q_i, n}$

- 1: $\bar{A}_{j,k} = 1$
- 2: **for** i from 0 by 1 to $r-2$ **do**
- 3: **for** j from 0 by 1 to $r-1$ **do**
- 4: **for** k from 0 by 1 to 1 **do**
- 5: $a_{i,j,k} = [NTT_{n,q_j}(c_i[2]) \odot evk_i^j[k]]_{q_j}$
- 6: $\bar{A}_{j,k} = [\bar{A}_{j,k} + a_{i,j,k}]_{q_j}$
- 7: **end for**
- 8: **end for**
- 9: **end for**
- 10: **for** j from 0 by 1 to r **do**
- 11: $A_{j,0} = INTT_{n,q_j}(\bar{A}_{j,0})$
- 12: $A_{j,1} = INTT_{n,q_j}(\bar{A}_{j,1})$
- 13: **end for**
- 14: $half = \lfloor \frac{q_{r-1}}{2} \rfloor$
- 15: **for** i from 0 by 1 to $r-1$ **do**
- 16: $halfmod = [half]_{q_i}$
- 17: **for** k from 0 by 1 to 1 **do**
- 18: $tmp = [[A_{r-1,k} + half]_{q_{r-1}} - halfmod]_{q_i}$
- 19: $tmp = [tmp \times q_r^{-1}]_{q_i}$
- 20: $ct_i[k] = [c_i[k] + tmp]_{q_i}$
- 21: **end for**
- 22: **end for**

16 are used to implement the arithmetic with the half modulus as previously described in Figure 2. Lastly, the operation between lines 14 to 22 in Algorithm 9 is implemented with a single kernel.

4) ROTATION

The BFV rotation operation uses the so-called `apply_galois` method, whose pseudo-code is given in Algorithm 10 and the `switchkey` operation in Algorithm 9. Before the rotation operation, `galois_elt` algorithm for a given shift amount is executed in CPU using Algorithm 11 and the result `galois_elt` is sent to GPU. Then, a single kernel is used to implement `apply_galois` algorithm. Finally, another kernel function is used to implement `switchkey` operation as explained in part IV-B3.

V. EXPERIMENTAL RESULTS

In this section, we present our GPU implementation results and their comparison with state-of-the-art works in the literature. Also, we present the implementation of gradient boosting framework (XGBoost) [30] using our GPU library to show its performance in practical real-world applications.

For a fair comparison with GPU and CPU implementations of NTT and of the homomorphic operations of the BFV scheme, we use a powerful CPU and two GPU devices, whose configurations are listed in Table 2. For the CPU

Algorithm 10 Apply Galois

Input: $galois_elt, c_i^j[k] \in \mathbf{R}_{q_i, n}$, where $0 \leq i < (r-1)$
 $0 \leq j < n, k = 0, 1$
Output: $c_i^j[k] \in \mathbf{R}_{q_i}$

- 1: **for** i from 0 by 1 to $r-2$ **do**
- 2: **for** j from 0 by 1 to $n-1$ **do**
- 3: $index_raw = j \times galois_elt$
- 4: $index = index_raw \& (n-1)$
- 5: **for** k from 0 by 1 to 1 **do**
- 6: $r_val = c_i^j[k]$
- 7: **if** $(index_raw \gg \log_2(n)) \& 1$ **then**
- 8: $non_zero = int(r_val \neq 0)$
- 9: $r_val = (q_i - r_val) \& (-non_zero)$
- 10: **end if**
- 11: $c_i^j[k] = r_val$
- 12: **end for**
- 13: **end for**
- 14: **end for**

Algorithm 11 Galois Elt

Input: $steps, n$
Output: $galois_elt$

- 1: $m32 = n \times 2$
- 2: **if** $steps == 0$ **then**
- 3: **return** $m32 - 1$
- 4: **else**
- 5: $pop_steps = abs(steps)$
- 6: **if** $steps < 0$ **then**
- 7: $steps = (n \gg 1) - pop_steps$
- 8: **else**
- 9: $steps = pop_steps$
- 10: **end if**
- 11: $gen = 3$
- 12: $galois_elt = 1$
- 13: **for** i from 0 by 1 to $steps$ **do**
- 14: $galois_elt = galois_elt \times gen$
- 15: $galois_elt = galois_elt \& (m32 - 1)$
- 16: **end for**
- 17: **return** $galois_elt$
- 18: **end if**

implementation, we use the Microsoft SEAL library³ for the BFV scheme, which is one of the fastest and highly optimized software implementations of BFV developed in C++ language. We use a combination of the build options given by the SEAL developers,⁴ which yields the fastest executable (More specifically, the build options set the `CMAKE_CXX_FLAGS_RELEASE` variable to include the “-DNDEBUG -flto -O3” flags by default.). Finally, our 32-bit implementation uses no assembly optimization while the 64-bit implementation utilizes assembly language to

³<https://github.com/microsoft/SEAL>

⁴<https://github.com/microsoft/SEAL#basic-cmake-options>

TABLE 2. Hardware features of the Testbed environment.

Feature	CPU	GPU	
		RTX3060Ti	GTX1080
Model	Ryzen7 3800X	RTX3060Ti	GTX1080
Threads	16	4864	2560
Freq.	4.20 GHz	1665 MHz	1733 MHz
RAM	32 GB (3600 MHz)	8 GB	8 GB
Mem. Type	-	GDDR6	GDDR5X
Mem. Bus	-	256 bits	256 bits
Bandwidth	-	448 GB/s	320 GB/s

Mem.: Memory.

Operating system and version: Windows 10 21H2.

CUDA version: 11.6.2

optimize only some of the 64-bit arithmetic operations such as 64-bit modular multiplication. This is solely needed for fast carry detection as GPU cores are essentially of 32-bit architecture. On the other hand, the same optimization is not needed for CPU as it is already 64-bit architecture that performs 64-bit arithmetic operations in hardware.

A. GPU IMPLEMENTATION OF NTT RESULTS AND COMPARISON WITH RELATED WORKS

Since the BFV scheme used here employs RNS, NTT must be concurrently calculated for each modulus in RNS. Therefore, it is essential to simultaneously perform multiple NTT operations in batches. Naturally, the throughput of an NTT operation is as important as (if not more than) the latency of a single NTT operation on GPU. In our GPU implementation, we aim to optimize both throughput and latency and we favor the former over the latter most of the time. In the literature, there are few works that report results for batch execution of NTT operations. Thus, in Figure 10 we include results from [29], which represents the state-of-the-art in GPU implementation of NTT and is the only work in the literature that reports batch computation results comparable to ours to the best of our knowledge. The source codes of the GPU implementation used in [29] are publicly available on GitHub⁵ and to provide a fair comparison for batch execution, we run them on our GPU device (i.e., RTX 3060Ti) in order to compare them with our results obtained from the same GPU device in Figure 10. Note that a more detailed version of Figure 10 is provided in Table 7 in Appendix A.

While most NTT GPU implementations in the literature use special form moduli to accelerate NTT operation, our implementation works with any NTT-friendly modulus and it is still faster. Furthermore, our implementation, which is optimized for performing NTT operations in batches, outperforms those that report only the timings for a single NTT operation in the literature. To compare our work with those that report only single NTT and inverse NTT timings, we include Table 3, which shows that, our implementation also outperforms all works in the literature except for one case when a single NTT (iNTT) operation is executed.

TABLE 3. Timings of GPU implementation of single NTT and single inverse NTT operations and their comparison with the works in literature.

Work	Device	n	$\log_2 q$	NTT	iNTT
[31]	Titan V	2^{14}	60	44.1 μs	- μs
		2^{15}	60	84.2 μs	- μs
[32]	RTX 2080 Ti	2^{15}	64*	83.3 μs	96 μs
[33]	GTX 1070	2^{14}	64*	57.8 μs	- μs
[34]	GTX 1070	2^{14}	64*	66.8 μs	- μs
		2^{15}	55	51 μs	41 μs
[29]	GTX 1080	2^{14}	55	33 μs	20 μs
		2^{15}	55	36 μs	24 μs
	Tesla V100	2^{14}	55	29 μs	21 μs
		2^{15}	55	39 μs	23 μs
This Work	RTX 3060 Ti	2^{12}	32	10.2 μs	10.2 μs
		2^{13}	32	10.9 μs	11.1 μs
		2^{14}	32	13.8 μs	14.2 μs
		2^{15}	32	19.4 μs	20.0 μs
		2^{12}	64	14.0 μs	15.0 μs
		2^{13}	64	14.9 μs	17.2 μs
		2^{14}	64	19.1 μs	23.1 μs
		2^{15}	64	35.9 μs	37.1 μs

*: Actual q_i is restricted by $q_i^2 n < 2^{64} - 2^{32} + 1$

In [29], the inverse NTT operation is faster than ours for ring sizes 2^{14} and 2^{15} . For the ring size 2^{14} , the total time of NTT and inverse NTT operations of our implementation is less than that in [29] (compare 42.4 μs and 50 μs). For 2^{15} , the implementation in [29], on the other hand, outperforms ours. Nevertheless, as Table 7 shows that our batch implementation outperforms the one in [29] for every case. The performance of batch NTT is much more important as NTT (and inverse NTT) operations are always executed in batches in all homomorphic encryption applications.

We also note that the works [32], [33], and [34] use the special modulus, $Q = 2^{64} - 2^{32} + 1$ known as goldilock prime, to perform NTT operations faster. However, Q serves as the carrier modulus for the actual moduli used in RNS arithmetic in homomorphic encryption applications. Thus, the actual moduli are much smaller due to the constraint $q_i^2 n < Q$ [35]. For example, for the ring size $n = 2^{14}$, each moduli in RNS arithmetic can be at most 25-bit. As our work can employ 64-bit RNS moduli, our actual performance is much better than the implementations in [32], [33], and [34]. For example, to match our size the implementations in those works should use at least twice as many RNS moduli.

The implementations in [32] and [33] take 83.3 μs and 57.8 μs , respectively, for 32768 ring size. And the implementation in [34] takes 66.8 μs for 16384 ring size. Our GPU implementation takes either 19.4 μs (32-bit implementation) or 35.9 μs (64-bit implementation) for the the ring dimension of 32768. On the other hand, when the ring dimension is 16384, our GPU implementation takes either 13.8 μs (32-bit implementation) or 19.4 μs (64-bit implementation) to perform single NTT operation. As the

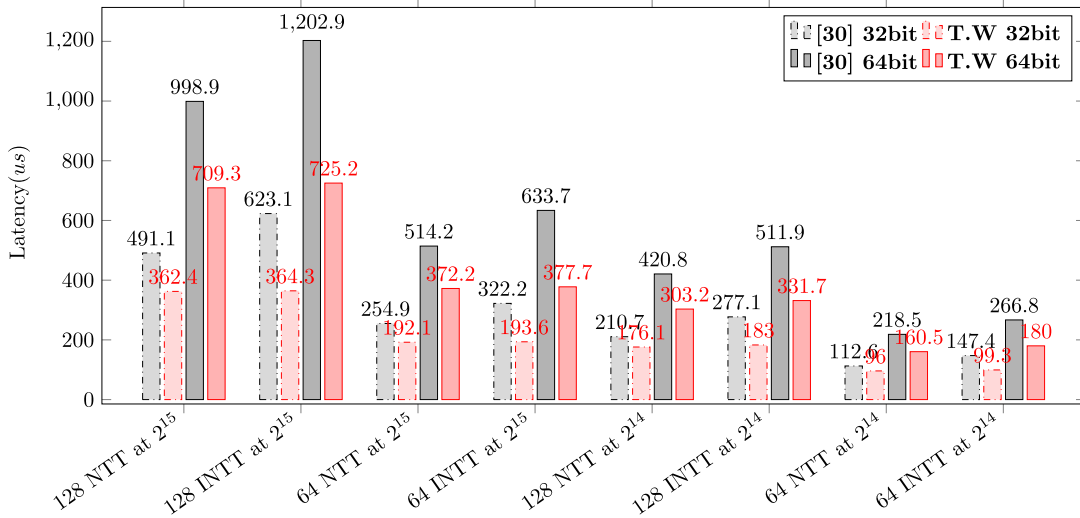


FIGURE 10. Timing bar graph of GPU implementation of NTT and inverse NTT operations on RTX 3060Ti and their comparison with [29].

TABLE 4. Comparison results of SEAL BFV scheme operations and literature with our GPU implementations of BFV scheme operations.

Operation	n	log ₂ q	GPU with [29] NTT		GPU with new NTT		[22]	SEAL	T
			RTX3060Ti	GTX1080	RTX3060Ti	GTX1080	Tesla V100	CPU	T _s
Add.	2 ¹²	109	4 µs	4.6 µs	4 µs	4.6 µs	-	14 µs	3.5×
	2 ¹³	218	5.1 µs	6.2 µs	5.1 µs	6.1 µs	-	58 µs	11.37×
	2 ¹⁴	438	12.3 µs	19.4 µs	12.3 µs	19.4 µs	-	233 µs	18.94×
	2 ¹⁵	881	44 µs	64.2 µs	44 µs	64.2 µs	-	778 µs	17.68×
Mult.	2 ¹²	109	172 µs	259 µs	86 µs	155.8 µs	-	3212 µs	37.3×
	2 ¹³	218	297 µs	532 µs	202 µs	423.4 µs	-	11883 µs	58.8×
	2 ¹⁴	438	1037 µs	2294 µs	768 µs	1856.1 µs	-	48757 µs	63.4×
	2 ¹⁵	881	5372 µs	10657 µs	3757 µs	-	-	205295 µs	54.6×
Relin.	2 ¹²	109	46 µs	82.7 µs	39.51 µs	59.3 µs	-	625 µs	15.81×
	2 ¹³	218	104 µs	145 µs	88.54 µs	143.4 µs	-	3100 µs	35.01×
	2 ¹⁴	438	462 µs	1013 µs	376.61 µs	825.3 µs	-	18295 µs	48.57×
	2 ¹⁵	881	3530 µs	6651 µs	3150 µs	-	-	111736 µs	35.47×
Rot.	2 ¹²	109	51 µs	87 µs	42.1 µs	59.4 µs	-	642 µs	15.24×
	2 ¹³	218	116 µs	172 µs	103.3 µs	162.7 µs	-	3157 µs	30.56×
	2 ¹⁴	438	544 µs	1339 µs	458.7 µs	1067.2 µs	-	18338 µs	39.97×
	2 ¹⁵	881	3879 µs	10504 µs	3464.5 µs	-	-	113437 µs	32.74×
Mult. + Relin.	2 ¹²	60	-	-	136 µs	-	859 µs	-	6.31×
	2 ¹³	120	-	-	170 µs	-	1012 µs	-	5.95×
	2 ¹⁴	360	-	-	661 µs	-	2010 µs	-	3.04×
	2 ¹⁵	600	-	-	2875 µs	-	4826 µs	-	1.67×

Add.:BFV Addition. Mult.:BFV Multiplication. Relin.:BFV Relinearization. T_s: speed up

works in [31], [32], [33], and [34] do not report timing results for batch execution, these works are not included in Table 7, which includes only comparison of batch NTT execution.

Table 7 lists the GPU timings for NTT and inverse NTT operations, which are organized into two main columns. On the left are the GPU timings when the modular multiplication (see Barret reduction in Algorithm 1) is done using 32-bit modulus q, whereas, on the right, the same timings are listed for a 64-bit q. Note that as the BFV scheme works with integers, only integer arithmetic is employed for these modular multiplications.

Note that as the 64-bit implementation uses twice the size of the modulus than the 32-bit implementations, it can be advantageous for homomorphic operations. For instance, for an acceptable level of security, one must use a 218-bit size for q when n = 2¹³. When we set the sizes of moduli q_i to 32-bit, we use seven RNS moduli q_i. On the other hand, if the size of each q_i is 64 bits, we use only five such moduli.

In the table, NTT_count represents the number of independent NTT operations performed simultaneously. We compare forward NTT and inverse NTT separately.

Our timing results show a significant acceleration in comparison with those of the state-of-the-art GPU

implementation in the literature [29]. For the 32 bit case, the new forward NTT and inverse NTT are $1.36\times$ and $1.71\times$, faster than their counterparts in the work [29], respectively, where $n = 2^{15}$ and $\text{NTT_count} = 128$. The sum of NTT and INTT timings, which is $726.7\mu\text{s}$, is $1.53\times$ faster than that of [29]. For the 64 bit case with $n = 2^{15}$ and $\text{NTT_count} = 128$, the new forward NTT and inverse NTT are $1.4\times$ and $1.66\times$ faster than their counterparts in [29], respectively. These performance achievements obtained for NTT and INTT operations help accelerate the operations of the BFV scheme and any application using homomorphic encryption as shown in the following section.

B. GPU IMPLEMENTATIONS OF BFV HE OPERATIONS RESULTS AND COMPARISON WITH RELATED WORKS

There are not many prior works in the literature that present GPU implementations of homomorphic operations of the BFV scheme, and the existing ones do not give performance results for all homomorphic operations let alone the homomorphic application results. Therefore, the comparison of our work with other works in the literature cannot be comprehensive. The work in [34] provides timing results on GPU for homomorphic applications of an old and completely different homomorphic encryption scheme, LTV [36], which is not in use today. We compare the results of our GPU implementation of the BFV-scheme operations with the work [22], which represents the state-of-the-art in the literature for GPU implementation of the BFV scheme. The work [22] provides only the timing results of homomorphic multiplication, which include those of the following relinearization operation.

We first provide our timing results separately in Table 4, which includes all major homomorphic arithmetic operations, typically used in many homomorphic applications. Our GPU implementation shows significant improvements over the CPU implementation of the SEAL library running on a powerful CPU. The Microsoft SEAL library, which is open-source and written in C++ programming language, provides highly optimized NTT implementation and therefore, represents state-of-the-art implementations of both the NTT and the BFV algorithms. As shown in Table 4 the proposed GPU library provides up to $63.4\times$ faster BFV multiplication operation, $48.57\times$ faster BFV relinearization operation, $39.97\times$ faster BFV rotation operation, $18.94\times$ faster BFV addition operation, when $n = 2^{14}$ and $q = 438$ with respect to the SEAL library, which is running on AMD Ryzen7 3800X. The total number of threads available in our GPU devices accounts for the optimum results obtained at $n = 2^{14}$. Since we parallelized all operations, RTX 3060Ti's threads become fully utilized and, the system cannot be parallelized more. Since the number threads on GTX 1080 is much fewer than RTX 3060 Ti, the best scenario for GTX 1080 is obtained at $n = 2^{13}$.

TABLE 5. Power consumption of BFV multiplication on CPU and GPU.

n	Count	CPU		GPU	
		Power*	Time	Power*	Time
2^{12}	1	55.5 W	3025 μs	44.37 W	92 μs
	10	60.49 W	30430 μs	44.5 W	830 μs
	100	60.65 W	317682 μs	44.54 W	9890 μs
	500	63.55 W	1369910 μs	44.49 W	51246 μs
2^{13}	1	55.71 W	9121 μs	44.15 W	161 μs
	10	59.52 W	89266 μs	44.25 W	1589 μs
	100	61.39 W	852098 μs	44.24 W	16275 μs
	500	64.78 W	3970960 μs	47.05 W	81714 μs
2^{14}	1	57.23 W	38414 μs	44.35 W	829 μs
	10	61.31 W	381921 μs	44.37 W	7515 μs
	100	60.8 W	3763901 μs	44.61 W	71858 μs
	500	64.16 W	18786647 μs	48.60 W	337236 μs
2^{15}	1	59.55 W	186334 μs	44.29 W	3382 μs
	10	59.46 W	1796976 μs	45.48 W	36301 μs
	100	60.89 W	17822724 μs	55.71 W	338212 μs
	500	66.93 W	89749372 μs	65.13 W	342199 μs

Power* = Max power observation during operation.

Idle CPU Power = 24 W

Idle GPU Power = 17 W

Care must be taken when evaluating CPU and GPU timing results as a fair comparison of GPU and CPU implementations can be very difficult due to differences in their architectures and applicable optimization techniques. As pointed out in [37] acceleration figures can overestimate the performance of a GPU device when compared to a CPU. The Microsoft SEAL library is probably the fastest CPU implementation of the BFV scheme in the open source and we are unable to identify any further optimization or parallelization technique that outperforms its NTT or homomorphic operation implementation. On the other hand, we find that CPU parallelization when deployed in the right way can be extremely useful in acceleration at application level as shown in Section V-D. Indeed, the speedup of GPU implementation turns out to be more moderate when all computing power of CPU is utilized for XGBoost classification application in Section V-D; a result which is in line with those in [37].

Then, we compare our results with the work in [22] only for homomorphic multiplication including the following relinearization operation as it is the only one reported. The GPU used in [22] has 5120 cores and 16 GB of memory operating at the clock frequency of 1.380 GHz, which is comparable to RTX 3060Ti used in our measurements. The execution times are also measured for the same ciphertext modulus sizes and the ring dimensions used in the work [22] for a fair comparison. As observable from Table 4, our GPU implementation outperforms that in [22] for all cases. For instance, our multiplication including relinearization implementation results are $6.31\times$ faster for $n = 2^{12}$, $5.95\times$ faster for ring size $n = 2^{13}$, $3.04\times$ faster for ring size $n = 2^{14}$, and $1.67\times$ faster for ring size $n = 2^{15}$ than the work [22], respectively.

C. POWER CONSUMPTION OF BFV HOMOMORPHIC MULTIPLICATION ON CPU AND GPU

This section reports on power consumption of homomorphic multiplication on GPU and CPU devices in Table 2 to evaluate the power efficiency of GPU implementations. The software tool HWMonitor, which is free and publicly available,⁶ is used for the experimental setup in this section. HWMonitor instantly shows the amount of power consumed by each unit of a computer system individually. It also shows the maximum and minimum power consumption achieved during its operation. In this experiment, only the power consumption of the CPU and GPU are considered (i.e., RAM or PCI-e power consumption are excluded).

For comparison we used homomorphic multiplication of the BFV scheme as the benchmark, which is the primary, and also the most time and resource consuming, operation in all homomorphic evaluations. We simply compared power requirements of the Microsoft SEAL implementation of BFV homomorphic multiplication running on CPU and that of our GPU implementation of the same operation running on GPU. In order to measure the power consumption on both devices in a fair way, different numbers of operations are performed using different ring sizes (whereby operations run sequentially). In order to create a reliable experimental environment, the power required by the CPU and GPU in the IDLE state is found using the same software tool before the measurements. As stated in Table 5, the IDLE power consumption of the CPU and the GPU are 24 W and 17 W, respectively. All observed power usage and latency measurements are enumerated in Table 5, which capture the maximum power values observed on the HWMonitor display, when the test codes are run. The tabulated power measurements include IDLE and execution power consumption. For instance, a maximum of 64.16 W (24 W + 40.16 W) is observed during the BFV multiplication, which is run 500 times on the CPU for the ring size of 2^{14} , while 48.6 W (17 W + 31.6 W) is observed on the GPU for the same case.

As a result, we can conclude that the GPU implementation seems to be a more power efficient alternative for single BFV multiplication than an optimized CPU implementation. On the other hand, care must be taken when interpreting the figures reported in this work as the power measurement and the accuracy of the measurement tools are always a research subject as computer devices with all their subsystems are overly involved.

D. IMPLEMENTATION RESULTS OF PRIVACY-PRESERVING INFERENCE FOR GENOME DATA USING XGBoost TREES

Mağara et al. [30] introduce a privacy-preserving gradient boosting inference framework (XGBoost) algorithm using homomorphic encryption for the classification of the encrypted genome data of different tumor types. We implemented their framework using our GPU library of the BFV scheme. XGBoost is a learning algorithm, which uses

TABLE 6. Implementation of gradient boosting framework(XGBoost) results.

n	$\log_2 q$	SEAL		T.W.		
		S.T.	M.T.	RTX 3060Ti	T	S
2^{13}	218	25.62 s	3.4 s	0.596 s	$42.98\times$	$5.7\times$
2^{14}	438	127.028 s	19.27 s	2.41 s	$52.7\times$	$8\times$

S.T.:Single Thread M.T.:Multi Thread.(16 threads) T.W.:This Work. T: The ratio of single-thread results over this work. S: The ratio of multi-thread results over this work.

Extreme Gradient Boosting ensembles. The model consists of classification trees that are constructed by training data. Trees of the ensemble evaluate the test data that are classified into one of the leaves. Lastly, a final prediction score is formed by summing up the numerical scores obtained from each tree. To decrease the complexity of the model and the depth of the corresponding circuit to be homomorphically evaluated, shallow trees are selected.

As explained in [30], test data is encrypted, and the XGBoost trees are homomorphically evaluated for a total of 258 test data points. The total number of homomorphic multiplications, rotations, subtractions, plain multiplications, addition, and relinearization operations are 1290, 1806, 1806, 1290, 3354, and 2322, respectively.

We run the inference framework both on GPU and CPU, and all known possible optimization and parallelization techniques to us are employed for the CPU implementation. As shown in Table 6, our GPU library accelerates the classification operation at least 42.98 times with respect to the results obtained from AMD Ryzen7 3800X CPU with a single thread while the speedup is 5.7 when multi-threaded version of the CPU implementation is used.

VI. CONCLUSION

In this paper, we presented a GPU library that features highly parallelized and optimized implementations of NTT and inverse NTT operations and homomorphic operations of the BFV scheme. Although the library can be independently used, it is also integrated with the Microsoft SEAL library and its functions can be called from any application code using SEAL. Therefore, the library is truly an accelerator for homomorphic encryption applications.

By reducing the number of GPU kernel function calls and optimizing the use of fast memory on GPU, the library offers the best timing performance for NTT and inverse NTT operations in the literature. For instance, concurrent executions of 128 NTT and INTT operations for the ring degree of 2^{14} take 303.19 μs and 331.7 μs , respectively, on RTX3060Ti GPU, which are 1.39 and 1.54 times faster than those of the state-of-the-art GPU implementation reported in the literature.

Then, all homomorphic operations of the BFV scheme are also implemented on GPU and compared against the SEAL library running on a CPU. When compared with CPU implementation for the ring size of 2^{14} and the modulus bit

⁶<https://hwmonitor.softonic.com>

TABLE 7. Timings of GPU implementation of NTT and inverse NTT operations and their comparison with [29].

		32 Bit (Implemented On RTX 3060Ti)						64 Bit (Implemented On RTX 3060Ti)					
		Forward NTT			Inverse NTT			Forward NTT			Inverse NTT		
n	NTT_count	[29]	T.W.	T	[29]	T.W.	T	[29]	T.W.	T	[29]	T.W.	T
2 ¹²	4	12.3 μ s	11 μs	1.12 \times	11.2 μ s	11.1 μs	1.11 \times	19.5 μ s	14.3 μs	1.36 \times	16 μ s	15.4 μs	1.03 \times
	16	13 μ s	11.2 μs	1.16 \times	12.2 μ s	12.2 μ s	1 \times	22.5 μ s	17.2 μs	1.30 \times	17.8 μ s	19.4 μ s	0.91 \times
	32	13.9 μ s	17.9 μ s	0.77 \times	18.4 μ s	19.2 μ s	0.96 \times	23.5 μ s	25.2 μ s	0.93 \times	29.3 μ s	26.6 μs	1.10 \times
	64	23.1 μ s	28.6 μ s	0.80 \times	29.5 μ s	29.3 μs	1 \times	39.9 μ s	43 μ s	0.93 \times	52.9 μ s	50.1 μs	1.05 \times
	128	38.9 μ s	46.7 μ s	0.83 \times	48.1 μ s	48.7 μ s	0.98 \times	75.7 μ s	81.6 μ s	0.92 \times	96.5 μ s	91.1 μs	1.06 \times
2 ¹³	4	17.1 μ s	12.2 μs	1.40 \times	14.1 μ s	14.3 μ s	0.98 \times	24.4 μ s	16.6 μs	1.47 \times	21.1 μ s	20.5 μs	1.03 \times
	16	18.4 μ s	18.4 μ s	1 \times	22.2 μ s	20.3 μs	1.09 \times	28.2 μ s	25.6 μs	1.10 \times	34.8 μ s	28.6 μs	1.21 \times
	32	28.6 μ s	29.4 μ s	0.97 \times	34.9 μ s	30.3 μs	1.15 \times	47.9 μ s	44.4 μs	1.08 \times	59.4 μ s	49.9 μs	1.19 \times
	64	49.8 μ s	46.7 μs	1.07 \times	57.3 μ s	50.7 μs	1.13 \times	97.1 μ s	82.1 μs	1.18 \times	121.2 μ s	93.9 μs	1.29 \times
	128	88 μ s	91.2 μ s	0.96 \times	124 μ s	96.1 μs	1.29 \times	170.7 μ s	156.4 μs	1.09 \times	224.1 μ s	173.6 μs	1.29 \times
2 ¹⁴	4	20.1 μ s	15.2 μs	1.32 \times	17.6 μ s	16.3 μs	1.08 \times	29.98 μ s	21.76 μs	1.37 \times	24.5 μ s	24.1 μs	1.01 \times
	16	33.7 μ s	30.5 μs	1.05 \times	40.9 μ s	30.1 μs	1.36 \times	54.4 μ s	46.54 μs	1.17 \times	65.9 μ s	53.2 μs	1.23 \times
	32	57.3 μ s	50.1 μs	1.14 \times	64.5 μ s	51.8 μs	1.24 \times	121.8 μ s	84.6 μs	1.44 \times	143.3 μ s	97.2 μs	1.47 \times
	64	112.6 μ s	96 μs	1.17 \times	147.4 μ s	99.3 μs	1.48 \times	218.5 μ s	160.5 μs	1.36 \times	266.8 μ s	180 μs	1.48 \times
	128	210.7 μ s	176.1 μs	1.19 \times	277.1 μ s	183 μs	1.51 \times	420.8 μ s	303.19 μs	1.38 \times	511.9 μ s	331.7 μs	1.54 \times
2 ¹⁵	4	25.6 μ s	26.4 μ s	0.96 \times	28.7 μ s	25.9 μs	1.10 \times	35.8 μ s	41.2 μ s	0.86 \times	47.3 μ s	50.3 μ s	0.94 \times
	16	64.1 μ s	52.2 μs	1.22 \times	74.7 μ s	53.2 μs	1.40 \times	147.1 μ s	100 μs	1.47 \times	170.1 μ s	95.6 μs	1.78 \times
	32	136.3 μ s	100.3 μs	1.35 \times	173 μ s	102.4 μs	1.69 \times	266.2 μ s	191.8 μs	1.38 \times	325.5 μ s	193.2 μs	1.68 \times
	64	254.9 μ s	192.1 μs	1.32 \times	322.2 μ s	193.6 μs	1.66 \times	514.2 μ s	372.2 μs	1.38 \times	633.7 μ s	377.7 μs	1.67 \times
	128	491.1 μ s	362.4 μs	1.35 \times	623.1 μ s	364.3 μs	1.71 \times	998.9 μ s	709.3 μs	1.41 \times	1202.9 μ s	725.2 μs	1.66 \times

T.W:This Work. T: speed up

size of 438, the GPU library running on RTX3060Ti achieves speedups of 18.94, 63.4, 48.57, and 39.97 for homomorphic addition, homomorphic multiplication, relinearization, and homomorphic rotation, respectively. We also compared our homomorphic multiplication followed by a relinearization operation with that of the state-of-the-art GPU implementation in the literature, and found that ours is up to 6.31 times faster than the latter.

We also showed that the proposed GPU library is profitably used in the homomorphic processing of real data such as the classification of encrypted genome data for tumor types and reported at least a speedup of 5 in comparison with a powerful CPU running 16 threads.

The reported performance gains establish that GPU implementations of homomorphic encryption prove to be useful to help privacy-preserving data processing applications become more practicable.

As future work, we envision integrating our GPU accelerator into other HE libraries and using it to accelerate other more challenging operations such as bootstrapping and scheme switching. We can achieve these goals by joining recent open-source efforts in the development of HE software libraries such as OpenFHE [13].

APPENDIX

See Table 7.

ACKNOWLEDGMENT

(Ali Şah Özcan, Can Ayduman, and Enes Recep Türkoğlu are co-first authors.)

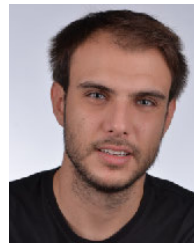
REFERENCES

- [1] C. Gentry, “Fully homomorphic encryption using ideal lattices,” in *Proc. 41st Annu. ACM Symp. Theory Comput.*, New York, NY, USA, May 2009, pp. 169–178, doi: 10.1145/1536414.1536440.
- [2] 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, doi: 10.1145/2535925.
- [3] C. Gentry and S. Halevi, “Implementing gentry’s fully-homomorphic encryption scheme,” in *Advances in Cryptology—EUROCRYPT (Lecture Notes in Computer Science)*, K. G. Paterson, Ed. Berlin, Germany: Springer, 2011, pp. 129–148.
- [4] Z. Brakerski and V. Vaikuntanathan, “Efficient fully homomorphic encryption from (standard) LWE,” *SIAM J. Comput.*, vol. 43, no. 2, pp. 831–871, 2014, doi: 10.1137/120868669.
- [5] 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, doi: 10.1145/2633600.
- [6] J. Fan and F. Vercauteren, “Somewhat practical fully homomorphic encryption,” *IACR Cryptol. ePrint Arch.*, vol. 2012, p. 144, Mar. 2012.
- [7] (Nov. 2020). *Microsoft SEAL (Release 3.6)*. Redmond, WA, USA. [Online]. Available: <https://github.com/Microsoft/SEAL>
- [8] (2021). *PALISADE Lattice Cryptography Library (Release 1.11.5)*. [Online]. Available: <https://palisade-crypto.org/>
- [9] S. Halevi and V. Shoup, “Algorithms in helib,” in *Advances in Cryptology—CRYPTO (Lecture Notes in Computer Science)*, J. A. Garay and R. Gennaro, Eds. Berlin, Germany: Springer, 2014, pp. 554–571.
- [10] W. Wang, Z. Chen, and X. Huang, “Accelerating leveled fully homomorphic encryption using GPU,” in *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS)*, Jun. 2014, pp. 2800–2803.
- [11] A. C. Mert, E. Öztürk, and E. Savas, “Design and implementation of encryption/decryption architectures for BFV homomorphic encryption scheme,” *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 28, no. 2, pp. 353–362, Feb. 2020.
- [12] Y. Doroz, E. Ozturk, E. Savas, and B. Sunar, “Accelerating LTV based homomorphic encryption in reconfigurable hardware,” in *Proc. Int. Workshop Cryptograph. Hardw. Embedded Syst.* Cham, Switzerland: Springer, 2015, pp. 185–204.

- [13] A. Al Badawi, J. Bates, F. Bergamaschi, D. B. Cousins, S. Erabelli, N. Genise, S. Halevi, H. Hunt, A. Kim, Y. Lee, and Z. Liu, "OpenFHE: Open-source fully homomorphic encryption library," in *Proc. 10th Workshop Encrypted Comput. Appl. Homomorphic Cryptogr.*, 2022, pp. 53–63. [Online]. Available: <https://eprint.iacr.org/2022/915>
- [14] *NVIDIA CUDA C Programming Guide*, NVIDIA Corporation, Santa Clara, CA, USA, 2010.
- [15] C. Boura, N. Gama, M. Georgieva, and D. Jetchev, "CHIMERA: Combining ring-LWE-based fully homomorphic encryption schemes," *J. Math. Cryptol.*, vol. 14, no. 1, pp. 316–338, 2020, doi: [10.1515/jmc-2019-0026](https://doi.org/10.1515/jmc-2019-0026).
- [16] J.-C. Bajard, J. Eynard, M. A. Hasan, and V. Zucca, "A full RNS variant of FV like somewhat homomorphic encryption schemes," in *Selected Areas in Cryptography—SAC (Lecture Notes in Computer Science)*, vol. 10532, R. Avanzi and H. M. Heys, Eds. St. John's, NL, Canada: Springer, 2017, pp. 423–442, doi: [10.1007/978-3-319-69453-5_23](https://doi.org/10.1007/978-3-319-69453-5_23).
- [17] P. Barrett, "Implementing the Rivest Shamir and Adleman public key encryption algorithm on a standard digital signal processor," in *Advances in Cryptology—CRYPTO*, (Lecture Notes in Computer Science), vol. 263. Santa Barbara, CA, USA: Springer, 1986, pp. 311–323.
- [18] P. L. Montgomery, "Modular multiplication without trial division," *Math. Comput.*, vol. 44, no. 170, pp. 519–521, Apr. 1985.
- [19] S. Antao, J.-C. Bajard, and L. Sousa, "RNS-based elliptic curve point multiplication for massive parallel architectures," *Comput. J.*, vol. 55, no. 5, pp. 629–647, May 2012, doi: [10.1093/comjnl/bxr119](https://doi.org/10.1093/comjnl/bxr119).
- [20] J. Bajard, J. Eynard, N. Merkiche, and T. Plantard, "RNS arithmetic approach in lattice-based cryptography: Accelerating the 'rounding-off' core procedure," in *Proc. 22nd IEEE Symp. Comput. Arithmetic*, Lyon, France, 2015, pp. 113–120, doi: [10.1109/ARITH.2015.30](https://doi.org/10.1109/ARITH.2015.30).
- [21] Z. Brakerski, "Fully homomorphic encryption without modulus switching from classical gapped," in *Advances in Cryptology—CRYPTO (Lecture Notes in Computer Science)*, vol. 7417, R. Safavi-Naini and R. Canetti, Eds. Santa Barbara, CA, USA: Springer, 2012, pp. 868–886, doi: [10.1007/978-3-642-32009-5_50](https://doi.org/10.1007/978-3-642-32009-5_50).
- [22] A. Al Badawi, Y. Polyakov, K. M. M. Aung, B. Veeravalli, and K. Rohloff, "Implementation and performance evaluation of RNS variants of the BFV homomorphic encryption scheme," *IEEE Trans. Emerg. Topics Comput.*, vol. 9, no. 2, pp. 941–956, Apr. 2021.
- [23] J. W. Cooley and J. W. Tukey, "An algorithm for the machine calculation of complex Fourier series," *Math. Comput.*, vol. 19, no. 90, pp. 297–301, 1965.
- [24] J.-C. Bajard and T. Plantard, "RNS bases and conversions," in *Proc. SPIE*, vol. 5559, pp. 60–69, Oct. 2004.
- [25] S. Halevi, Y. Polyakov, and V. Shoup, "An improved RNS variant of the BFV homomorphic encryption scheme," in *Proc. Cryptographers Track RSA Conf.* Cham, Switzerland: Springer, 2019, pp. 83–105.
- [26] H. Chen, W. Dai, M. Kim, and Y. Song, "Efficient multi-key homomorphic encryption with packed ciphertexts with application to oblivious neural network inference," in *Proc. ACM SIGSAC Conf. Comput. Commun. Secur.* New York, NY, USA: Association for Computing Machinery, Nov. 2019, pp. 395–412, doi: [10.1145/3319535.3363207](https://doi.org/10.1145/3319535.3363207).
- [27] K. Laine, "Simple encrypted arithmetic library 2.3.1," Microsoft Research, WA, USA, Tech. Rep., 2017.
- [28] Z. Jia, M. Maggioni, B. Staiger, and D. P. Scarpa, "Dissecting the NVIDIA Volta GPU architecture via microbenchmarking," 2018, *arXiv:1804.06826*.
- [29] O. Ozerk, C. Elgezen, A. Mert, E. Ozturk, and E. Savas, "Efficient number theoretic transform implementation on GPU for homomorphic encryption," *J. Supercomput.*, vol. 78, pp. 2840–2872, Jul. 2021, doi: [10.1007/978-3-319-45744-4_15](https://doi.org/10.1007/978-3-319-45744-4_15).
- [30] S. S. Magara, C. Yildirim, F. Yaman, B. Dilekoglu, F. Tutas, E. Öztürk, K. Kaya, O. Tastan, and E. Savas, "ML with he: Privacy preserving machine learning inferences for genome studies," in *Proc. ACM CCS Privacy Preserving Mach. Learn. Workshop*, 2021, pp. 1–5.
- [31] S. Kim, W. Jung, J. Park, and J. H. Ahn, "Accelerating number theoretic transformations for bootstrappable homomorphic encryption on GPUs," in *Proc. IEEE Int. Symp. Workload Characterization (IISWC)*, Oct. 2020, pp. 264–275.
- [32] Z. Zheng, "Encrypted cloud using GPUs," M.S. thesis, 2020.
- [33] J.-Z. Goey, W.-K. Lee, B.-M. Goi, and W.-S. Yap, "Accelerating number theoretic transform in GPU platform for fully homomorphic encryption," *J. Supercomput.*, vol. 77, no. 2, pp. 1455–1474, Feb. 2021, doi: [10.1007/s11227-020-03156-7](https://doi.org/10.1007/s11227-020-03156-7).
- [34] W. Dai and B. Sunar, "cuHE: A homomorphic encryption accelerator library," in *Proc. Int. Conf. Cryptography Inf. Secur. Balkans*. Cham, Switzerland: Springer, 2015, pp. 169–186.
- [35] W. Dai, Y. Doroz, and B. Sunar, "Accelerating NTRU based homomorphic encryption using GPUs," in *Proc. IEEE High Perform. Extreme Comput. Conf. (HPEC)*, Waltham, MA, USA, Sep. 2014, pp. 1–6, doi: [10.1109/hpec.2014.7041001](https://doi.org/10.1109/hpec.2014.7041001).
- [36] A. López-Alt, E. Tromer, and V. Vaikuntanathan, "Multikey fully homomorphic encryption and applications," *SIAM J. Comput.*, vol. 46, no. 6, pp. 1827–1892, Jan. 2017, doi: [10.1137/14100124x](https://doi.org/10.1137/14100124x).
- [37] V. W. Lee, C. Kim, J. Chhugani, M. Deisher, D. Kim, A. D. Nguyen, N. Satish, M. Smelyanskiy, S. Chennupati, P. Hammarlund, R. Singhal, and P. Dubey, "Debunking the 100X GPU vs. CPU myth: An evaluation of throughput computing on CPU and GPU," *SIGARCH Comput. Archit. News*, vol. 38, no. 3, pp. 451–460, Jun. 2010, doi: [10.1145/1816038.1816021](https://doi.org/10.1145/1816038.1816021).



ALİ ŞAH ÖZCAN received the B.S. degree in electronics engineering program from Sabancı University, Istanbul, Turkey, in 2020, where he is currently pursuing the M.S. degree in electronics engineering program. His research interests include RFIC design, homomorphic encryption, lattice-based cryptography, and cryptographic hardware/software design.



CAN AYDUMAN received the B.S. degree in electronics engineering program from Sabancı University, Istanbul, Turkey, in 2020, where he is currently pursuing the M.S. degree in electronics engineering program. His research interests include RFIC design, homomorphic encryption, lattice-based cryptography, and cryptographic hardware/software design.



ENES RECEP TÜRKOĞLU received the B.S. degree in electronics engineering program from Sabancı University, Istanbul, Turkey, in 2020, where he is currently pursuing the M.S. degree in electronics engineering program. His research interests include RFIC design, homomorphic encryption, lattice-based cryptography, and cryptographic hardware/software designs.



ERKEY SAVAS (Member, IEEE) received the B.S. and M.S. degrees in electrical engineering from the Electronics and Communications Engineering Department, Istanbul Technical University, in 1990 and 1994, respectively, and the Ph.D. degree from the Department of Electrical and Computer Engineering, Oregon State University, in 2000. He has been a Faculty Member with Sabancı University, since 2002. His research interests include applied cryptography, data and communication security, security and privacy in data mining applications, embedded systems security, and distributed systems.

• • •