

# Implementation and Performance Evaluation of RNS Variants of the BFV Homomorphic Encryption Scheme

Ahmad Al Badawi, Yuriy Polyakov, Khin Mi Mi Aung,  
Bharadwaj Veeravalli, and Kurt Rohloff

DECEMBER 5, 2018

## Abstract

Homomorphic encryption is an emerging form of encryption that provides the ability to compute on encrypted data without ever decrypting them. Potential applications include aggregating sensitive encrypted data on a cloud environment and computing on the data in the cloud without compromising data privacy. There have been several recent advances resulting in new homomorphic encryption schemes and optimized variants. We implement and evaluate the performance of two optimized variants, namely Bajard-Eynard-Hasan-Zucca (BEHZ) and Halevi-Polyakov-Shoup (HPS), of the most promising homomorphic encryption scheme in CPU and GPU. The most interesting (and also unexpected) result of our performance evaluation is that the HPS variant in practice scales significantly better (typically by 15%-30%) with increase in multiplicative depth of the computation circuit than BEHZ, implying that the HPS variant will always outperform BEHZ for most practical applications. For the multiplicative depth of 98, our fastest GPU implementation performs homomorphic multiplication in 51 ms for 128-bit security settings, which is faster by two orders of magnitude than prior results and already practical for cloud environments supporting GPU computations. Large multiplicative depths supported by our implementations are required for applications involving deep neural networks, logistic regression learning, and other important machine learning problems.

## Index Terms

Secure Computation, Lattice-based Cryptography, Homomorphic Encryption, Residue Number System, BFV SHE, Parallel Processing.

## I. INTRODUCTION

**H**OMOMORPHIC encryption (HE) has been a topic of active research since the first design of a Fully Homomorphic Encryption (FHE) scheme by Gentry in 2009 [1]. FHE allows performing arbitrary secure computations over encrypted data without ever decrypting them. One of the potential applications is to upload encrypted data to a public cloud computing environment and then outsource computations over these data to the cloud without sharing secret keys or compromising data privacy by decrypting the data in the cloud.

At a high level, FHE is based on a Somewhat Homomorphic Encryption (SHE) scheme that provides at least two homomorphic operations: addition and multiplication. The encryption and homomorphic operations require adding some “noise” to guarantee a certain level of security based on the underlying hardness assumption, e.g., Ring Learning With Errors. This noise grows at some controlled rate with each homomorphic operation. As long as the noise is under a certain level, the decryption returns the correct result for a given computation circuit.

The literature includes a number of FHE and SHE schemes that vary in construction, functionality and performance [2], [3], [4], [5], [6], [7]. One of the most promising schemes is the Fan-Vercauteren variant of Brakerski’s

A. Al Badawi and B. Veeravalli are with the Department of Electrical and Computer Engineering, National University of Singapore, 117576; E-mail: ahmad@u.nus.edu and elebv@nus.edu.sg

Y. Polyakov and K. Rohloff are with the Department of Computer Science and Cybersecurity Research Center, New Jersey Institute of Technology, USA 07102; E-mail: polyakov@njit.edu and rohloff@njit.edu

K. M. M. Aung is with A \* STAR, Data Storage Institute (DSI), Singapore 138632; E-mail: mi\_mi\_aung@dsi.a-star.edu.sg

scale-invariant scheme [3], [5], which we refer to as BFV in this paper. The scheme has an elegant/simple structure and provides promising practical performance [8].

One common aspect of many HE schemes is the need to manipulate large algebraic structures with multi-precision coefficients. In BFV, this multi-precision arithmetic is performed for polynomials of large degrees (several thousand) with relatively large integer coefficients (several hundred bits). Implementing the necessary multi-precision modular arithmetic is computationally expensive. A way to make these operations faster is to use the Residue Number System (RNS) to decompose large coefficients into vectors of smaller, native (machine-word-size) integers. RNS allows some arithmetic operations to be performed completely in parallel using native instructions and data types, thus potentially improving the efficiency.

Recently, two RNS variants of the BFV scheme have been proposed to reduce the computational complexity of decryption and homomorphic multiplication<sup>1</sup>: Bajard-Eynard-Hasan-Zucca (BEHZ) [9] and Halevi-Polyakov-Shoup (HPS) [10]. Both variants share a common design but differ in some aspects such as implementation complexity, effect on noise growth, and performance. The BEHZ variant employs integer arithmetic and RNS techniques to provide an asymptotically better performance over the textbook BFV scheme [5] at the expense of higher noise growth. The HPS variant employs a combination of integer and floating-point arithmetic in addition to RNS techniques. HPS is relatively simpler to implement and has essentially the same noise growth as the BFV scheme [10].

There are multiple publications on software implementations of the BFV scheme. For instance, SEAL is an open source C++ library implementing the BFV scheme and its RNS variant BEHZ [11]. Earlier versions of SEAL used RNS for some operations of the BFV scheme but relied on multi-precision arithmetic for decryption and homomorphic multiplication. The latter require the divide-and-round and base decomposition operations that are incompatible with RNS and require a computationally expensive conversion to the positional (multi-precision) representation. However, the recent versions of SEAL (> v2.3.0) provide a full RNS implementation of the BEHZ variant [9].

Other implementations of the BFV scheme can be found in PALISADE [12], an open source C++ lattice cryptography library that includes the implementations of multiple HE and proxy re-encryption schemes [13]; digital signature, identity-based encryption, and attribute-based encryption constructions [14], [15]; and conjunction obfuscation scheme [16]. Starting with v1.1, PALISADE provides an implementation of the HPS variant.

The BFV scheme has also become a subject for hardware acceleration studies. For instance, Al Badawi et al. [17] provide a GPU-accelerated implementation of BEHZ. Another recent effort dealt with accelerating the textbook BFV performance using FPGA [18].

In this work we evaluate and compare the practical performance of BEHZ and HPS. Although Halevi et al. provide a theoretical comparison [10], it is not clear from that analysis how the noise growth and performance compare in practical implementations. As seen later in this paper, each variant employs different elementary operations that cannot be compared easily without experiments in the same settings.

## A. Our Contributions

In this work we implement the BEHZ variant and present an optimized implementation of the HPS scheme in PALISADE. We also implement the HPS variant in GPU. The new algorithmic optimizations added to the HPS implementations in both CPU and GPU include lazy reduction and several precomputations.

We examine the RNS techniques of both variants and compare their computational complexity and theoretical noise growth. We also provide recommendations for achieving the best practical performance of both variants.

We evaluate the performance of both variants in CPU and GPU. Our analysis suggests that the HPS decryption and homomorphic multiplication runtimes are typically smaller (up to 30%) than those for BEHZ for most settings.

We discover from the analysis of experimental noise growth that the HPS variant scales significantly better with increase in multiplicative depth of the computation circuit than BEHZ: the depth supported by the HPS variant is typically 15%-30% larger for the same values of parameters. We provide an interpretation of the faster noise growth for the BEHZ variant.

The comparison of the homomorphic multiplication runtimes of the HPS variant for CPU and GPU demonstrates that our best GPU performance results are 3x-33x faster than our best multi-threaded results for a modern server

<sup>1</sup>Homomorphic multiplication is known to be the most performance-critical primitive in HE schemes

CPU environment. This implies that we improve the prior implementation results for the HPS variant [10] by more than one order of magnitude.

For the multiplicative depth of 98, we are able to reduce the decryption time to 0.5 ms and homomorphic multiplication to 51 ms for 128-bit security settings, which is already practical for cloud environments supporting GPU computations. Large multiplicative depths supported by our implementations are required for applications involving deep neural networks [19], logistic regression learning [20], and other important machine learning problems. Our best runtime results are at least two orders of magnitude faster than all previous results for the BFV scheme in the literature.

## B. Organization

The paper is organized as follows: Section II provides some preliminaries on the RNS tools and textbook BFV scheme. Sections III and IV review, analyze, and theoretically compare the decryption and homomorphic multiplication procedures of both RNS variants. Implementation details are discussed in Section V. In Section VI, we benchmark both variants and analyze their performance for different platforms. Finally, Section VII concludes the work and marks out the future work.

## II. PRELIMINARIES

### A. Cyclotomic Rings

The BFV scheme implemented in our work employs the polynomial ring  $R = \mathbb{Z}[X]/(X^n + 1)$ , where the ring dimension  $n$  is a power of 2. The ring can be viewed as a set of polynomials of degree less than  $n$ . The arithmetic in  $R$  is always done modulo  $(X^n + 1)$ . In some BFV primitives, the polynomials are sampled from predefined distributions. We use the symbol  $a \stackrel{\mathcal{U}}{\leftarrow} \mathcal{S}$  to refer to uniform sampling of  $a$  from the set  $\mathcal{S}$ , whereas the symbol  $a \stackrel{\mathcal{G}}{\leftarrow} \mathcal{S}$  is used to denote sampling from a Gaussian distribution.

The plaintext space in BFV is  $R_t$ , where  $t \geq 2$  is an integer plaintext modulus. The polynomials in plaintext space are reduced both modulo  $t$  and  $(X^n + 1)$ . A plaintext is normally a single element in  $R_t$  encoding the original plaintext message. Likewise, the ciphertext space  $R_q$  has  $q \gg t$  as the coefficient modulus. For practical implementations,  $q$  is usually a  $k$ -smooth number, s.t.  $q = \prod_{i=1}^k p_i$ , where  $p_i$  is a prime that fits in the underlying machine word. Similar to  $R_t$ , polynomials in  $R_q$  are reduced modulo  $q$  as well. We remark that unlike a plaintext, a ciphertext  $c$  is a pair of two elements in  $R_q$ , denoted by  $(c[0], c[1])$ .

### B. Residue Number System and Chinese Remainder Theorem

*a) Residue Number System:* RNS is a non-positional numbering system in which a number is represented by a tuple of residues modulo some predefined pairwise co-prime moduli, known as the RNS base. RNS is used to distribute a computation in some relatively large domain to a set of smaller sub-domains. Computations in sub-domains are completely independent and, therefore, can be performed in parallel. Moreover, the problem size in sub-domains can be highly controlled so that the computation can be done without using multi-precision arithmetic. The results of independent computations in sub-domains can be interpolated via the Chinese Remainder Theorem (CRT) to construct a solution to the problem in the original domain. Although this approach is more complicated and less intuitive, as compared to a direct solution, it typically provides better performance [9], [10].

To use RNS, we need to first define the RNS base  $B = \{m_1, \dots, m_k\}$ , with  $k$  moduli s.t.  $m_i$  is an integer and  $\gcd(m_i, m_j) = 1, \forall i \neq j$ . The latter condition is required to guarantee the number system is not redundant and to allow using the CRT for RNS-to-positional-number-system conversion. An integer  $x$  can be represented in RNS with base  $B$  via the tuple  $x = \{x_1, \dots, x_k\}$ , where  $x_i = x \pmod{m_i}$ , also denoted by  $[x_i]_{m_i}$ . The quantity  $M = \prod_{i=1}^k m_i$  is known as the RNS dynamic range. As long as  $x < M$ , or  $[-M/2] \leq x < [M/2]$  as in our case, there is a unique RNS representation of  $x$ .

*b) Chinese Remainder Theorem:* CRT can be used to do the backward conversion, i.e., converting a number represented in RNS to its equivalent in a positional numbering system. Equation (1) can be used to perform this conversion. We remark that this procedure is not only serial, but also requires multi-precision arithmetic. Equation (1)

can be reformulated into two forms as shown in equations (2) and (3). It is straightforward to show that  $v$  and  $v'$ , which represent the extra multiples of  $M$ , have the upper bounds of  $k$  and  $k \cdot \max\{m_i\}$ , respectively.

$$[x]_M = \left[ \sum_{i=1}^k \left[ x_i \cdot \left( \frac{M}{m_i} \right)^{-1} \right]_{m_i} \cdot \frac{M}{m_i} \right]_M \quad (1)$$

$$x = \left( \sum_{i=1}^k \left[ x_i \cdot \left( \frac{M}{m_i} \right)^{-1} \right]_{m_i} \cdot \frac{M}{m_i} \right) - v \cdot M \quad (2)$$

$$x = \left( \sum_{i=1}^k x_i \cdot \left[ \left( \frac{M}{m_i} \right)^{-1} \right]_{m_i} \cdot \frac{M}{m_i} \right) - v' \cdot M \quad (3)$$

### C. Elementary Operations

In complexity analysis, we only focus on two elementary operations: Modular Multiplication (MM) and Floating-Point (FP) operations. We compute the total number of elementary operations in core procedures. We use Single Precision (SP) to refer to elementary operations that fit in one word, and Double Precision (DP) for those that require two words.

### D. Efficient RNS Base Extension

As seen later, BFV requires division and rounding in decryption and homomorphic multiplication. Since these operations are incompatible with RNS, they are handled by base extension techniques. In this subsection, we introduce two different techniques to extend an RNS base. Both techniques form the bases of the BFV RNS variants.

1) *Base Extension via Integer Arithmetic:* The first method utilizes Equation (2) as shown in Equation (4).

$$\text{FastBaseConv\_I}(x, B, B') = \left[ \sum_{i=1}^k \left[ x_i \cdot \left( \frac{M}{m_i} \right)^{-1} \right]_{m_i} \cdot \frac{M}{m_i} \bmod m'_j \right]_{m'_j \in B'} \quad (4)$$

The conversion is fast since it does not require multi-precision operations, given that the quantities  $Mm_i^{-1} \pmod{m'_j}$ ,  $\forall m'_j \in B'$  are precomputed. `FastBaseConv_I` can be thought of as computing the CRT without reduction modulo the dynamic range  $M$ . This means that the converted value may have multiples of  $M$  as an overflow, i.e.,  $\text{FastBaseConv\_I}(x, B, m'_j) = x + v \cdot M \pmod{m'_j}$ . The overflow can be controlled and sometimes eliminated via the techniques described in later sections. `FastBaseConv_I` forms the basic building block of BEHZ [9].

a) *Complexity:* `FastBaseConv_I` requires for each residue two modular multiplications: one modulo  $m_i \in B$  and another modulo  $m'_j \in B'$ . Note that the quantity  $\left[ x_i \cdot \left( \frac{M}{m_i} \right)^{-1} \right]_{m_i}$  can be temporarily stored and reused. Hence, the overall complexity is  $(kk' + k)$  MM, where  $k'$  denotes the number of moduli in base  $B'$ .

b) *A Practical Note:* We note that one may invoke a sequence of multiplications and additions without modulo reduction, deferring the latter to as late as possible. Eventually, a reduction should be applied. The length of the sequence can be determined if the moduli sizes are known in advance, which is typically the case. This optimization is similar to Harvey's lazy reduction [21] used to improve the computation of Number Theoretic Transform (NTT).

2) *Base Extension via Integer and Floating-Point Arithmetic:* The second method also utilizes Equation (2). However, instead of doing approximate conversion (conversion with overflow), it estimates  $v$  and uses the estimated value for exact conversion. One can reformulate Equation (2) to find  $v$  as follows:

$$v = \left[ \sum_{i=1}^k \frac{1}{m_i} \left[ x_i \cdot \left[ \left( \frac{M}{m_i} \right)^{-1} \right]_{m_i} \right]_{m_i} \right]$$

Given that  $v$  is estimated correctly, we can directly apply Equation (2) modulo  $m'_j$  to do exact conversion as follows:

$$\text{FastBaseConv\_F}(x, B, B') = \left[ \left( \sum_{i=1}^k \left[ x_i \cdot \left[ \left( \frac{M}{m_i} \right)^{-1} \right]_{m_i} \right]_{m_i} \cdot \frac{M}{m_i} \right) - v \cdot [M]_{m'_j} \right]_{m'_j \in B'} \quad (5)$$

The key issue in this method is to ensure that the estimated  $v$  is correct. Since  $v$  is estimated via floating-point operations, errors due to limited precision may occur and produce an off-by-one value of  $v$ . However, such cases can be controlled, detected, and corrected [22], [23], [10]. This base extension method forms the basic building block of HPS [10].

a) *Complexity*: FastBaseConv\_F can be decomposed into two steps: 1) estimating  $v$ , and 2) base extension to the new RNS base. Estimating  $v$  is done once and requires  $k$  MM and  $k+1$  FP. On the other hand, base conversion requires  $k' \cdot (k+1)$  MM. Overall complexity is:  $(kk' + k + k')$  MM +  $(k+1)$  FP.

b) *A Practical Note*: While computing  $v$ , the numerator can be temporarily stored and reused in the computation of Equation (5).

c) *Comparison*: It is straightforward to notice that FastBaseConv\_I is slightly more efficient than FastBaseConv\_F. However, it introduces extra multiples of base  $B$ , therefore, a correction procedure needs to be applied afterwards. On the other hand, FastBaseConv\_F does not require any correction step given that enough precision is used in estimating  $v$ . Our experiments later in this paper show that the correction techniques significantly increase the total runtime.

### E. The Textbook BFV

The textbook BFV scheme, described in [5], is a tuple of 5 procedures: key generation, encryption, decryption, homomorphic addition and multiplication. The scheme defines a set of parameters as follows:

- $\lambda$ : security parameter.
- $w$ : a decomposition base used to express a polynomial in  $R_q$  in terms of  $l+1$  polynomials in base  $w \in \mathbb{Z}$ , where  $l = \lfloor \log_w^q \rfloor$ .
- $\mathcal{X}_{err}$ : a zero-mean discrete Gaussian distribution used to sample error polynomials. The distribution is parameterized by the standard deviation  $\sigma$  and error bound  $\beta_{err}$ .
- $t \geq 2$ : a plaintext modulus.
- $q \gg t$ : a ciphertext modulus.

The main five procedures of the scheme are as follows:

- **KeyGen**( $\lambda, w$ ): The Secret Key (sk) is a ternary polynomial  $sk \xleftarrow{\mathcal{U}} R_2$  taking values from the set  $\{-1, 0, 1\}$ . The Public Key (pk) is a pair of polynomials  $(pk_0, pk_1) = (-[a \cdot sk + e]_q, a)$ , where  $a \xleftarrow{\mathcal{U}} R_q$  and  $e \xleftarrow{\mathcal{G}} \mathcal{X}_{err}$ . The Evaluation Key (evk) is a set of  $(l+1)$  pairs of polynomials generated as follows: for  $0 \leq i \leq l$ , sample  $a_i \xleftarrow{\mathcal{U}} R_q$  and  $e_i \xleftarrow{\mathcal{G}} \mathcal{X}_{err}$ .  $evk[i] = ([w^i s^2 - (a_i \cdot sk + e_i)]_q, a_i)$ . The procedure outputs the tuple:  $(sk, pk, evk)$ .
- **Enc**( $m, pk$ ): takes a plaintext message  $m \in R_t$ , and samples  $u \xleftarrow{\mathcal{U}} R_2$  and  $e_1, e_2 \xleftarrow{\mathcal{G}} \mathcal{X}_{err}$ . It produces the ciphertext  $ct = ([\Delta m + pk[0]u + e_1]_q, [pk[1]u + e_2]_q)$ , where  $\Delta = \lfloor q/t \rfloor$ .
- **Dec**( $ct, sk$ ): computes  $m = \left\lfloor \left[ \frac{t}{q} [ct[0] + ct[1]sk]_q \right] \right\rfloor_t$ .
- **EvalAdd**( $ct_0, ct_1$ ): homomorphic addition takes two ciphertexts and produces:  $ct_{add} = ([ct_0[0] + ct_1[0]]_q, [ct_0[1] + ct_1[1]]_q)$ .
- **EvalMul**( $ct_0, ct_1, evk$ ): homomorphic multiplication takes two ciphertexts and performs
  - 1) Tensoring: compute  $c_\tau$ , with  $\tau \in \{0, 1, 2\}$ , such that:

$$\begin{aligned} c_0 &= \left[ \left[ \frac{t}{q} ct_0[0] \cdot ct_1[0] \right] \right]_q, \\ c_1 &= \left[ \left[ \frac{t}{q} (ct_0[0] \cdot ct_1[1] + ct_0[1] \cdot ct_1[0]) \right] \right]_q, \\ c_2 &= \left[ \left[ \frac{t}{q} ct_0[1] \cdot ct_1[1] \right] \right]_q. \end{aligned}$$

2) Relinearization:

2.1) decompose  $c_2$  in base  $w$  as  $c_2 = \sum_{i=0}^l c_2^{(i)} w^i$ .

2.2) return  $ct_{mul}[j]$ , with  $j \in \{0, 1\}$ , such that:

$$ct_{mul}[j] = \left[ c_j + \sum_{i=0}^l evk[i][j] c_2^{(i)} \right]_q.$$

Key generation, encryption and homomorphic addition can be computed easily in RNS. In contrast, decryption and homomorphic multiplication are more involved as they require scaling by the ratio  $t/q$  and rounding. The second step in homomorphic multiplication, known as base decomposition, is incompatible with RNS and requires a computationally expensive conversion to the positional (multi-precision) representation. The variants studied in our work employ different techniques to overcome these problems.

### III. EFFICIENT RNS VARIANTS OF BFV DECRYPTION

#### A. Scaling in Decryption

Decryption includes scaling by the factor  $t/q$  followed by rounding. Both operations require special treatment so that they can be performed directly in RNS. We show below how BEHZ and HPS tackle them.

1) *RNS Decryption in BEHZ*: BEHZ uses an efficient simple scaling procedure to compute an approximate value of decrypted result. The approximation affects the noise threshold for successful decryption. Note that, in the textbook BFV, decryption works as long as  $\|e\|_\infty < (\Delta - |q|_t)/2$ , where  $e$  is a polynomial representing the noise contained in the ciphertext.

The main objective is to compute  $\lfloor [t/q \cdot [x]_q]_t \rfloor$ , where  $x = ct[0] + ct[1]sk$ . In BEHZ, exact flooring is used instead of rounding as follows:

$$\left\lfloor \frac{t}{q} [x]_q \right\rfloor = \frac{t[x]_q - |t \cdot x|_q}{q}.$$

This converts the problem into integer division that is compatible with RNS. Since we are computing modulo  $t$ , the term  $t[x]_q$  cancels out. Then they compute the term  $-|t \cdot x|_q/q$  using FastBaseConv\_I from base  $q$  to  $t$ . Since the conversion is not exact, multiples of  $q$  may be generated. The BEHZ variant introduces a redundant modulus  $\gamma$  to remove this overflow. Moreover,  $\gamma$  helps in correcting the errors due to the use of flooring instead of rounding. Algorithm 1 shows the RNS decryption function. The reader should note that the notation  $|\cdot|_m$  is represented by the standard interval  $\{0, \dots, m-1\}$ , whereas  $[\cdot]_m$  is represented by the centered interval  $\{[-m/2], \dots, [m/2]\}$ .

---

#### Algorithm 1 Dec<sub>RNS</sub>: BEHZ RNS decryption [9]

---

**Input:** ciphertext  $ct$ , secret key  $sk$  and a redundant modulus  $\gamma \in \mathbb{Z}$ .

**Output:** the plaintext message  $[m]_t$ .

```

1: for  $j \in \{t, \gamma\}$  do
2:    $s^{(j)} \leftarrow \text{FastBaseConv\_I}(|\gamma t \cdot (ct[0] + ct[1]sk)|_q, q, j) \times | -q^{-1}|_j \pmod{j}$ 
3:    $\tilde{s}^{(\gamma)} \leftarrow [s^{(\gamma)}]_\gamma$ 
4:    $m^{(t)} \leftarrow [(s^{(t)} - \tilde{s}^{(\gamma)}) \times |\gamma^{-1}|_t]_t$ 
5: return  $m^{(t)}$ 

```

---

a) *Correctness*: Algorithm 1 works as long as  $\|e\|_\infty \leq \Delta(\frac{1}{2} - \frac{k}{\gamma}) - \frac{|q|_t}{2}$ . To ensure that noise threshold is close to that in the textbook BFV,  $\gamma$  can be chosen much larger than  $k$  so that the quantity  $k/\gamma$  approaches 0.

b) *Complexity*: Algorithm 1 requires the invocation of FastBaseConv\_I on two moduli  $t$  and  $\gamma$ . It also requires one extra modular multiplication by  $|\gamma^{-1}|_t$ . The overall complexity is  $3k + 3$  MM for each coefficient.

2) *RNS Decryption in HPS*: HPS employs a simpler decryption variant that does not introduce an auxiliary modulus. It utilizes Equation (3) as follows:

$$\left\lfloor \frac{t}{q} [x]_q \right\rfloor = \left[ \left[ \left( \sum_{i=1}^k x_i \cdot \left[ \left( \frac{q}{q_i} \right)^{-1} \right]_{q_i} \cdot \frac{t}{q_i} \right) \right] \right]_t \quad (6)$$

Note that, the term  $v' \cdot t$  cancels out due to the reduction modulo  $t$ . The procedure computes the summation of residues multiplied with precomputed floating-point constants. For maximum precision, the authors suggest to split the floating-point quantities into integer and fractional parts, i.e.,  $\left[ \left( \frac{q}{q_i} \right)^{-1} \right]_{q_i} \cdot \frac{t}{q_i} = \omega_i + \theta_i$  with  $\omega_i \in \mathbb{Z}_t$  and  $\theta_i \in [-1/2, 1/2)$ .

a) *Correctness*: RNS decryption works as long as the rounding errors due to floating-point operations are controlled. This can be done by limiting the size of  $q_i$ 's, the number of  $q_i$ 's and the precision used to store  $\theta_i$ 's. Given that  $\theta'_i$  is a limited-precision approximation of the actual value of  $\theta_i$ , the approximation error is given by  $\epsilon_i = \theta'_i - \theta_i$ . If the total approximation error less than  $1/4$ , the decryption procedure works correctly. Halevi et al. provide a detailed analysis of the approximation error. We provide concrete bounds on the number and size of  $q_i$ 's.

Let  $y = \lfloor \sum_i x_i \theta_i \rfloor$  and  $y' = \lfloor \sum_i x_i (\theta_i + \epsilon_i) \rfloor$ , then the total error term  $\epsilon = \sum_i x_i \epsilon_i$ . Let  $\nu$  denote the precision used to store  $\theta_i$ 's, i.e.,  $\epsilon_i < 2^{-\nu}$ . We know that  $|x_i|_\infty < q_i/2$ , hence the total error  $\|\epsilon\|_\infty < 2^{-\nu} \sum_i \frac{q_i}{2} < 2^{-\nu-1} k \cdot \max(q_i)$ . As we need to ensure that  $\|\epsilon\|_\infty < 1/4$ , the constraint for  $q_i$  can be written as  $q_i \leq 2^{\nu-1}/k$  for the case when the modular arithmetic is implemented using signed integers. It can be similarly shown that in the case of an unsigned-integer implementation, the constraint changes to  $q_i \leq 2^{\nu-3}/k$ .

There is a trade-off between  $\nu$ ,  $\max(q_i)$ , and  $k$ . A higher precision  $\nu$  implies costlier computations. The values of  $k$  and  $\max(q_i)$  limit the highest value of  $q$  the implementation can support. For a given ciphertext modulus  $q$ , it is better to minimize  $k$  as much as possible to reduce the number of NTT invocations. In order to optimize the performance of floating-point operations, we identify four cases, as shown in Table I.

Table I: Common C++ floating-point data types, total size in bits, floating-point precision ( $\nu$ ) in bits, bit size ( $\log_2^{q_i}$ ) and maximum number ( $K$ ) of supported moduli for HPS decryption (for a signed-integer implementation of modular arithmetic).

C++ data type	size	$\nu$	$\log_2^{q_i}$	$K$
float	32	24	16	128
double	64	53	44	64
long double	80	65	59	32
double double	128	113	107	32

We remark that `float`, `double` and `long double` are native data types in C++ 1999 standard and implemented natively (in hardware) in modern systems. On the other hand, `double double` is not supported in hardware, and is usually implemented in software. In our implementation, we utilize Shoup's NTL `quad_float` data type [24]. We also remark that our CPU implementation supports  $30 \leq \log_2^{q_i} \leq 60$  bits, and we use an unsigned-integer implementation of modular arithmetic, i.e., we employ `quad_float` only when  $\log_2^{q_i} > 57$  bits. Our GPU implementation, on the other hand, supports only 30-bit moduli and hence uses only `double`.

b) *Complexity*: The procedure requires  $k$  MM and  $(k + 1)$  FP operations for each coefficient.

## B. Comparison and Evaluation of RNS Decryption

We have seen how each variant handles the scaling problem in decryption. The appealing feature in BEHZ is the use of integer arithmetic, therefore; one does not need to worry about rounding errors as in HPS. However, it has some drawbacks since it requires a redundant modulus which triples the computation since we work in different rings:  $\mathbb{Z}_{q_i}$ ,  $\mathbb{Z}_t$ , and  $\mathbb{Z}_\gamma$ . It also affects the noise level for correct decryption due to the use of `FastBaseConv_I` and flooring instead of rounding. In contrast, HPS does not require an extra modulus nor it affects the noise level. HPS is also more efficient when native precision floating-point operations are used, as seen in Figure 1. The main drawback, however, is the need for high-precision (`double double`) floating-point operations when the moduli size is higher than 57 (59) bits. Table II shows the computational complexity of decryption for each variant.

Table II: Decryption computational complexity for BEHZ and HPS.

BEHZ	HPS
$n(3k + 3)$ MM	$nk$ MM + $n(k + 1)$ FP

Now, we compare the storage complexity of precomputed constants for both RNS variants. Although this may not be important in CPU implementations, it is crucial for GPU. The reason is that we use the constant memory in GPU to store any precomputed quantities in our implementation. The GPU constant memory is very fast but limited in size. Current devices typically include 64 KB of constant memory [25]. In addition, the size of any allocated buffer in constant memory must be defined at compile time. Hence, we limit our GPU implementations to 64 CRT moduli. Table III lists the precomputed quantities and their sizes for the decryption in each RNS variant.

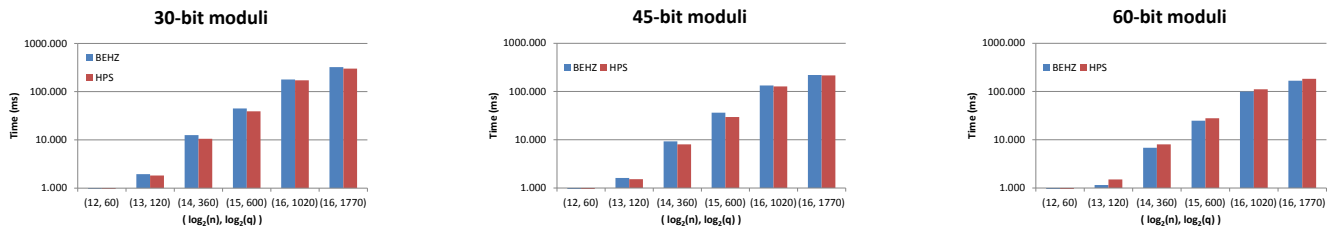


Figure 1: Effect of floating-point precision on decryption performance in PALISADE. C++ data types `double`, `long double`, and `NTL double double` are used with 30-bit, 45-bit, and 60-bit moduli, respectively. The horizontal axis represents  $\log_2$  of polynomial degree ( $n$ ) and ciphertext modulus ( $q$ ). Note that the vertical axes are in log-scale.

Table III: Decryption precomputed constants in both RNS variants. Size is given in terms of number of words and  $K$  denotes the maximum number of CRT moduli supported.  $\{\cdot\}$  denotes a fraction of  $\cdot$ .

Quantity	# of words	Required by BEHZ	Required by HPS
$t \cdot \gamma \cdot [(q/q_i)^{-1}]_{q_i}$	$K$	✓	-
$q/q_i \bmod t$	$K$	✓	-
$q/q_i \bmod \gamma$	$K$	✓	-
$\lfloor \frac{t \cdot [(q/q_i)^{-1}]_{q_i}}{q_i} \rfloor$	$K$	-	✓
$\{\frac{t \cdot [(q/q_i)^{-1}]_{q_i}}{q_i}\}$	$K$	-	✓

We see that in most cases HPS requires less memory for precomputed constants. However, if `double double` is used to store the floating-point quantities, then both procedures have the same storage complexity.

In order to show the effect of using a different floating-point precision on performance, Figure 1 shows decryption runtimes for the CPU implementations at various CRT moduli sizes. Three main ranges of CRT moduli sizes are identified for our CPU implementation of HPS: 1)  $< 45$ , 2)  $45 - 57$ , and 3)  $58 - 60$  bits. HPS decryption is faster for the CRT moduli range of  $30 - 57$  bits. However, for larger moduli sizes, the HPS performance degrades due to the use of `NTL double double`, and BEHZ becomes faster. A logical observation is that the HPS performance is faster as long as native floating-precision data types are used. We remark that the choice of parameters corresponds to at least 128 bits of security. The parameter selection is discussed in more detail in Section VI.

#### IV. EFFICIENT RNS VARIANTS OF BFV HOMOMORPHIC MULTIPLICATION

##### A. Scaling in Homomorphic Multiplication

Homomorphic multiplication is more complex than decryption and includes two main steps: (1) multiplication itself and (2) relinearization.

Step 1 in multiplication requires tensoring the input ciphertexts and scaling by the factor  $t/q$ , followed by rounding. Unlike decryption, the result should also be in base  $q$ , therefore, an RNS scaling algorithm is needed. Tensoring requires lifting the ciphertexts first from base  $q$  to a larger base. For that reason, an auxiliary RNS base  $q'$  with  $k'$  moduli is introduced, which is roughly as big as  $q$ . More concretely,  $q \cdot q'$  should be large enough to include the largest coefficients of tensored ciphertexts without any modular reduction. Lifting can be done using the RNS base extension techniques described previously. Next, the extended tensored ciphertexts are scaled down by  $q$ , which represents a subset of the moduli in the extended RNS base  $\{q \cup q'\}$ . An RNS scaling algorithm can be used for this purpose. Note that the scaled-down result is represented in base  $q'$ , therefore, another round of RNS base extension should be carried out to retain the result in base  $q$ .

In Step 2, relinearization is performed to reduce the ciphertext size. For relinearization, Bajard et al. [9] suggested to use the RNS representation of ciphertexts in base  $q$  instead of the standard digit decomposition procedure using the positional base  $\omega$ . The idea is that  $q_i$ 's and  $\omega$  are of the same size in practice; therefore, they can control the noise growth (due to relinearization) similarly. As ciphertexts are already represented in the RNS base  $q$ , the decomposition is free. However, the evaluation key must be modified to be compatible with RNS relinearization.



<b>Input:</b> $c_\tau = \{[c_\tau]_q, [c_\tau]_{q' \cup m_{sk}}\}$ in extended base $q \cup q' \cup m_{sk}$ , and plaintext modulus $t$	
<b>Output:</b> $[c_\tau^*]_{q' \cup m_{sk}}$ , with $c_\tau^* = \left\lfloor \frac{t}{q} c' \right\rfloor + b_\tau$ in $q' \cup m_{sk}$ , where $\ b_\tau\ _\infty \leq k$ , and $\tau \in \{0, 1, 2\}$	
Operations in base $q$	Operations in base $q' \cup m_{sk}$
0: $[c_\tau]_q$	$[c_\tau]_{q' \cup m_{sk}}$
1: $t \cdot [c_\tau]_q$	$t \cdot [c_\tau]_{q' \cup m_{sk}}$
2: FastBaseConv_I( $[t \cdot c_\tau]_q, q, q' \cup m_{sk}$ )	$\rightarrow \frac{[t \cdot c_\tau]_q + v_\tau q]_{q' \cup m_{sk}}}{q}$
3: —	$\frac{[t \cdot c_\tau]_{q' \cup m_{sk}} - ([t \cdot c_\tau]_q + v_\tau q)}{q}$

Figure 2: Fast RNS flooring in BEHZ

This can be done at the initialization phase without additional cost. We remark that HPS employs a similar approach for this step, but with minor optimizations.

The above logic is used as the blueprint for RNS homomorphic multiplication. Both variants follow this blueprint but rely on different RNS base extension techniques.

---

**Algorithm 2** SmMRq $_{\tilde{m}}$ : Small Montgomery Reduction mod  $q$  [9]

---

**Input:** polynomial  $c'' = [\tilde{m}c]_q + v \cdot q$  in  $q' \cup \{m_{sk}, \tilde{m}\}$

**Output:**  $c'$  in  $q' \cup m_{sk}$ , with  $c' \equiv c'' \tilde{m}^{-1} \pmod{q}$ ,  $\|c'\|_\infty \leq \frac{\|c''\|_\infty}{\tilde{m}} + \frac{q}{2}$ .

- 1:  $r_{\tilde{m}} \leftarrow [-c''_{\tilde{m}} \cdot q^{-1}]_{\tilde{m}}$
  - 2: **for**  $j \in q' \cup m_{sk}$  **do**
  - 3:  $c'_j \leftarrow |(c''_j + q r_{\tilde{m}}) \tilde{m}^{-1}|_j$
  - 4: **return**  $c'$  in  $q' \cup m_{sk}$
- 

1) *Lift-and-Scale in BEHZ:* Lift-and-scale in BEHZ is performed in 5 steps as follows:

- 1) Base extension from base  $q$  to base  $q' \cup \{m_{sk}, \tilde{m}\}$ : FastBaseConv\_I( $ct_i[j], q, q' \cup \{m_{sk}, \tilde{m}\}$ ),  $\forall 0 \leq i, j < 2$  is used to do the conversion efficiently. Since this may generate additional multiples of  $q$  (referred to as  $q$ -overflows), namely  $v \cdot q$ , the redundant modulus  $\tilde{m}$  is introduced for correction. Note that the correction procedure requires the input ciphertexts to be multiplied by  $\tilde{m}$  before the fast base conversion.
- 2)  $q$ -overflow correction: This procedure is known as the small Montgomery reduction (shown in Algorithm 2). Basically, the  $\tilde{m}$  residue is used to find  $v$ . Then we subtract  $v \cdot q \pmod{q'_i}$  from  $q' \cup m_{sk}$  residues. The result is a corrected extended RNS representation but with an extra term affecting the noise growth. Note that the choice of  $\tilde{m}$  is not arbitrary. It should satisfy Equation (7). In terms of complexity, Algorithm 2 requires  $2(k' + 1) + 1$  MM for each coefficient.

$$\tilde{m}\rho \geq 2k + 1. \quad (7)$$

We remark that in our implementation, we use a standard CRT modulus for  $\tilde{m}$ , i.e., similar in size to the CRT moduli in  $q$  and  $q'$ . Therefore,  $\rho$  tends to approach 0. As will be shown later,  $\rho$  is used to estimate the noise growth due to the fast conversion.

- 3) Computation of tensor product: At this point, we have the ciphertexts in base  $\{q \cup q' \cup m_{sk}\}$ . Hence we can perform polynomial multiplication in  $\mathbb{Z}_{q \cdot q' \cdot m_{sk}} / (X^n + 1)$ . We remark that  $2(k + k' + 1)$  NTT and  $(k + k' + 1)$  INTT invocations are required for each polynomial multiplication. Namely, we need 4 polynomial multiplications but they can be reduced to 3 using the Karatsuba algorithm. Note that  $\|c_\tau\|_\infty \leq \delta \frac{q}{2} (1 + \rho)^2$ , where  $\tau \in \{0, 1, 2\}$  and  $\delta = \sup \|a \cdot b\|_\infty / (\|a\|_\infty \cdot \|b\|_\infty)$  is known as the ring expansion factor. Following a conservative (worst-case) approach,  $\delta$  is set to  $n$ ; however, in practical settings  $\delta$  can often be set to  $2\sqrt{n}$  if the requirements for the Central Limit Theorem are met [10].
- 4) Approximate rounding in base  $q'$ : This step is used to scale down  $c_\tau$  by the factor  $t/q$ . This gives us flooring

instead of rounding, i.e.,  $\lfloor \frac{t}{q} c_\tau \rfloor$ . Although this produces approximate results, Bajard et al. show that the error due to approximation is very small. This procedure works as shown in Figure 2. In terms of complexity, the procedure is invoked 3 times. In each invocation, FastBaseConv\_I is invoked  $(k' + 1)$  times. The multiplication by  $t$  incurs  $(k + k')$  modular multiplications. Lastly,  $(k' + 1)$  modular multiplications are required for the flooring. Note that the above is repeated for each coefficient.

- 5) Base extension from base  $q'$  to base  $q$ : This step is similar to step (1) above to convert to the original RNS base. The extra modulus  $m_{sk}$  is used to correct the overflows generated by FastBaseConv\_I. This is done via Shenoy-Kumaresan exact base extension [26].

*a) Remarks:* The choice of  $\rho, \tilde{m}, m_{sk}$ , and  $\gamma$  is not quite straightforward. For instance, one may choose  $\rho \approx 2k$  to avoid the usage of Algorithm 2. However, this can only be applied to a limited number of circuits where the multiplicative depth is small. On the other hand, higher values of  $\rho$  proportionally increase the noise growth and may reduce the multiplicative depth. Another aspect one needs to consider is that some choices of  $\tilde{m}$  may require very large values of  $m_{sk}$  and  $\gamma$ , which may not fit in the machine word size. In our implementation, we choose  $\tilde{m}, m_{sk}$ , and  $\gamma$  as regular CRT moduli, i.e., close in size to  $q_i$ 's.

2) *Lift-and-Scale in HPS:* HPS lift-and-scale follows the same blueprint as BEHZ. However, it is much simpler and requires no correction tools. The reason is that base extension is exact and does not generate extra multiples of the dynamic range. Therefore, the small Montgomery reduction algorithm is not required. Likewise, the Shenoy-Kumaresan CRT extension is not necessary in step 5. One only needs to worry about the rounding errors due to floating-point operations. Below we provide an overview of the lift-and-scale in HPS.

Lift-and-scale in HPS is performed in 4 steps as follows:

- 1) Base extension from base  $q$  to  $q'$ : FastBaseConv\_F ( $ct_i[j], q, q'$ ),  $\forall 0 \leq i, j < 2$  is used to lift the ciphertexts to a larger RNS base  $q'$ .
- 2) Computation of tensor product: This step is similar to step 3 in BEHZ. We remark that the number of NTTs in this step is less by 7, as compared to BEHZ, since there is no redundant modulus  $m_{sk}$ .
- 3) Exact rounding in base  $q'$ : In this step,  $c_\tau$  is scaled down by the factor  $t/q$ . This can be done by a procedure similar to the one used in HPS decryption.  $c_\tau$  is bounded by  $qq'/2t^2$ . The scaling is done by replacing  $t$  and  $q$  in Equation (6) by  $tq'$  and  $qq'$  respectively. Applying Equation (6) gives us  $\lfloor [tq'/qq' \cdot c_\tau] \rfloor_{q'}$ . Since  $c_\tau$  is bounded by  $qq'/2t$ , then  $\lfloor tq'/qq' \cdot c_\tau \rfloor \in [-q'/2, q'/2)$ . The complete procedure is given in Equation (8), where  $x$  is a single coefficient in  $c_\tau$ . Note that the quantity  $\frac{1}{q_i} \cdot tq' \lfloor (\frac{qq'}{q_i})^{-1} \rfloor_{q_i}$  is precomputed and broken into integral and fractional parts. The same applies to  $\left\lfloor t \left( \frac{qq'}{q_j} \right)^{-1} \cdot \frac{q'}{q_j} \right\rfloor_{q_j}$ . In terms of complexity, this procedure requires  $(k + 1)$  FP +  $k'(k + 1)$  MM for each coefficient.

$$\lfloor [t/q \cdot x] \rfloor_{q_j'} = \left[ \left[ \sum_{i=1}^k x_i \cdot \frac{1}{q_i} \cdot tq' \lfloor (\frac{qq'}{q_i})^{-1} \rfloor_{q_i} \right] + x_j' \cdot \left[ t \left( \frac{qq'}{q_j'} \right)^{-1} \cdot \frac{q'}{q_j'} \right]_{q_j'} \right]_{q_j'} \quad (8)$$

- 4) Exact base extension from base  $q'$  to base  $q$  using FastBaseConv\_F. Note that there is no extra correction step here.

## B. Comparison and Evaluation of RNS Lift-and-Scale

*a) Complexity:* It is not straightforward to compare the computational complexity between the variants. The employed tools use different elementary operations. However, we provide below an estimated analysis of the computational complexity, leaving the more precise empirical analysis for Section VI.

We only evaluate the lift-and-scale complexity as both variants use the same relinearization procedure. Table IV summarizes the computational complexity of the RNS and tensoring operations in the lift-and-scale procedure for BEHZ and HPS. We assume that the Karatsuba multiplication algorithm is used in tensoring, hence the factor 3 instead of 4 in tensoring, RNS rounding, and FastBaseConv( $x', q', q$ ). Note that FastBaseConv( $x', q', q$ ) in BEHZ includes the computational complexity of the Shenoy-Kumaresan algorithm as well.

<sup>2</sup>Note that  $t$  in the denominator is due to plaintext scaling in encryption.

Table IV: The computational complexity of the RNS tools in lift-and-scale (assuming the Karatsuba multiplication technique is used in tensoring).

Proc	BEHZ	HPS
FastBaseConv( $x, q, q'$ )	$4nk(k' + 2)$ MM	$4n(k(k' + 1) + k')$ MM + $4n(k + 1)$ FP
SmMR $q_{\tilde{m}}$	$4n(2(k' + 1) + k + 1)$ MM	—
Tensoring	$7(k + k' + 1)$ NTT + $3n(k + k' + 1)$ MM	$7(k + k')$ NTT + $3n(k + k')$ MM
RNS Rounding	$3n(k'(k + 2) + 2(k + 1))$ MM	$3nk'(k + 1)$ MM + $3n(k + 1)$ FP
FastBaseConv( $x', q', q$ )	$3n(k(k' + 1) + 2k' + 1)$ MM	$3n(k'(1 + k) + k)$ MM + $3n(k' + 1)$ FP
Total	$n(10k'k + 24k + 23k' + 24)$ MM + $7(k + k' + 1)$ NTT	$n(10k'k + 10k + 13k')$ MM + $7(k + k')$ NTT + $n(7k + 3k' + 10)$ FP

Table IV suggests that the number of modular multiplications for HPS is less by  $n(14k + 10k' + 24)$  and the number of NTTs is less by 7. However, HPS has a floating-point cost of  $n(7k + 3k' + 10)$  operations. We discuss the effect of these differences on homomorphic multiplication runtime in Section VI.

*b) Effect on Noise Growth:* Noise growth can be analyzed using the same logic as applied to the textbook BFV in [8] and YASHE in [27]. Both the BEHZ and HPS papers provide a detailed noise analysis. We only review some closed-form solutions that describe the noise growth in each variant.

**Noise growth in BEHZ [9]:** Initial noise in a fresh ciphertext is  $V = \beta_{err}(1 + 2\delta\beta_{key})$ . To ensure the correctness of a depth- $L$  binary tree multiplication, the maximum noise  $C_1^L V + LC_1^{L-1}C_2$  must be less than  $\frac{q}{t}(\frac{1}{2} - \frac{k}{\gamma}) - \frac{|q|_t}{2}$ , where

$$C_1 = \delta^2 t(1 + \rho)\beta_{key} + \delta t(4 + \rho) + \frac{\delta}{2},$$

$$C_2 = (1 + \delta\beta_{key})(\delta t|q|_t \frac{(1 + \rho)}{2} + \delta\beta_{key}(k + \frac{1}{2})) +$$

$$2\delta t|q|_t + k(2\delta\beta_{err}l_{\omega, 2^\nu}\omega + 1) + \frac{1}{2}(3 + |q|_t).$$

Here,  $l_{\omega, 2^\nu}$  is the number of base- $w$  digits in a CRT modulus  $q_i$  (for the second level of decomposition in relinearization).

**Noise growth in HPS [10]:** To guarantee the correctness of a depth- $L$  binary tree multiplication for HPS, the maximum noise  $C'_1 V + LC'_1^{L-1}C'_2$  must be less than  $(\Delta - r_t(q))/4$ , where  $r_t(q) = t(q/t - \Delta)$  and

$$C'_1 = (1 + \frac{5}{\delta\beta_{key}})\delta^2 t\beta_{key},$$

$$C'_2 = \delta^2 \beta_{key}((1 + 2\|\epsilon_s\|_\infty)\beta_{key} + t^2) + \delta\beta_{err}l_{\omega, 2^\nu}\omega k.$$

This constraint is similar to the textbook BFV case [8]; HPS adds at most two extra bits to the textbook BFV constraint, as shown in [10].

**Noise Growth Comparison:** The most significant quantity in noise growth is  $C_1$  (or  $C'_1$  for HPS). In BEHZ, an extra factor of  $(1 + \rho)$  is introduced, which can be minimized if  $\rho$  is small. This can be controlled by choosing a large  $\tilde{m}$  using Equation (7). Note that  $\beta_{key} = 1$ . We use the noise growth bounds above to find the maximum multiplicative depth  $L_\circ$  supported under a given parameter set as shown in Table V. Note that in this experiment, we set  $\delta = n$ , which corresponds to the worst-case analysis. It can be seen that the BEHZ and HPS RNS techniques have almost no effect on  $L_\circ$  for these settings. Note that the behavior in BEHZ is attributed to our choice of the redundant moduli, which is different from the parameter selection in [9]. We use standard CRT moduli for all redundant moduli to minimize the effect of the moduli on noise growth. We remark that the practical noise growths observed in our experiments are significantly different, as discussed later in this paper. We remark that we performed the same comparison of BEHZ and HPS noise growth as in Table V for 60-bit moduli, and we found the results to be very similar to the 30-bit case.

*c) Precomputed Constants:* The number of precomputed constants for homomorphic multiplication is quite large. Due to space constraints, we only list the numbers of vectors and matrices required by each variant in Table VI. Note that our GPU implementations include other precomputed constants related to CRT and NTT computation. We maximized  $K$  to use the largest amount of the GPU constant memory. It can be seen that BEHZ includes a larger number of parameters but requires less storage.

Table V: Parameters of the BFV scheme and its RNS variants with plaintext modulus  $t$ , and maximum theoretical multiplicative depth  $L_{\circ}$  for each variant. Parameters are generated to provide at least 128 bits of security, with the worst-case bound for the expansion factor  $\delta = n$ . The choice of  $t = 65537$  is motivated by the fact that the plaintext space for most of the practical settings can be decomposed into linear factors modulo 65537, which provides maximum number of slots for SIMD-like execution.

$\log_2^n$	$\lceil \log_2^q \rceil$	$t$	$L_{\circ}$		
			BEHZ	HPS	Textbook
12	60	2	1	1	1
		65537	0	0	0
13	120	2	3	3	3
		65537	2	2	2
14	360	2	11	11	11
		65537	7	7	7
15	600	2	18	18	18
		65537	12	12	12
16	1020	2	29	29	29
		65537	20	20	20
16	1770	2	52	52	52
		65537	36	36	36

## V. IMPLEMENTATION

### A. CPU Implementation

We implemented the BEHZ variant in the PALISADE library, which already provides an implementation of HPS (starting with version 1.1). We added the parameter generation, decryption, and homomorphic multiplication and RNS tools for the BEHZ variant. Other primitives, such as key generation, encryption, and homomorphic addition, were borrowed from the existing HPS implementation. The details of the HPS implementation are provided in [10]. Our implementation of the BEHZ variant is publicly accessible (included in PALISADE starting with version 1.2).

Table VI: Memory requirement for precomputed constants used in homomorphic multiplication.  $K$  denotes the largest number of CRT moduli supported by the GPU implementation. Memory size is given in KB.

Item	BEHZ	HPS
$K$	64	61
# of vectors size	12 3 KB	7 1.67 KB
# of matrices size	2 32 KB	3 43.61 KB
total size	35 KB	45.27 KB

Multi-threading in our CPU implementations is achieved via OpenMP<sup>3</sup>. The loop parallelization in the scaling and RNS base extension operations is applied at the level of single-precision polynomial coefficients (w.r.t.  $n$ ). The loop parallelization for NTT and component-wise vector multiplications (polynomial multiplication in the evaluation representation) is applied at the level of CRT moduli (w.r.t.  $k$ ).

We also added a new optimization in this work to improve the performance of both BEHZ and HPS in PALISADE:

a) *Lazy Reduction*: Lazy reduction can be used to reduce the number of modular reductions in a sum of products. Suppose we have a summation  $\sum_i a \cdot b \pmod{p}$ . If we can determine the upper bound of  $a \cdot b$ , we may use regular multiplications and additions as long as the bound is not reached. As soon as we approach the bound, we can apply one modular reduction modulo  $p$  and continue. In PALISADE, the largest efficient CRT modulus is 60 bits long, i.e.,  $a$  and  $b$  are also 60-bit numbers. If we use a 128-bit multiplier, we can do 128 multiply-add (MULADD) operations before we overflow a 128-bit register. Thus, a single modular reduction after 128 MULADDs is sufficient to compute the sum. To use lazy reduction, we extended PALISADE by adding two

<sup>3</sup><http://www.openmp.org/>

---

**Algorithm 3** 128-bit Barrett reduction modulo 60-bit integer.

---

**Input:** 128-bit  $a \in \mathbb{Z}^+$ , 60-bit  $p \in \mathbb{Z}^+$ , and  $\mu = \lfloor 2^{128}/p \rfloor$ .

**Output:**  $r = a \pmod{p}$ .

```
1:  $l_{hi} = \text{Mul128}(a_{lo}, \mu_{lo}) \ggg 64$ 
2:  $m = \text{Mul128}(a_{lo}, \mu_{hi})$ 
3:  $(s_0, c_{out}) = \text{AddwCarry}(m_{lo}, l_{hi})$  ▷ addition with carry
4:  $s_1 = m_{hi} + c_{out}$ 
5:  $m = \text{Mul128}(a_{hi}, \mu_{lo})$ 
6:  $(NULL, c_{out}) = \text{AddwCarry}(m_{lo}, s_0)$  ▷ check carry
7:  $l_{hi} = m_{hi} + c_{out}$ 
8:  $s_0 = a_{hi} * \mu_{hi} + s_1 + l_{hi}$ 
9:  $r = a_{lo} - s_0 * p$ 
10: while ( $r \geq p$ ) do
11:    $r -= p$ 
12: return  $r$ 
```

---

procedures: a 64-bit-by-64-bit multiplier to produce a 128-bit result, and a 128-bit modular reducer using Barrett algorithm. Although we replace a 64-bit modular MULADD with a 128-bit MULADD, experiments show that the latter is more efficient. Algorithm 3 shows our 128-bit Barrett reduction algorithm.

### B. GPU Implementation

The GPU library, referred to as DSI\_BFV [17], already includes a full implementation of BEHZ. We report here some of the recent features added to DSI\_BFV and details of our implementation of HPS.

DSI\_BFV includes two main components: 1) lattice cryptography library and 2) implementation of BEHZ. We should note that a key feature in DSI\_BFV is that it executes the entire BFV computation on GPU. The CPU is merely responsible for launching kernels and memory allocations/deallocations. Even the cryptographic keys are computed on GPU. This is slightly different from the normal processor-coprocessor model where only intensive tasks are offloaded to coprocessor. As we saw previously, BFV and its RNS variants include a large level of parallelism that is suitable for vector processors such as GPUs. We avoid frequent costly memory copying between CPU and GPU by performing the entire computation on GPU. This paradigm is also suitable for cluster GPUs where the CPU may distribute a large homomorphic circuit to multiple GPUs.

The lattice library in DSI\_BFV includes a set of tools described as follows.

a) *CRT/RNS*: Currently, for CRT/RNS, DSI\_BFV includes 30-bit moduli generated in a special form to support lazy reductions in NTT computations [21]. It also employs fixed-size primes generated according to the design in NTLlib [28]. For CRT reconstruction, DSI\_BFV uses Garner’s mixed radix algorithm since it requires less memory for precomputed constants and is faster than the classic CRT reconstruction [17]. For RNS arithmetic, we launch  $k \cdot n$  threads over 2D thread blocks to exploit the maximum parallelism provided by RNS, which is a scalable solution that can benefit from as many computational cores as are available in the GPU.

We decided to support only 30-bit CRT moduli on GPU for the following reasons:

- GPUs support natively 32-bit instructions as the internal register size is 32 bits. Though 64-bit operations are supported, they are simulated via 32-bit instructions [29].
- Certain types of memory in NVIDIA GPUs are optimized for 32-bit words. For instance, shared memory in GPU is optimized for 32-bit access. 64-bit access can cause shared memory bank conflicts. Although these bank conflicts can be eliminated by padding extra unused shared memory, this will consume more hardware resources and limit the number of concurrent thread blocks. Note that our GPU implementation uses heavily this special type of memory to facilitate intra-block communication among threads in a block.
- GPUs do not include a native 128-bit data type that can be used to multiply 64-bit operands. Hence, this should be implemented manually and simulated via 32-bit operations. Efficient implementation would require an optimized PTX (assembly-like) instructions that are complex, error prone and not hardware-oblivious.

*b) Discrete Galois Transform (DGT):* The DGT is used for efficient polynomial multiplication using negacyclic convolution. It was found suitable for GPUs and memory-bound platforms as it cuts the transform length into half and requires less amount of memory for precomputed twiddle factors. The DGT algorithm was originally proposed by Crandall [30] for fast negacyclic convolution. It works in the field  $GF(p^2)$  where  $p$  is a Gaussian prime, i.e.,  $p \equiv 3 \pmod{4}$ . Unfortunately, NTLlib primes are non-Gaussian and hence they are not compatible with vanilla DGT. However, Al Badawi et al. [31] showed how to extend Crandall’s DGT algorithm to work with non-Gaussian primes as well. We remark that the current version of DSI\_BFV includes the NTT optimizations from David Harvey [21] adapted for DGT computations.

*c) Uniform and Gaussian Random Samplers:* Random polynomials are generated on GPU using CUDA cuRAND. We need to sample polynomials from three random distributions: 1)  $\mathcal{X}_2$ , 2)  $\mathcal{X}_q$ , and 3)  $\mathcal{X}_{err}$ . The first two are uniform distributions while the third is a discrete Gaussian. DSI\_BFV includes efficient parallel implementations for these distributions.

*d) Memory Pool:* DSI\_BFV includes a GPU memory pool. Memory allocations are done once when needed, and the used memory is deallocated only when the program terminates. This is important for performance since we are dealing with large polynomials that require substantial amounts of memory. Frequent memory allocations and deallocations on GPU are costly and should be avoided whenever possible.

*e) Implementation of HPS:* The main component of HPS is FastBaseConv\_F. For a polynomial of degree  $n$ , we launch  $n$  threads to extend a polynomial coefficient in base  $B$  to base  $B'$ . The required constants are precomputed and stored in the GPU constant memory. We note here that although HPS requires less precomputed parameters, their size is larger than that required by BEHZ. We had to reduce  $K$  from 64 to 61 in order to fit the precomputed parameters in the GPU constant memory. Since we are only dealing with 61 30-bit CRT moduli, we use the native C++ double for floating-point operations. This provides enough precision as shown in Table I. Our FastBaseConv\_F kernel launches  $n$  threads, one for each coefficient residue. We also apply the lazy reduction technique similar to the CPU implementation. The only difference is that we work with 30-bit moduli; therefore, the GPU Barrett reduction algorithm is simpler than Algorithm 3. Key generation, encryption, homomorphic addition, and relinearization procedures are borrowed from the implementation of BEHZ that is already included in DSI\_BFV.

*f) Performance Tuning:* In order to explain our performance fine-tuning methodology for GPU, we first provide an overview of some concepts related to CUDA programming and NVIDIA GPU hardware. A GPU card includes a set of streaming multi-processors (SMs), for example, V100 includes 80 SMs. Each SM is essentially a vector processor that comprises a set of scalar processors (SPs) (each SM in V100 includes 64 SPs), register file and shared memory. The SPs are the GPU cores, or hardware threads, that perform mathematical and logical operations.

An important performance metric in CUDA is the achieved occupancy factor [32] defined as the ratio between the active warps<sup>4</sup> and the maximum active warps a device can support. The latter is a hardware specification that cannot be controlled, whereas the former depends on the availability of hardware resources, memory instructions that stall warps, the number of thread blocks launched and thread block size (number of threads per block). What can generally be controlled among these factors are the number of blocks and their size. For each kernel, one should aim at maximizing the achieved occupancy until at some point the performance starts degrading (since more active warps may limit the resources a thread can use). In our implementation, we fine-tuned the GPU kernels to ensure optimum performance on both V100 and K80. The fine-tuning is a manual process that can be done by intensive profiling studies while varying the number of blocks used and their dimensions, and observing the performance.

We provide below an example showing an estimate for the average number of cores used over all kernels for V100. We profiled the implementation for running the five HPS primitives (KeyGen, Enc, Dec, EvalAdd and EvalMul) under the settings  $(n, \log q, t) = (2^{16}, 1770, 2^{16} + 1)$ . The average achieved occupancy was found to be 0.628, i.e., 40.192 active warps per a SM (since the maximum active warps for V100 is 64). Hence we have  $40.192 * (\text{warp\_size} = 32) \approx 1286$  concurrent threads on average per SM. Note that  $1286 > (\text{number of SPs in a SM, that is, } 64 \text{ in V100})$ , which suggests that cores on average are fully utilized and they are all being used. Moreover, we can conclude that increasing the number of cores (with proportionally added hardware resources) is expected to improve the performance.

<sup>4</sup>A warp is a bundle of 32 threads that forms the minimal execution unit in CUDA.

## VI. PERFORMANCE EVALUATION

### A. Methodology

We report the execution times for decryption and homomorphic multiplication. On CPU, we measure the time via the C++ library chrono. On GPU, CUDA events are used instead.

We perform experiments in different platform settings. For CPUs, we use single-threaded and multi-threaded (OpenMP) settings. For GPUs, we run our experiments on one GPU card.

### B. Experimental Setup

Our GPU implementation was developed via CUDA 9.0 on a 64-bit server equipped with 2 sockets, 26 cores per socket and 2 logical CPUs per core, i.e., 104 CPU threads in total. The machine also hosts two NVIDIA cards: 1) Tesla K80 and 2) V100-PCIe. Table VII describes the hardware configuration of both CPU and GPU.

The OS was ArchLinux version (4.15.13-1-ARCH), and the compilers were g++ (GCC) 7.3.1 and nvcc (8.0.61). We disabled the PALISADE library OpenMP support in single-threaded experiments and used 26 threads in multi-threaded experiments.

Table VII: CPU and GPU hardware configurations

Feature	CPU	GPU	
		K80	V100
Model	Intel(R) Xeon(R) Platinum	K80	V100-PCIe
# Cores	104	2496	5120
Frequency	2.10 GHz	0.82 GHz	1.380 GHz
RAM	187.5 GB	12 GB	16 GB

### C. Parameter Selection

To choose the ring dimension  $n$ , we ran the Learning With Errors security estimator<sup>5</sup> (commit f59326c) [33] to find the lowest security levels for the uSVP, decoding, and dual attacks following the standard homomorphic encryption security recommendations [34]. We selected the least value of the number of security bits  $\lambda$  for all 3 attacks on classical computers based on the estimates for the BKZ sieve reduction cost model.

Table VIII: Decryption latency in (milliseconds) of BFV RNS variants for single- and multi-threaded CPU settings, and GPUs with different moduli sizes ( $\nu$ ).

Configuration	Variant	$(\log n, \log q)$					
		(12,60)	(13,120)	(14,360)	(15,600)	(16,1020)	(16,1770)
Single-threaded CPU ( $\nu = 30$ )	BEHZ	0.535	2.079	13.330	48.547	196.266	336.617
	HPS	<b>0.472</b>	<b>1.857</b>	<b>11.374</b>	<b>43.334</b>	<b>172.916</b>	<b>318.942</b>
Multi-threaded CPU ( $\nu = 30$ )	BEHZ	<b>0.298</b>	0.831	3.799	12.397	49.850	88.845
	HPS	0.302	<b>0.819</b>	<b>3.607</b>	<b>12.042</b>	<b>49.741</b>	<b>87.774</b>
Single-threaded CPU ( $\nu = 60$ )	BEHZ	<b>0.401</b>	<b>1.153</b>	<b>7.379</b>	<b>26.389</b>	<b>103.034</b>	<b>173.313</b>
	HPS	0.504	1.632	8.735	30.314	120.576	195.398
Multi-threaded CPU ( $\nu = 60$ )	BEHZ	0.278	<b>0.64</b>	<b>2.311</b>	<b>7.189</b>	<b>24.204</b>	<b>46.708</b>
	HPS	<b>0.273</b>	0.669	2.392	7.475	25.811	48.195
K80 GPU ( $\nu = 30$ )	BEHZ	0.115	0.139	0.304	0.558	1.564	2.630
	HPS	<b>0.111</b>	<b>0.123</b>	<b>0.235</b>	<b>0.455</b>	<b>1.207</b>	<b>2.025</b>
V100 GPU ( $\nu = 30$ )	BEHZ	0.057	0.063	0.101	0.134	0.329	0.516
	HPS	<b>0.054</b>	<b>0.059</b>	<b>0.087</b>	<b>0.111</b>	<b>0.298</b>	<b>0.457</b>

<sup>5</sup><https://bitbucket.org/malb/lwe-estimator>

Table IX: Homomorphic multiplication (including relinearization) latency in (milliseconds) of BFV RNS variants in single- and multi-threaded CPU settings, and GPUs with different moduli sizes ( $\nu$ ).

Configuration	Variant	$(\log n, \log q)$					
		(12,60)	(13,120)	(14,360)	(15,600)	(16,1020)	(16,1770)
Single-threaded CPU ( $\nu = 30$ )	BEHZ	10.157	<b>40.675</b>	<b>382.149</b>	1899.136	11761.077	35173.622
	HPS	<b>10.065</b>	41.417	385.057	<b>1873.935</b>	<b>10815.283</b>	<b>27229.210</b>
Multi-threaded CPU ( $\nu = 30$ )	BEHZ	4.270	10.838	74.716	351.296	1986.586	5697.640
	HPS	<b>4.054</b>	<b>10.179</b>	<b>74.420</b>	<b>351.140</b>	<b>1984.211</b>	<b>5553.058</b>
Single-threaded CPU ( $\nu = 60$ )	BEHZ	6.952	23.300	155.365	<b>670.946</b>	3526.113	9260.904
	HPS	<b>6.326</b>	<b>22.088</b>	<b>154.350</b>	673.132	<b>3464.492</b>	<b>8605.612</b>
Multi-threaded CPU ( $\nu = 60$ )	BEHZ	3.343	7.325	32.244	<b>124.117</b>	<b>585.080</b>	<b>1653.966</b>
	HPS	<b>2.844</b>	<b>6.834</b>	<b>31.744</b>	126.945	603.554	1667.919
K80 GPU ( $\nu = 30$ )	BEHZ	2.166	2.754	9.603	25.885	112.062	307.531
	HPS	<b>1.977</b>	<b>2.517</b>	<b>7.834</b>	<b>21.924</b>	<b>96.151</b>	<b>255.502</b>
V100 GPU ( $\nu = 30$ )	BEHZ	0.997	1.178	2.412	5.705	22.848	59.473
	HPS	<b>0.859</b>	<b>1.012</b>	<b>2.010</b>	<b>4.826</b>	<b>18.725</b>	<b>50.779</b>

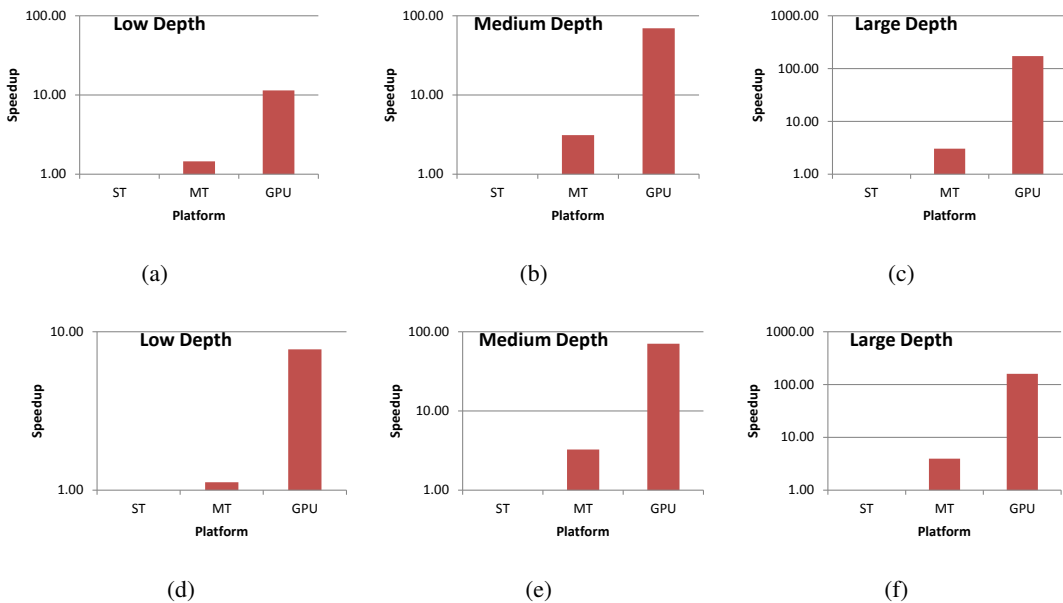


Figure 3: Speedup factors of homomorphic multiplication on different platforms (CPU single-threaded (ST), CPU multi-threaded (MT), and GPU) for different parameter sets:  $(\log_2 n, \log_2 q) = (12, 60)$  [low],  $(14, 360)$  [medium],  $(16, 1770)$  [large]. Upper three are for BEHZ, while the lower three are for HPS. Note that vertical axes are in log-scale.

The secret-key polynomials were generated using discrete ternary uniform distribution over  $\{-1, 0, 1\}^n$ . In all of our experiments, we selected the minimum ciphertext modulus bitwidth that satisfied the correctness constraint for the lowest ring dimension  $n$  corresponding to the security level  $\lambda \geq 128$ .

We set the Gaussian distribution parameter  $\sigma$  to  $8/\sqrt{2\pi}$  [34], the error bound  $B_e$  to  $6\sigma$ , and the lower bound for  $q'$  to  $2tnq$ . In the relinearization procedure, we utilized only the CRT decomposition (and did not use the second-level digit decomposition of residues).

#### D. Benchmarking

Tables VIII and IX show the results for decryption and homomorphic multiplication, respectively, in BEHZ and HPS. As PALISADE supports CRT moduli sizes  $\nu \in \{30, \dots, 60\}$  bits, we include the runtimes for  $\nu = 30$  and  $\nu = 60$ . Note that DSI\_BFV only supports 30-bit moduli.

It can be clearly seen that HPS outperforms BEHZ in decryption for  $\nu = 30$  bits on all platforms. When  $\nu = 60$  bits, HPS decryption performance degrades due to the use of quad-float (double double) floating-point



arithmetic. Given that HPS performs faster in three different platforms when  $\nu = 30$  bits, it can be concluded that HPS’s decryption function requires less computational overhead compared to BEHZ (unless quad-float or higher-precision floating-point arithmetic is needed). More concretely, HPS decryption can achieve the following speedup factors: 1.06 to 1.17 (CPU<sub>ST</sub>,  $\nu = 30$ ), 0.71 to 0.89 (CPU<sub>ST</sub>,  $\nu = 60$ ), 0.99 to 1.05 (CPU<sub>MT</sub>,  $\nu = 30$ ), 0.94 to 1.02 (CPU<sub>MT</sub>,  $\nu = 60$ ), 1.04 to 1.30 (GPU<sub>K80</sub>,  $\nu = 30$ ), and 1.06 to 1.21 (GPU<sub>V100</sub>,  $\nu = 30$ ), as compared to BEHZ.

For homomorphic multiplication, HPS typically outperforms BEHZ in all platforms regardless of  $\nu$ . Significant improvements can be noticed for large parameters. For instance, in CPU<sub>ST</sub> and  $\nu = 30$  bits, a 7.95-second difference is recorded. More concretely, HPS homomorphic multiplication can achieve the following speedup factors: 0.98 to 1.29 (CPU<sub>ST</sub>,  $\nu = 30$ ), 1.00 to 1.10 (CPU<sub>ST</sub>,  $\nu = 60$ ), 1.00 to 1.06 (CPU<sub>MT</sub>,  $\nu = 30$ ), 0.97 to 1.18 (CPU<sub>MT</sub>,  $\nu = 60$ ), 1.09 to 1.23 (GPU<sub>K80</sub>,  $\nu = 30$ ), and 1.16 to 1.22 (GPU<sub>V100</sub>,  $\nu = 30$ ), as compared to BEHZ. The fluctuations in the runtimes for the multi-threaded CPU experiments can be attributed to specifics of the OpenMP setup in our benchmarking environment.

Tables VIII and IX suggest that the use of  $\nu = 60$  bits provides roughly 0.9x to 2x (resp. 1.5x to 3.9x) improvements for decryption (resp. homomorphic multiplication) on 64-bit machines. This is logical since these machines perform 64-bit operations natively. It can also be noticed that V100 outperforms K80 by 2x to 5x for both decryption and homomorphic multiplication. This can be attributed to a combined effect of multiple factors, such as compute capability (3.7 for K80) vs. (7.0 for V100), number of cores and clock rate (2,496 cores at 824MHz for K80 vs 5,120 cores at 1,380MHz for V100) and memory technology and bandwidth (GDDR5 240 GB/sec for K80 vs. HBM2 900GB/sec for V100). It should be noted that K80 and V100 were released by NVIDIA in November 2014 and June 2017, respectively [35].

To examine the improvements that can be achieved from parallel implementations on different platforms, Figure 3 shows the speedup factors computed as the ratio of CPU<sub>ST</sub> and either CPU<sub>MT</sub> or GPU. We use the best performance between  $\nu = 30$  bits and  $\nu = 60$  bits for multi-threaded speedups and the best between K80 and V100 for GPU speedups. We observe that GPU can improve the performance by one to two orders of magnitudes whereas a multi-threaded CPU implementation can hardly achieve one order of magnitude. It should be remarked that the multi-threaded CPU implementation achieved a sub-linear speedup because (1) the loop parallelization for NTTs was done at the CRT level by executing each CRT channel in a separate CPU thread, i.e., no intra-NTT parallelization was exploited in PALISADE, and 2) the workloads were not equally distributed between the CPU sockets.

We remark that our best GPU results, namely the homomorphic multiplication runtime of 51 ms for  $n = 2^{16}$  and  $\log_2 q = 1,770$  and 18.7 ms for  $n = 2^{16}$  and  $\log_2 q = 1,020$ , **are more than two orders of magnitude faster than best previously reported runtimes for other implementations of the BFV scheme**. For instance, the FPGA-based implementation HEPCloud in [18] of the textbook BFV scheme computed a homomorphic multiplication for  $n = 2^{15}$  and  $\log_2 q = 1,228$  in 26.67 seconds (with 3.36 seconds spent on the computation and the rest on the off-chip memory access). The BEHZ variant NFLlib CPU implementation in [9] ran a homomorphic multiplication for  $n = 2^{15}$  and  $\log_2 q = 1,590$  in 4.9 seconds.

### E. Practical Noise Growth

We showed previously that BEHZ and HPS can theoretically achieve the same multiplicative depth (see Table V) as the textbook BFV scheme. The noise analysis provided for both variants was conservative (worst-case) and used  $\delta = n$  to estimate the noise growth. However, in practice it was shown [10] that a lower value of  $\delta$  can be used, specifically,  $\delta = 2\sqrt{n}$ . The authors used the Central Limit Theorem (CLT) to derive a practical heuristic estimate of noise growth. We ran an experiment to verify this analysis and measure the maximum multiplicative depth each variant can achieve.

We wrote a simple procedure that encrypts a plaintext message  $\mu$  and iteratively multiplies it with an encryption of 1. In each iteration, we decrypt the product and check if it is equal to  $\mu$  counting the number of sequential multiplications. Table X shows the lowest value of the maximum multiplicative depth that can be reached without decryption failures by each variant in this experiment vs. the maximum depth estimated using the heuristic CLT-based technique of [10] for the same parameters. It suggests that HPS can achieve a higher multiplicative depth than BEHZ for almost all parameters. The HPS multiplicative depth conforms to the heuristic noise analysis having  $\delta = O(\sqrt{n})$ , where the constant is always less than two (which corresponds to the column “Max. Est.  $L_o$ ” in Table X), whereas the BEHZ depth requires a higher value for  $\delta$ . We found by fitting the experimental results to

Table X: Maximum multiplicative depth  $L_\circ$  experimentally observed for each RNS variant ( $\lambda_\circ \geq 128$ ) over 1024 runs vs. maximum heuristic estimate for  $\delta = 2\sqrt{n}$  for the same lattice parameters.

$\log_2^n$	$\lceil \log_2^q \rceil$	$t$	Max. Est. $L_\circ$	Max. Exp. $L_\circ$	
				BEHZ	HPS
12	60	2	1	2	2
		65537	0	0	1
13	120	2	5	5	6
		65537	2	2	3
14	360	2	18	16	21
		65537	10	9	10
15	600	2	31	26	35
		65537	16	15	19
16	1020	2	51	43	56
		65537	28	25	30
16	1770	2	90	75	98
		65537	50	44	52

the noise expression that  $\delta \approx O(n^{0.7})$  provides an adequate estimate of practical noise growth in BEHZ. We ran each experiment  $2^{10}$  times. Hence the practical maximum depth numbers presented in Table X correspond to an estimated decryption failure probability of  $2^{-10}$ . For the heuristic correctness constraint of HPS, the authors used a much lower probability estimate [10] (close to  $2^{-40}$  for all products of random polynomials), which explains why the HPS correctness constraint gives more conservative depth estimates. We remark that we performed the same comparison of practical noise growth for both schemes as in Table X for 60-bit moduli, and we found the results to be very similar to the 30-bit case.

Our interpretation of this behavior is that the RNS techniques used in the BEHZ variant transform the ciphertexts in such a way that CLT is no longer valid, i.e., we no longer deal with sums of zero-centered independent random variables in certain polynomial multiplications. We claim that the deviation from uncorrelated zero-centered random distribution is introduced by the step described in (Lemma 4 in [9]), i.e., by the small Montgomery representation/reduction shown in Algorithm 2. In the textbook BFV and HPS variants, ciphertexts are uncorrelated and zero-centered (with respect to the interval  $[-q/2] \leq x < [q/2]$ ). When we introduce the overflow ( $q \cdot v$ ) term, where  $\|v\|_\infty \leq k$ , in fast base conversion, we get ciphertexts that are no longer zero-centered but are biased (correlated). Therefore, the noise growth is faster than what CLT predicts (except for the case of  $k = 2$  at  $t = 2$  in Table X, when the overflow term contribution is the smallest). In contrast to the HPS case, we cannot use heuristic (average-case) estimates for BEHZ and can guarantee the correctness only with worst-case estimates, i.e.,  $\delta = n$ .

To verify this interpretation, we experimentally examined the distribution of  $v$  for different RNS bases. The variable  $v$  can be best approximated by a generalized Pareto (power-law) distribution that is not only different from the uniformly random distribution of coefficients in the ciphertext polynomials, but also depends on these coefficients. This violates the conditions for the independence between random variables required for CLT analysis, and hence the square-root average-case noise growth cannot be achieved in this case.

This implies there is a major practical difference between BEHZ and HPS, which is far more significant than the incremental performance improvements we observed when comparing the runtimes for same lattice parameters. Although the effect of  $\rho$  is small for large  $\tilde{m}$  (for the worst-case analysis), the deviation from zero-centered random distribution has a more profound effect in practice. We note that this behavior has been observed in both our CPU and GPU implementations of BEHZ. We also observed a similar noise growth behavior in the SEAL implementation of BEHZ (SEAL version 2.3.0-4).

## VII. CONCLUSION

Our work presents the implementation and performance evaluation of two RNS variants (BEHZ [9] and HPS [10]) of the BFV SHE scheme. We have analyzed the performance of both variants theoretically and experimentally using several flavors of implementations (CPU single- and multi-threaded, and GPU). Our analysis shows that HPS outperforms BEHZ in almost all settings on different platforms (for same values of lattice parameters). However,

HPS decryption is outperformed by BEHZ when the moduli size is 60 bits, which is due to the multi-precision (double double) floating-point operations in HPS.

Our experiments show that our multi-threaded CPU implementation using OpenMP can hardly attain a one-order-of-magnitude improvement over the single-threaded setting, whereas GPUs can achieve up to two orders of magnitude.

We have also demonstrated that the practical noise growth in BEHZ is much faster than that of HPS, which can significantly increase the runtime and storage requirements, and also limit the maximum multiplicative depth supported by BEHZ. We provide a possible explanation for this behavior; however, a more careful analysis of BEHZ noise growth would be needed to characterize this behavior formally.

## VIII. ACKNOWLEDGMENT

This work was supported by A \* STAR Institute for Infocomm Research (I2R) and the National University of Singapore. Polyakov and Rohloff's work was supported in part by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA). The views expressed are those of the authors and do not necessarily reflect the official policy or position of the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein.

## REFERENCES

- [1] Craig Gentry. *A Fully Homomorphic Encryption Scheme*. PhD thesis, Stanford University, 2009. [crypto.stanford.edu/craig](http://crypto.stanford.edu/craig).
- [2] Marten Van Dijk, Craig Gentry, Shai Halevi, and Vinod Vaikuntanathan. Fully Homomorphic Encryption over the Integers. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pages 24–43. Springer, 2010.
- [3] Zvika Brakerski. Fully Homomorphic Encryption without Modulus Switching from Classical GapSVP. In *Advances in cryptology—crypto 2012*, pages 868–886. Springer, 2012.
- [4] Adriana López-Alt, Eran Tromer, and Vinod Vaikuntanathan. On-the-Fly Multiparty Computation on the Cloud via Multikey Fully Homomorphic Encryption. In *Proceedings of the forty-fourth annual ACM symposium on Theory of computing*, pages 1219–1234. ACM, 2012.
- [5] Junfeng Fan and Frederik Vercauteren. Somewhat Practical Fully Homomorphic Encryption. *IACR Cryptology ePrint Archive*, 2012:144, 2012.
- [6] Craig Gentry, Amit Sahai, and Brent Waters. Homomorphic encryption from Learning with Errors: Conceptually-Simpler, Asymptotically-Faster, Attribute-Based. In *Advances in Cryptology—CRYPTO 2013*, pages 75–92. Springer, 2013.
- [7] Zvika Brakerski, Craig Gentry, and Vinod Vaikuntanathan. (leveled) Fully Homomorphic Encryption without Bootstrapping. *ACM Transactions on Computation Theory (TOCT)*, 6(3):13, 2014.
- [8] Tancrede Lepoint and Michael Naehrig. A Comparison of the Homomorphic Encryption Schemes FV and YASHE. In *International Conference on Cryptology in Africa*, pages 318–335. Springer, 2014.
- [9] Jean-Claude Bajard, Julien Eynard, M Anwar Hasan, and Vincent Zucca. A Full RNS Variant of FV Like Somewhat Homomorphic Encryption Schemes. In *International Conference on Selected Areas in Cryptography*, pages 423–442. Springer, 2016.
- [10] Shai Halevi, Yuriy Polyakov, and Victor Shoup. An Improved RNS Variant of the BFV Homomorphic Encryption Scheme. <https://eprint.iacr.org/2018/117>.
- [11] Hao Chen Kim Laine and Rachel Player. Simple Encrypted Arithmetic Library-SEAL (v2.1). Technical report, Technical report, Microsoft Research, 2016.
- [12] Yuriy Polyakov, Kurt Rohloff, and Gerard W. Ryan. PALISADE Lattice Cryptography Library. <https://git.njit.edu/palisade/PALISADE>, Accessed June 2018.
- [13] Yuriy Polyakov, Kurt Rohloff, Gyana Sahu, and Vinod Vaikuntanathan. Fast Proxy Re-Encryption for Publish/Subscribe Systems. *ACM Trans. Priv. Secur.*, 20(4):14:1–14:31, 2017.
- [14] Kamil Doruk Gür, Yuriy Polyakov, Kurt Rohloff, Gerard W. Ryan, Hadi Sajjadpour, and Erkey Savaş. Practical Applications of Improved Gaussian Sampling for Trapdoor Lattices. *Cryptology ePrint Archive*, Report 2017/1254, 2017. <https://eprint.iacr.org/2017/1254>.
- [15] W. Dai, Y. Doröz, Y. Polyakov, K. Rohloff, H. Sajjadpour, E. Savaş, and B. Sunar. Implementation and Evaluation of a Lattice-Based Key-Policy ABE Scheme. *IEEE Transactions on Information Forensics and Security*, 13(5):1169–1184, May 2018.
- [16] D. B. Cousins, G. Di Crescenzo, K. D. Gür, K. King, Y. Polyakov, K. Rohloff, G. W. Ryan, and E. Savaş. Implementing Conjunction Obfuscation under Entropic Ring Lwe. In *2018 IEEE Symposium on Security and Privacy (SP)*, volume 00, pages 68–85.
- [17] Ahmad Al Badawi, Bharadwaj Veeravalli, and Mi Mi Khin Aung. High-Performance FV Somewhat Homomorphic Encryption on GPUs: An Implementation using CUDA. *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2018(2), 2018.
- [18] Sujoy Sinha Roy, Kimmo Jarvinen, Ingrid Verbauwhede, Frederik Vercauteren, and Jo Vliegen. HEPCloud: An FPGA-based Multicore Processor for FV Somewhat Homomorphic Function Evaluation. *IEEE Transactions on Computers*, 2017.
- [19] Nathan Dowlin, Ran Gilad-Bachrach, Kim Laine, Kristin Lauter, Michael Naehrig, and John Wernsing. Cryptonets: Applying Neural Networks to Encrypted Data with High Throughput and Accuracy. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, ICML'16*, pages 201–210. JMLR.org, 2016.
- [20] Miran Kim, Yongsoo Song, Shuang Wang, Yuhou Xia, and Xiaoqian Jiang. Secure Logistic Regression Based on Homomorphic Encryption: Design and Evaluation. *JMIR Med Inform*, 6(2):e19, Apr 2018.

- [21] David Harvey. Faster Arithmetic for Number-Theoretic Transforms. *Journal of Symbolic Computation*, 60:113–119, 2014.
- [22] Karl C Posch and Reinhard Posch. Modulo Reduction in Residue Number Systems. *IEEE Transactions on Parallel and Distributed Systems*, 6(5):449–454, 1995.
- [23] Shinichi Kawamura, Masanobu Koike, Fumihiko Sano, and Atsushi Shimbo. Cox-Rower Architecture for Fast Parallel Montgomery Multiplication. In *International Conference on the Theory and Applications of Cryptographic Techniques*, pages 523–538. Springer, 2000.
- [24] Victor Shoup et al. NTL: A Library for Doing Number Theory, 2001.
- [25] CUDA Nvidia. Toolkit Documentation. *NVIDIA CUDA Getting Started Guide for Linux*, 2014.
- [26] AP Shenoy and Ramdas Kumaresan. Fast Base Extension using a Redundant Modulus in RNS. *IEEE Transactions on Computers*, 38(2):292–297, 1989.
- [27] Joppe W. Bos, Kristin Lauter, Jake Loftus, and Michael Naehrig. Improved Security for a Ring-Based Fully Homomorphic Encryption Scheme. In *IMACC 2013*, pages 45–64, 2013.
- [28] Carlos Aguilar-Melchor, Joris Barrier, Serge Guelton, Adrien Guinet, Marc-Olivier Killijian, and Tancrede Lepoint. NFLlib: NTT-based Fast Lattice Library. In *Cryptographers? Track at the RSA Conference*, pages 341–356. Springer, 2016.
- [29] How to Tell If GPU Cores are Actually 32/64-bit Processors. Available at: <https://devtalk.nvidia.com>. Accessed Nov 2 2018.
- [30] Richard E Crandall. Integer Convolution via Split-Radix Fast Galois Transform. *Center for Advanced Computation Reed College*, 1999.
- [31] Ahmad Al Badawi, Veeravalli Bharadwaj, and Khin Mi Mi Aung. Efficient Polynomial Multiplication via Modified Discrete Galois Transform and Negacyclic Convolution. In *Future of Information and Communications Conference (FICC)*. IEEE, 2018.
- [32] Achieved Occupancy. Available at: <https://docs.nvidia.com/gameworks/content/developertools/desktop/analysis/report/cudaexperiments/kernellevel/achievedoccupancy.htm>. Accessed Nov 3 2018.
- [33] Martin Albrecht, Samuel Scott, and Rachel Player. On the Concrete Hardness of Learning with Errors. *Journal of Mathematical Cryptology*, 9(3):169–203, 10 2015.
- [34] Melissa Chase, Hao Chen, Jintai Ding, Shafi Goldwasser, Sergey Gorbunov, Jeffrey Hoffstein, Kristin Lauter, Satya Lokam, Dustin Moody, Travis Morrison, Amit Sahai, and Vinod Vaikuntanathan. Security of Homomorphic Encryption. Technical report, HomomorphicEncryption.org, Redmond WA, July 2017.
- [35] List of Nvidia Graphics Processing Units. Available at: [https://en.wikipedia.org/wiki/List\\_of\\_Nvidia\\_graphics\\_processing\\_units](https://en.wikipedia.org/wiki/List_of_Nvidia_graphics_processing_units). Accessed May 26th 2018.