




Revisiting Homomorphic Encryption Schemes for Finite Fields

Andrey Kim^{1,2}, Yuriy Polyakov^{1,3}(✉) , and Vincent Zucca^{4,5,6}

¹ New Jersey Institute of Technology, Newark, USA

² Samsung Advanced Institute of Technology, Suwon, Republic of Korea
andrey.kim@samsung.com

³ Duality Technologies, Newark, USA
ypolyakov@dualitytech.com

⁴ DALI, Université de Perpignan Via Domitia, Perpignan, France
vincent.zucca@univ-perp.fr

⁵ LIRMM, Univ Montpellier, Montpellier, France

⁶ imec-COSIC, KU Leuven, Leuven, Belgium

Abstract. The Brakerski-Gentry-Vaikuntanathan (BGV) and Brakerski/ Fan-Vercauteren (BFV) schemes are the two main homomorphic encryption (HE) schemes to perform exact computations over finite fields and integers. Although the schemes work with the same plaintext space, there are significant differences in their noise management, algorithms for the core homomorphic multiplication operation, message encoding, and practical usability. The main goal of our work is to revisit both schemes, focusing on closing the gap between the schemes by improving their noise growth, computational complexity of the core algorithms, and usability. The other goal of our work is to provide both theoretical and experimental performance comparison of BGV and BFV.

More precisely, we propose an improved variant of BFV where the encryption operation is modified to significantly reduce the noise growth, which makes the BFV noise growth somewhat better than for BGV (in contrast to prior results showing that BGV has smaller noise growth for larger plaintext moduli). We also modify the homomorphic multiplication procedure, which is the main bottleneck in BFV, to reduce its algorithmic complexity. Our work introduces several other novel optimizations, including lazy scaling in BFV homomorphic multiplication and an improved BFV decryption procedure in the Residue Number System (RNS) representation. We also develop a usable variant of BGV as a more efficient alternative to BFV for common practical scenarios.

We implement our improved variants of BFV and BGV in PALISADE and evaluate their experimental performance for several benchmark computations. The experimental results suggest that our BGV implementation is faster for intermediate and large plaintext moduli, which are often used in practical scenarios with ciphertext packing, while our BFV implementation is faster for small plaintext moduli.

The full version of the paper is available at <https://eprint.iacr.org/2021/204>.

© International Association for Cryptologic Research 2021
M. Tibouchi and H. Wang (Eds.): ASIACRYPT 2021, LNCS 13092, pp. 608–639, 2021.
https://doi.org/10.1007/978-3-030-92078-4_21

1 Introduction

Homomorphic encryption (HE) is a powerful cryptographic primitive that enables performing computations over encrypted data without having access to the secret key. The HE research area has seen a lot of progress since the formulation of the first fully homomorphic encryption construction by Gentry in 2009 [17], and the schemes implemented in modern HE libraries are multiple orders of magnitude faster than the initial implementation of Gentry’s scheme [18]. The most common HE schemes are typically grouped into three classes based on the data types they support computations on. The first class primarily works with Boolean circuits and decision diagrams, similar to the original Gentry scheme, and includes the FHEW and TFHE schemes [12, 15]. The second class supports modular arithmetic over finite fields, which typically correspond to vectors of integers mod t , where t is a prime power commonly called as the plaintext modulus. The second class is also sometimes used for small-integer arithmetic. This class includes Brakerski-Gentry-Vaikuntantan (BGV) and Brakerski/Fan-Vercauteren (BFV) schemes [8, 9, 16]. The third, and most recent, class supports approximate computations over vectors of real and complex numbers, and is represented by the Cheon-Kim-Kim-Song (CKKS) scheme [11]. All these schemes are based on the hardness of the Ring Learning With Errors (RLWE) problem, where noise is added during encryption and key generation to achieve the hardness properties. The noise grows as encrypted computations are performed, and the main functional parameter in all these schemes, the ciphertext modulus Q , needs to be large enough to accommodate the noise growth, or a special bootstrapping procedure may be used to reset the noise and keep the value of Q relatively small.

Our work focuses on the HE schemes of the second class. Although the BGV and BFV schemes work with the same plaintext algebra, they use different strategies for encoding the message composed of integers in \mathbb{Z}_t and controlling the noise. The BGV scheme encodes the message in the least significant digit (LSD) of integers in \mathbb{Z}_Q and applies the modulus switching technique to keep the noise magnitude constant, i.e., it scales down Q by a factor that corresponds to the noise added after the previous modulus switching call. The BFV scheme encodes the message in the most significant digit (MSD) of integers in \mathbb{Z}_Q and uses a special form of homomorphic multiplication, where ciphertext polynomials are multiplied without modular reduction and then scaled down by Q/t . In BFV, the value of Q is typically constant and the noise magnitude increases at a rate similar to how Q decreases in BGV. The difference in noise management strategies between BGV and BFV affects the noise growth and efficiency of the schemes. Costache and Smart performed a noise growth comparison, which suggested that BGV has better noise growth for larger t than BFV [13]. However, the authors did not examine the computational complexity difference, and it has not been clear up to this moment how the schemes compare in terms of practical performance, both from the perspective of computational complexity and actual experimental measurements.

The main goal of this paper is to present improved variants of BFV and BGV schemes, which also close the gap between the schemes. The other goal is to compare the theoretical complexity of their primitive operations, and exper-

imental performance of BGV and BFV for several different scenarios using our software implementation in the PALISADE library [2].

Modified BFV Scheme. We propose two modifications for the BFV scheme. The first modification deals with encryption, and the second modification revises the homomorphic multiplication operation. The net effects of these modifications are smaller noise growth and faster homomorphic multiplication in BFV.

The encryption in BFV can be represented as $\mathbf{a} \cdot \mathbf{s} + \mathbf{e} + \Delta \mathbf{m}$ (for simplicity, we focus here on the secret-key formulation), where \mathbf{a} is a uniformly random ring element in cyclotomic ring \mathcal{R}_Q , \mathbf{s} and \mathbf{e} are the secret key and Gaussian noise ring elements in \mathcal{R} , \mathbf{m} is a message in \mathcal{R}_t , and $\Delta = \lfloor Q/t \rfloor$ is the scaling factor. Our analysis shows that the difference between Δ and Q/t , which is often described in terms of $r_t(Q) := Q - t\Delta$, brings about a significant error (proportional to $r_t(Q)$) that affects the first homomorphic multiplication and increases the noise growth in BFV as compared to BGV for larger t . If this error is removed, i.e., $r_t(Q) \approx 0$, the noise growth in BFV becomes the same, or actually somewhat better, as in BGV. In view of this, our first modification suggested for BFV is to replace the encryption operation in BFV with $\mathbf{a} \cdot \mathbf{s} + \mathbf{e} + \left\lfloor \frac{Q}{t} \mathbf{m} \right\rfloor$, which is a more natural choice as compared to the one in the original BFV. This encryption function also significantly simplifies the noise analysis and estimates for BFV homomorphic multiplication. Note that this modification can be likewise applied to the original Brakerski LWE scheme [8].

The most expensive operation in BFV is homomorphic multiplication as it requires a multiplication of two ciphertexts \mathbf{c}_1 and \mathbf{c}_2 without modular reduction, followed by scaling the results of the tensor product by t/Q . Algorithmically, this requires extending both ciphertexts to modulus QP , where P is sufficiently larger than Q , performing a tensor product which involves expensive Number Theoretic Transforms (NTTs), scaling down the result by Q/t , and finally switching the scaling result from P to Q . We propose a more efficient procedure for homomorphic multiplication where the values of $P \approx Q$, which saves some expensive modulus extension operations and NTTs. The main idea is to apply modulus switching to one of the ciphertexts, e.g., \mathbf{c}_2 , to change it from Q to P (denote it as \mathbf{c}'_2), and then do scaling by t/P after the tensor product. This removes the requirement for extending the scaling result from P to Q (it will already be in Q) at the expense of doing a smaller number of modulus extensions during the modulus switching of \mathbf{c}_2 . The other benefit is that the tensor product of \mathbf{c}_1 and \mathbf{c}'_2 can be scaled by t/P directly in PQ , i.e., we have a tensor product mod PQ instead of a tensor product without modular reduction.

We also introduce a leveled version of BFV homomorphic multiplication, where ciphertexts modulo a larger modulus Q are internally scaled down to a smaller modulus Q_ℓ (or P_ℓ), the standard homomorphic multiplication operations are performed, and then the results are scaled back up to Q . The benefit of this approach is that the ciphertexts still look the same (modulo Q) outside the homomorphic multiplication operations, but we get BGV-like benefits of working with smaller moduli in multiplication. The combined effect of our improvements in homomorphic multiplication is the speed-up of up to 4x, as compared to a

prior state-of-the-art BFV implementation, when dealing with multiplications at deeper levels of computation.

BFV Scheme Optimizations. We also introduce several algorithmic optimizations that equally apply to the classical BFV and our modified variant. The first optimization is for the scenarios where we need to add multiple BFV ciphertexts that were just obtained by BFV multiplication. The standard way is to perform many expensive BFV multiplications and then add up the result. However, we can delay the scaling by t/Q (or by t/P in our BFV variant) in each homomorphic multiplication until the sum is computed, and then just do one scaling at the end. This saves many expensive NTTs and modulus extension operations. We denote this optimization as *lazy scaling*. The lazy scaling can be combined with previously known lazy relinearization to push most of the expensive computations in a homomorphic multiplication, i.e., scaling and relinearization, to the end, after the aggregation is done.

Some of the other optimizations apply to Residue Number Systems (RNS) variants of BFV, where multi-precision integers in \mathbb{Z}_Q are split into vectors of smaller integers using the Chinese Remainder Theorem (CRT) to perform operations efficiently using native (64-bit) integer types. The RNS variants are now predominately used in practice, and are implemented in the SEAL [29], PALISADE [2], and Lattigo [1] software libraries. There are two main RNS variants of BFV: the Bajard-Eynard-Hasan-Zucca (BEHZ) variant based on modular integer arithmetic and Montgomery reductions and the Halevi-Polyakov-Shoup (HPS) variant based on a combination of modular integer arithmetic and floating-point approximations [5, 21].

A significant limitation of the HPS approach is that high-precision (“long double” or even quad-precision) floating-point arithmetic is required to support larger CRT moduli: long doubles are needed for CRT moduli from 47 to 58 bits, and quad-precision floats are needed for higher CRT moduli [21]. We introduce a general-purpose digit decomposition technique (inspired by digit decomposition in key switching) and apply it to the HPS decryption procedure to add support for arbitrary CRT moduli using only regular double-precision floating-point arithmetic, thus overcoming this limitation of the HPS variant. This digit decomposition technique can be applied to other mixed integer/floating-point RNS operations to reduce precision requirements for floating-point arithmetic.

We also apply the full RNS variant of hybrid key switching [24] recently proposed for the CKKS scheme to both BFV RNS variants, and demonstrate how some auxiliary CRT moduli needed for homomorphic multiplication can be reused for hybrid key switching. This key switching method has some benefits (smaller noise growth, better efficiency for deeper computations) over the residue decomposition key switching method previously used in both RNS variants of BFV.

BGV Scheme Optimizations and Usability Improvements. We use the Gentry-Halevi-Smart (GHS) variant of BGV as the basis for our BGV instantiation [20, 23]. Although the original GHS variant performs some operations in RNS, it still

uses multiprecision integer arithmetic for key switching and some scenarios of modulus switching. For instance, although the GHS paper originally introduced the hybrid key switching technique, the authors used multiprecision arithmetic for the digit decomposition step. We apply the full RNS version of hybrid key switching to our BGV instantiation and eliminate any multiprecision arithmetic from our BGV implementation, thus significantly improving its efficiency.

One of the challenges in the GHS variant is the need to perform dynamic noise estimation, which makes the BGV implementation less robust and usable as compared to the BFV variants where noise estimation is typically needed only at the parameter generation phase. We develop a more usable and robust variant of BGV that is essentially as simple to use as the current BFV implementations. This variant only needs to know the multiplicative depth and maximum number of additions per level for many common scenarios. The main advantage of this BGV variant is that it is significantly faster than our BFV implementations for certain practical scenarios, yet its usability matches that of BFV.

Implementation and Performance Comparison. We implement the improved variants of BFV and BGV in PALISADE, and provide their comparison for specific benchmark computations. To the best of our knowledge, this is the first publicly available implementation of both schemes in the same software library. We also perform theoretical comparison of the computational complexity for the operations that differ between BFV and BGV.

The comparison results can be summarized as follows:

- Our improved variant of BFV has somewhat better noise growth than BGV, in contrast to prior results for the original BFV scheme that showed better noise growth for BGV at larger plaintext moduli [13].
- Our best variant of BFV is faster than BGV for small plaintext moduli, while BGV is faster for intermediate and large plaintext moduli used in many practical scenarios.
- The speed-up in homomorphic multiplication of our best BFV variant compared to a prior state-of-the-art RNS implementation of BFV goes up to 4x for deeper computations.

Related work. Costache et al. further examine the difference between the noise growth in BFV and BGV [14] to improve the analysis presented in [13]. They explore an alternative heuristic noise analysis approach to obtain tighter noise bounds. We point out that this new analysis has some inaccuracies, e.g., the effect of extra noise due to $r_t(Q)$ in BFV homomorphic multiplication is not accounted for. We show that this extra noise determines the higher noise growth in BFV for large plaintext moduli, and demonstrate how this noise is removed in our BFV variant. Moreover, we show that this analysis can be carried out independently of the chosen heuristic for noise analysis. The authors also do not consider the difference in the complexity of homomorphic encryption operations between BGV and BFV, which affects their conclusions. In view of the above, we primarily compare our results with the prior work [13].

The encoding of a message in the MSD of a ciphertext as $\lceil \frac{Q}{t} \mathbf{m} \rceil$ was already used in the Key Encapsulation Mechanism (KEM) Kyber [7]. But in the case of Kyber, the plaintext modulus $t = 2$, i.e., the coefficients of the messages are either 0 or 1. Therefore, the messages can be recovered directly during decryption by checking whether the coefficients are closer to $\lceil Q/2 \rceil$ or 0, and the noise is not affected.

SEAL also independently added to v3.4.0 a modification of BFV encryption similar to what we describe in our work [29]. However, no underlying noise analysis was presented, and the prior paper related to SEAL [14] included noise analysis inaccuracies involving the rounding term $r_t(Q)$, suggesting that the full effect of this change was not well-understood.

Note that the FHEW and TFHE schemes can also support arithmetic over finite fields for small plaintext moduli (typically up to 4 bits) [28], and can be considered as an alternative to BGV and BFV for these scenarios. These schemes support fast bootstrapping (the latency is much lower than for BGV and BFV bootstrapping [22]), but their main limitation is the lack of support for CRT packing, which makes the BGV/BFV approach much more appealing when large arrays of numbers need to be computed on/bootstrapped because one ciphertext operation can perform thousands of integer operations at once.

Organization. The rest of the paper is organized as follows. In Sect. 2 we provide the necessary background on BGV and BFV. In Sects. 3 and 4, we present our improved variants of BFV and BGV, respectively. Section 5 includes the theoretical comparison of the schemes, and discussion of the experimental results. Section 6 provides the conclusions and outlines the ideas for future work.

2 Background

All logarithms are expressed in base 2 if not indicated otherwise. Let N be a power of two. We denote the $2N$ -th cyclotomic ring $\mathcal{R} = \mathbb{Z}[X]/(X^N + 1)$ and $\mathcal{R}_Q := \mathcal{R}/Q\mathcal{R}$.¹ Ring elements are indicated in bold, e.g. \mathbf{a} . For an integer $Q > 1$, we identify the ring \mathbb{Z}_Q with $(-Q/2, Q/2) \cap \mathbb{Z}$ as a representative interval and for $z \in \mathbb{Z}$, $[z]_Q \in \mathbb{Z}_Q$ denotes the centered remainder of z modulo Q , while $r_Q(z)$ denotes the classical Euclidean remainder in $[0, Q) \cap \mathbb{Z}$. For $x \in \mathbb{Q}$, $\lfloor x \rfloor$, $\lceil x \rceil$ and $\lceil x \rceil$ denote the rounding to the lower, closest and higher integer, respectively. We extend these notations to elements of \mathcal{R} by applying them coordinate-wise. For $\mathbf{a} = a_0 + a_1 \cdot X + \dots + a_{N-1} \cdot X^{N-1} \in \mathcal{R}$, we denote the ℓ_∞ norm of \mathbf{a} as $\|\mathbf{a}\|_\infty = \max_{0 \leq i < N} \{|a_i|\}$. There exists a constant $\delta_{\mathcal{R}}$ such that $\|\mathbf{a} \cdot \mathbf{b}\|_\infty \leq \delta_{\mathcal{R}} \|\mathbf{a}\|_\infty \|\mathbf{b}\|_\infty$ for any $(a, b) \in \mathcal{R}^2$. It is well-known that for $\mathcal{R} = \mathbb{Z}[X]/(X^N + 1)$, $\delta_{\mathcal{R}} = N$. However in practice this bound is only reached with exponentially low probability. As shown in [21], the bound $\delta_{\mathcal{R}} = 2\sqrt{N}$ is much closer to what we observe experimentally, and can be used to achieve

¹ More general cyclotomic rings are also supported, and all results of our work equally apply to these non-power-of-two rings; please see [23] for more details on general cyclotomic rings.

tighter noise bounds. Another approach consists in estimating the noise size using the canonical embedding norm [20], as currently done in HELib [23]. Nonetheless, in this work we estimate the noise size using the expansion factor $\delta_{\mathcal{R}}$ with the method of [21] as it is simpler and precise enough for our purpose.

We use $\mathbf{a} \leftarrow \chi$ to denote the sampling of $\mathbf{a} \in \mathcal{R}$ according to a distribution χ . χ_{key} denotes the *uniform ternary* distribution, i.e., all the coefficients of $\mathbf{a} \leftarrow \chi_{\text{key}}$ are selected uniformly and independently from $\{-1, 0, 1\}$. This distribution is commonly used for secret key generation as it is the most efficient option conforming to the HE standard [4]. χ_{err} denotes a *discrete Gaussian* distribution with standard deviation σ_{err} , i.e. all the coefficients of $\mathbf{a} \leftarrow \chi_{\text{err}}$ are selected independently from a truncated discrete Gaussian distribution with standard deviation σ_{err} . Truncated discrete Gaussian distributions are commonly used to generate error polynomials to meet the desired hardness requirement [4]. We assume that the polynomials sampled from χ_{key} and χ_{err} have their coefficients bounded by $B_{\text{key}} = 1$ and $B_{\text{err}} = 6\sigma_{\text{err}}$, respectively. Although a Gaussian distribution is not bounded by nature, the probability for a Gaussian coefficient to be larger than $B_{\text{err}} = 6\sigma_{\text{err}}$, is less than 2^{-30} , therefore the two distributions are very close in practice. \mathcal{U}_Q denotes the *uniform distribution* over \mathcal{R}_Q , where every coefficient of \mathbf{a} is sampled uniformly and independently from \mathbb{Z}_Q .

2.1 Plaintext Space

We are interested in the BGV and BFV homomorphic encryption schemes which both share the same plaintext space \mathcal{R}_t for some integer $t > 1$. Hence, the most natural way to represent plaintext messages of these schemes is to think of them as vectors of size N with their coefficients taken modulo t . However, \mathcal{R}_t has many algebraic properties, in particular when $t = p^r$ is a prime power with p coprime to $2N$. In this case \mathcal{R}_t is actually a \mathbb{Z}_t -algebra, which means that it contains a subring isomorphic to \mathbb{Z}_t . In this paper we focus on the case $r = 1$, where $t = p$ is a prime. The interested reader can nonetheless refer to [23] for further details regarding the general case. \mathbb{Z}_t -algebra supports efficient Single-Instruction Multiple-Data (SIMD) packing/batching. For more details on the packing, the reader is referred to [30].

2.2 Homomorphic Encryption Schemes for Finite Field Arithmetic

The two schemes studied in this work: BGV and BFV are actually two instantiations of the same idea, and share, therefore, many common features. First, according to the desired security level λ and the targeted application, one starts by selecting public parameters for the considered scheme: ring dimension $N = 2^d$, the plaintext modulus t , a ciphertext modulus Q and two probability distributions χ_{key} and χ_{err} on the ring \mathcal{R} . In both cases, the secret key will be an element $\mathbf{s} \leftarrow \chi_{\text{key}}$. Note that BGV and BFV may be viewed as different modes of a unified scheme, where the ciphertexts may be switched from one mode/scheme to the other (see the full version for details).

Original BGV Scheme. In 2011, Brakerski et al. designed a leveled homomorphic scheme, namely capable of evaluating circuits of arbitrary size, but known beforehand [9]. The key tool of their construction is the *modulus switching* procedure which allows to switch a ciphertext \mathbf{ct} encrypted under a modulus Q to a smaller modulus Q' in order to maintain the noise level “constant”. As a consequence, one must select a chain of $L + 1$ moduli $Q_0 \mid Q_1 \mid \dots \mid Q_L = Q$ such that t and Q_L are coprime. The public key is formed as:

$$\mathbf{pk} = \left([\mathbf{a} \cdot \mathbf{s} + t\mathbf{e}]_{Q_L}, -\mathbf{a} \right) \in \mathcal{R}_{Q_L}^2,$$

which is equivalent to the Ring-LWE sample $([\mathbf{a}/t \cdot \mathbf{s} + \mathbf{e}]_{Q_L}, [-\mathbf{a}/t]_{Q_L})$ (since t and Q_L are coprime) associated to \mathbf{s} and Q_L with $\mathbf{a} \leftarrow \mathcal{U}_{Q_L}$ and $\mathbf{e} \leftarrow \chi_{\text{err}}$.

A ciphertext $\mathbf{ct} = (\mathbf{c}_0, \mathbf{c}_1) \in \mathcal{R}_Q^2$ corresponds to a degree 1 polynomial whose coefficients lie in \mathcal{R}_Q . The message $\mathbf{m} \in \mathcal{R}_t$ is hidden in the LSD of the first coefficient \mathbf{c}_0 of the ciphertext as follows:

$$\mathbf{ct} = \left([[\mathbf{m}]_t + \mathbf{u} \cdot \mathbf{pk}_0 + t\mathbf{e}_0]_{Q_L}, [\mathbf{u} \cdot \mathbf{pk}_1 + t\mathbf{e}_1]_{Q_L} \right)$$

with $\mathbf{u} \leftarrow \chi_{\text{key}}$ and $\mathbf{e}_0, \mathbf{e}_1 \leftarrow \chi_{\text{err}}$. The noise contained in a ciphertext $\mathbf{ct} = (\mathbf{c}_0, \mathbf{c}_1)$ appears explicitly once the ciphertext is evaluated on the secret key \mathbf{s} :

$$\mathbf{c}_0 + \mathbf{c}_1 \cdot \mathbf{s} = [\mathbf{m}]_t + t(\mathbf{u} \cdot \mathbf{e} + \mathbf{e}_1 \cdot \mathbf{s} + \mathbf{e}_0) = [\mathbf{m}]_t + t\mathbf{v}_{\text{fresh}} \text{ mod } Q_L, \quad (1)$$

where the term $\mathbf{v}_{\text{fresh}} = \mathbf{u} \cdot \mathbf{e} + \mathbf{e}_1 \cdot \mathbf{s} + \mathbf{e}_0$ is the noise inherent to a “freshly” encrypted ciphertext. Since $Q_0 \mid Q_1 \mid \dots \mid Q_L$, encryptions can be performed equivalently at any level i , i.e., modulo Q_i .

To decrypt a ciphertext $\mathbf{ct} = (\mathbf{c}_0, \mathbf{c}_1) \in \mathcal{R}_{Q_i}^2$ with $i \in [0, L]$, one computes $\mathbf{m}' = [\mathbf{c}_0 + \mathbf{c}_1 \cdot \mathbf{s}]_{Q_i}$ and then outputs $[\mathbf{m}']_t$. To ensure correctness of the decryption, the noise \mathbf{v} must be “small enough” such that $\mathbf{m}' = [\mathbf{m}]_t + t\mathbf{v}$ does not wrap-around modulo Q_i . As a consequence, decryption remains correct as long as:

$$\|\mathbf{v}\|_\infty < \frac{Q_0}{2t} - \frac{1}{2}.$$

One can add two ciphertexts \mathbf{ct} and \mathbf{ct}' encrypting \mathbf{m} and \mathbf{m}' , respectively, at the same level i to yield:

$$\mathbf{c}_0 + \mathbf{c}'_0 + (\mathbf{c}_1 + \mathbf{c}'_1) \cdot \mathbf{s} = [\mathbf{m} + \mathbf{m}']_t + t(\mathbf{v} + \mathbf{v}' + \mathbf{u}) \text{ mod } Q_i,$$

with $\|\mathbf{u}\|_\infty \leq 1$. This means that

$$\mathbf{ct}_{\text{add}} = ([\mathbf{c}_0 + \mathbf{c}'_0]_{Q_i}, [\mathbf{c}_1 + \mathbf{c}'_1]_{Q_i})$$

is a level- i encryption of $[\mathbf{m} + \mathbf{m}']_t$ and its noise is almost the sum of the noises of \mathbf{ct} and \mathbf{ct}' :

$$\|\mathbf{v}_{\text{add}}\|_\infty = \|\mathbf{v} + \mathbf{v}' + \mathbf{u}\|_\infty \leq \|\mathbf{v}\|_\infty + \|\mathbf{v}'\|_\infty + 1.$$

Similarly to addition, we can multiply two level- i ciphertexts \mathbf{ct} and \mathbf{ct}' to obtain the following congruence modulo Q_i :

$$(\mathbf{c}_0 + \mathbf{c}_1 \cdot \mathbf{s}) \cdot (\mathbf{c}'_0 + \mathbf{c}'_1 \cdot \mathbf{s}) = [\mathbf{m} \cdot \mathbf{m}']_t + t([\mathbf{m}]_t \cdot \mathbf{v}' + \mathbf{v} \cdot [\mathbf{m}']_t + t\mathbf{v} \cdot \mathbf{v}' + \mathbf{r}_m)$$

with $[\mathbf{m}]_t \cdot [\mathbf{m}']_t = [\mathbf{m} \cdot \mathbf{m}']_t + t\mathbf{r}_m$ and $\|\mathbf{r}_m\|_\infty \leq \delta_{\mathcal{R}}t/2$. This means that

$$\mathbf{ct}_{\text{mult}} = ([\mathbf{c}_0 \cdot \mathbf{c}'_0]_{Q_i}, [\mathbf{c}_0 \cdot \mathbf{c}'_1 + \mathbf{c}_1 \cdot \mathbf{c}'_0]_{Q_i}, [\mathbf{c}_1 \cdot \mathbf{c}'_1]_{Q_i}) \in \mathcal{R}_{Q_i}^3$$

is a degree-2 ciphertext encrypting $[\mathbf{m} \cdot \mathbf{m}']_t$ and its noise is bounded by

$$\begin{aligned} \|\mathbf{v}_{\text{mult}}\|_\infty &= \|([\mathbf{m}]_t \cdot \mathbf{v}' + \mathbf{v} \cdot [\mathbf{m}']_t + t\mathbf{v} \cdot \mathbf{v}' + \mathbf{r}_m)\|_\infty \\ &\leq \frac{\delta_{\mathcal{R}}t}{2} (2\|\mathbf{v}\|_\infty \|\mathbf{v}'\|_\infty + \|\mathbf{v}\|_\infty + \|\mathbf{v}'\|_\infty + 1). \end{aligned}$$

Remark 1. The reader can notice that the degree, and thus the size, of a ciphertext increases after each multiplication, increasing therefore the future communication and computational costs. Since this is something one wants to avoid in practice, the degree-2 ciphertexts are “relinearized” after a homomorphic multiplication to degree-1 ciphertexts using a key-switching procedure (see the full version).

The main issue with homomorphic multiplication is its quadratic noise growth, which implies that by choosing $Q_L \approx \|\mathbf{v}_{\text{fresh}}\|_\infty^L$ one could only perform $\log_2 L$ consecutive multiplications. The idea of modulus switching is to reduce the size of the noise after each multiplication to keep it constant and prevent the quadratic blow-up. This is achieved by scaling the ciphertext \mathbf{ct} by Q_i/Q_j for $i < j$, which scales down the noise by roughly the same factor. More precisely, let $\mathbf{ct} = (\mathbf{c}_0, \mathbf{c}_1)$ be a level $j \in (0, L) \cap \mathbb{Z}$ encryption of a message \mathbf{m} with noise \mathbf{v} and let i be an integer smaller than j , then set:

$$\boldsymbol{\delta} = (t[-\mathbf{c}_0/t]_{Q_j/Q_i}, t[-\mathbf{c}_1/t]_{Q_j/Q_i}) \in \mathcal{R}^2.$$

Then one can compute

$$\mathbf{ct}' = \frac{Q_i}{Q_j} \cdot (\mathbf{c}_0 + \boldsymbol{\delta}_0, \mathbf{c}_1 + \boldsymbol{\delta}_1) \bmod Q_i.$$

Brakerski et al. showed that if $\mathbf{ct} = (\mathbf{c}_0, \mathbf{c}_1)$ is such that

$$\|[\mathbf{c}_0 + \mathbf{c}_1 \cdot \mathbf{s}]_{Q_j}\|_\infty < \frac{Q_j}{2} - \frac{tQ_j}{2Q_i} (1 + \delta_{\mathcal{R}}B_{\text{key}}),$$

then \mathbf{ct}' is an encryption of $[Q_i/Q_j \mathbf{m}]_t$ whose noise \mathbf{v}' is bounded by

$$\|\mathbf{v}'\|_\infty \leq \frac{Q_i}{Q_j} \|\mathbf{v}\|_\infty + \|\mathbf{v}_{\text{ms}}\|_\infty$$

with $\|\mathbf{v}_{\text{ms}}\|_\infty \leq (1 + \delta_{\mathcal{R}}B_{\text{key}})/2$. Therefore by choosing the encryption parameters such that performing modulus switching after a homomorphic multiplication (plus a key-switching) brings the noise back to its initial level, one can perform L consecutive multiplications instead of approximately $\log_2 L$ initially.

Gentry-Halevi-Smart (GHS) Variant of BGV. Since the modulus-switching procedure outputs an encryption of the message scaled by a factor $[Q_i/Q_j]_t$, Brakerski et al. proposed to choose the moduli $Q_i = 1 \pmod t$. This approach, although very convenient in theory, becomes challenging in practice when using a large t since it reduces significantly the range of possible moduli for the Q_i . As a consequence, Gentry, Halevi, and Smart proposed to keep track of the scaling factor for each ciphertext instead [20]. In particular, they suggested to encrypt $[Q_L \mathbf{m}]_t$ instead of $[\mathbf{m}]_t$ in Eq. (1), which provides natural downscaling to $[Q_i \mathbf{m}]_t$ as modulus switching operations are applied. However in this case, one has to pay attention when adding two ciphertexts with different scaling factors. Nonetheless this can be achieved without impacting significantly the noise by following the methodology of [25].

Gentry et al. also proposed several optimizations related to noise management. The first one is to perform modulus switching after encryption and before first multiplication, in order to reduce the noise from $\|\mathbf{v}_{\text{fresh}}\|_\infty$ to $\|\mathbf{v}_{\text{ms}}\|_\infty$. The second one is to perform modulus switching just before the next multiplication instead of just after a multiplication. This permits to reduce the noise accumulated due to other operations, such as additions or key switching, that are performed between two subsequent multiplications.

Original BFV Scheme. In [8], Brakerski proposed a *scale-invariant* construction that achieves asymptotically the same noise growth as BGV, but does not *explicitly* call the modulus-switching procedure, embedding it internally in the homomorphic multiplication. Fan and Vercauteren then ported Brakerski’s construction to the Ring-LWE setting [16], and the scheme is now commonly referred to as BFV. The BFV scheme uses a public key

$$\text{pk} = \left([\mathbf{a} \cdot \mathbf{s} + \mathbf{e}]_Q, -\mathbf{a} \right) \in \mathcal{R}_Q^2,$$

which corresponds exactly to a Ring-LWE sample associated to \mathbf{s} and Q with $\mathbf{a} \leftarrow \mathcal{U}_Q$ and $\mathbf{e} \leftarrow \chi_{\text{err}}$. The main difference between BGV and BFV is that BFV ciphertexts encode messages in their MSD instead of LSD:

$$\text{ct} = \left([\Delta[\mathbf{m}]_t + \mathbf{u} \cdot \mathbf{pk}_0 + \mathbf{e}_0]_Q, [\mathbf{u} \cdot \mathbf{pk}_1 + \mathbf{e}_1]_Q \right)$$

with $\Delta = \lfloor Q/t \rfloor$, $\mathbf{u} \leftarrow \chi_{\text{key}}$ and $\mathbf{e}_0, \mathbf{e}_1 \leftarrow \chi_{\text{err}}$. Similarly to BGV, the noise contained in a ciphertext $\text{ct} = (\mathbf{c}_0, \mathbf{c}_1)$ appears explicitly once the ciphertext is evaluated on the secret key \mathbf{s} :

$$\mathbf{c}_0 + \mathbf{c}_1 \cdot \mathbf{s} = \Delta[\mathbf{m}]_t + \mathbf{u} \cdot \mathbf{e} + \mathbf{e}_1 \cdot \mathbf{s} + \mathbf{e}_0 = \Delta[\mathbf{m}]_t + \mathbf{v}_{\text{fresh}} \pmod Q,$$

where the “fresh” noise $\mathbf{v}_{\text{fresh}} = \mathbf{u} \cdot \mathbf{e} + \mathbf{e}_1 \cdot \mathbf{s} + \mathbf{e}_0$ is the same as for BGV.

To decrypt the ciphertext ct , one needs to scale and round $\text{ct}(\mathbf{s})$ by t/Q to remove the factor Δ . Hence the decryption procedure requires computing

$$\mathbf{m}' = \left\lfloor \frac{t}{Q} [\mathbf{c}_0 + \mathbf{c}_1 \cdot \mathbf{s}]_Q \right\rfloor,$$

and the decryption will be correct as long as:

$$\|\mathbf{v}\|_\infty < \frac{Q}{2t} - \frac{r_t(Q)}{2}. \tag{2}$$

Note that the term $r_t(Q)/2$ is the error inherited from the difference between Δ^{-1} and t/Q : $\Delta t/Q = 1 - r_t(Q)/Q$. Addition of two ciphertexts \mathbf{ct} and \mathbf{ct}' is done like in BGV, but the noise growth is slightly different since the carry of the addition of the two messages is scaled by Δ :

$$\mathbf{c}_0 + \mathbf{c}'_0 + (\mathbf{c}_1 + \mathbf{c}'_1) \cdot \mathbf{s} = \Delta[\mathbf{m} + \mathbf{m}']_t + \mathbf{v} + \mathbf{v}' - r_t(Q)\mathbf{u} \bmod Q,$$

with $\|\mathbf{u}\|_\infty \leq 1$. This implies that

$$\mathbf{ct}_{\text{add}} = ([\mathbf{c}_0 + \mathbf{c}'_0]_Q, [\mathbf{c}_1 + \mathbf{c}'_1]_Q)$$

is an encryption of $[\mathbf{m} + \mathbf{m}']_t$, and its noise is bounded by

$$\|\mathbf{v}_{\text{add}}\|_\infty = \|\mathbf{v} + \mathbf{v}' + r_t(Q)\mathbf{u}\|_\infty \leq \|\mathbf{v}\|_\infty + \|\mathbf{v}'\|_\infty + r_t(Q). \tag{3}$$

The BFV multiplication of two ciphertexts \mathbf{ct} and \mathbf{ct}' is done differently, as compared to BGV, since once the product is computed, it gets scaled by Δ^2 , which has two important consequences. First, the product of \mathbf{ct} and \mathbf{ct}' cannot be reduced modulo Q , therefore it must be done in \mathcal{R} , i.e., without any modular reduction. Second, the product must be scaled down by $t/Q \approx \Delta^{-1}$ to remove the extra Δ factor and reduce the noise similarly to modulus switching in BGV. We describe the two steps of the homomorphic multiplication separately. The first part, called the *tensoring*, consists in computing the product of two ciphertexts in \mathcal{R} :

$$\mathbf{ct}_{\text{tensor}} = (\mathbf{c}_0 \cdot \mathbf{c}'_0, \mathbf{c}_0 \cdot \mathbf{c}'_1 + \mathbf{c}_1 \cdot \mathbf{c}'_0, \mathbf{c}_1 \cdot \mathbf{c}'_1) \in \mathcal{R}^3.$$

When evaluating $\mathbf{ct}_{\text{tensor}}$ on the secret key, one obtains:

$$\begin{aligned} (\mathbf{c}_0 + \mathbf{c}_1 \cdot \mathbf{s}) \cdot (\mathbf{c}'_0 + \mathbf{c}'_1 \cdot \mathbf{s}) &= (\Delta[\mathbf{m}]_t + \mathbf{v} + Q\mathbf{k}) \cdot (\Delta[\mathbf{m}']_t + \mathbf{v}' + Q\mathbf{k}') \\ &= \frac{Q}{t} \Delta [\mathbf{m} \cdot \mathbf{m}']_t + \frac{Q}{t} \mathbf{v}_{\text{tensor}} + \frac{Q^2}{t} \mathbf{k}_{\text{tensor}}, \end{aligned}$$

where

$$\begin{aligned} \mathbf{v}_{\text{tensor}} &= \frac{t\mathbf{v} \cdot \mathbf{v}'}{Q} + \frac{t\Delta}{Q} ([\mathbf{m}]_t \cdot \mathbf{v}' + [\mathbf{m}']_t \cdot \mathbf{v}) + t(\mathbf{v} \cdot \mathbf{k}' + \mathbf{v}' \cdot \mathbf{k}) \\ &\quad - r_t(Q) \left([\mathbf{m}]_t \cdot \mathbf{k}' + [\mathbf{m}']_t \cdot \mathbf{k} + \mathbf{r}_m + \frac{\Delta}{Q} [\mathbf{m}]_t \cdot [\mathbf{m}']_t \right), \\ \mathbf{k}_{\text{tensor}} &= [\mathbf{m}]_t \cdot \mathbf{k}' + [\mathbf{m}']_t \cdot \mathbf{k} + t\mathbf{k} \cdot \mathbf{k}' + \mathbf{r}_m \end{aligned}$$

with $\|\mathbf{r}_m\|_\infty \leq \delta_{\mathcal{R}}t/2$ like for BGV. Also note that $\mathbf{k} = (\mathbf{c}_0 + \mathbf{c}_1 \cdot \mathbf{s} - \Delta[\mathbf{m}]_t - \mathbf{v})/Q$ and $\mathbf{k}' = (\mathbf{c}'_0 + \mathbf{c}'_1 \cdot \mathbf{s} - \Delta[\mathbf{m}']_t - \mathbf{v}')/Q$ have their norm bounded by $(\delta_{\mathcal{R}}B_{\text{key}} + 3)/2$.

The scaling operation is done in \mathcal{R}_Q and outputs a result modulo Q :

$$\mathbf{ct}_{\text{scale}} = \left[\left[\frac{t}{Q} \mathbf{ct}_{\text{tensor}} \right] \right]_Q \in \mathcal{R}_Q^3.$$

The scaling leads to

$$\frac{t}{Q} (\mathbf{c}_{\text{tensor}_0} + \mathbf{c}_{\text{tensor}_1} \cdot s + \mathbf{c}_{\text{tensor}_2} \cdot s^2) = \Delta [m \cdot m'] + \mathbf{v}_{\text{tensor}} + Q \mathbf{k}_{\text{tensor}}.$$

The rounding of scaled terms introduces an additional error \mathbf{v}_r such that

$$\mathbf{ct}_{\text{scale}_0} + \mathbf{ct}_{\text{scale}_1} \cdot s + \mathbf{ct}_{\text{scale}_2} \cdot s^2 = \Delta [m \cdot m'] + \mathbf{v}_{\text{tensor}} + \mathbf{v}_r \pmod{Q},$$

with $\|\mathbf{v}_r\|_\infty \leq (1 + \delta_{\mathcal{R}} B_{\text{key}} + \delta_{\mathcal{R}}^2 B_{\text{key}}^2)/2$. Hence the total multiplication noise $\mathbf{v}_{\text{mult}} = \mathbf{v}_{\text{tensor}} + \mathbf{v}_r$ is bounded by

$$\begin{aligned} \|\mathbf{v}_{\text{mult}}\|_\infty \leq & \frac{\delta_{\mathcal{R}} t}{2} \left(\frac{2 \|\mathbf{v}\|_\infty \|\mathbf{v}'\|_\infty}{Q} + (4 + \delta_{\mathcal{R}} B_{\text{key}}) (\|\mathbf{v}\|_\infty + \|\mathbf{v}'\|_\infty) \right. \\ & \left. + r_t(Q) (\delta_{\mathcal{R}} B_{\text{key}} + 5) \right) + \frac{1 + \delta_{\mathcal{R}} B_{\text{key}} + \delta_{\mathcal{R}}^2 B_{\text{key}}^2}{2}. \end{aligned} \tag{4}$$

For the same reasons as for BGV, one needs to perform a key-switching operation to relinearize the resulting degree-2 ciphertext. The key-switching procedure is the same for both schemes, and we refer to the full version for further details.

2.3 RNS Representation

The Chinese Remainder Theorem (CRT) permits decomposing multi-precision integers in \mathbb{Z}_Q into vectors of smaller integers to perform operations efficiently using native (64-bit) integer data types. The integer CRT representation is also often referred to as the Residue-Number-System (RNS) representation. As a consequence, the ciphertext modulus is usually chosen as a product of “small”, i.e., fitting in a machine word, co-prime moduli so that elements of \mathcal{R}_Q are represented with their residues modulo the different q_i ’s. For BGV, we choose $Q = q_0 \cdots q_L$ and we denote $Q_i = q_0 \cdots q_i$ for $0 \leq i \leq L$, where each $q_i = 1 \pmod{2N}$, to use the efficient NTT algorithm for the multiplication of elements in \mathcal{R}_Q . For BFV, we choose $Q = q_1 \cdots q_k$, with $q_i = 1 \pmod{2N}$ for $1 \leq i \leq k$ and similarly to BGV we denote $Q_i = q_1 \cdots q_i$ for $1 \leq i \leq k$. Note that we have chosen different notations on purpose since in the original BGV the size of the moduli is directly related to the noise reduction we want to achieve by modulus switching, and is therefore dependent on the circuit one wants to evaluate. However, this constraint can be removed by considering a granular approach with dynamic noise estimation, as implemented in HELib [23] (see the discussion in Sect. 4 for more details on dynamic vs static noise estimation in BGV). On the other hand, in BFV the size of the moduli is independent of the circuit and, hence, the moduli are usually chosen as large as possible within the limit of a machine word.

When performing computations in RNS, and more particularly when implementing BGV and BFV, it is sometimes needed to switch the RNS basis, i.e., convert $\mathbf{a} \in \mathcal{R}_Q$ from its residues modulo $Q = q_1 \cdots q_k$ to $[\mathbf{a}]_Q$ modulo P for some $P = p_1 \cdots p_{k'}$. This can be achieved using a *basis extension* formulated as

$$\text{FastBaseExtension}(\mathbf{a}, Q, P) = \sum_{i=1}^k \left[\mathbf{a} \left(\frac{Q}{q_i} \right)^{-1} \right]_{q_i} \frac{Q}{q_i} \bmod p_j. \tag{5}$$

Note that the basis extension does not yield $[\mathbf{a}]_Q \bmod P$ but rather $[\mathbf{a}]_Q + \mathbf{u}Q \bmod P$ with $\|\mathbf{u}\|_\infty < k/2$. When the result of this extension is divided by Q , as in many procedures of BGV and BFV, the error caused by this Q -overflow \mathbf{u} can be neglected most of the times. However in certain cases, as in the BFV decryption procedure, this overflow cannot be tolerated and needs to be removed/corrected. This can be achieved either using integer instructions with the so-called γ -correction technique of [5], or using floating-point instructions to retrieve \mathbf{u} as in [21] since

$$\mathbf{u} = \left[\sum_{i=1}^k \left[\mathbf{a} \left(\frac{Q}{q_i} \right)^{-1} \right]_{q_i} \frac{1}{q_i} \right]. \tag{6}$$

The same problem occurs during BFV homomorphic multiplication, and if it is not handled using either of the techniques of [5] or [21], the impact of \mathbf{u} on the noise growth will be significant [6].

2.4 Hybrid Key Switching in RNS

Key switching transforms a ciphertext $\mathbf{ct} = (\mathbf{c}_0, \mathbf{c}_1) \in \mathcal{R}_Q^2$, which can be decrypted with \mathbf{s}_A , into a ciphertext $\mathbf{ct}' = (\mathbf{c}'_0, \mathbf{c}'_1) \in \mathcal{R}_Q^2$ encrypting the same message as \mathbf{ct} , but decryptable with another secret key \mathbf{s}_B . This procedure is needed to compute automorphisms (rotations) of the ciphertexts [19], or to relinearize ciphertexts after a homomorphic multiplication. Note that this procedure adds a noise \mathbf{v}_{ks} to the ciphertexts.

Several ways of performing the key-switching procedure have been found over the years. The first one was formulated by Brakerski and Vaikuntanathan (BV) in [10] and extended to RNS in [5]. This technique is based on digit decomposition of one ring element in the ciphertext. Unfortunately the BV key switching requires a quadratic number of NTTs to be computed, and hence becomes the main bottleneck of the scheme (asymptotically, and often in practice), and causes a relatively large noise growth. In [20], Gentry, Halevi, and Smart proposed another alternative for key switching. Their method, which we refer to as the GHS key switching, has a smaller noise growth than BV, and is also more efficient (asymptotically, and in many practical cases) since it only requires a linear number of NTTs. The drawback of this method is that one either needs to double the dimension N or reduce the size of Q by a factor of 2 for security reasons. Gentry, Halevi, and Smart also presented a hybrid key switching technique,

which combines BV digit decomposition and larger modulus from GHS to provide the best tradeoff between the two techniques. The RNS versions of hybrid key switching were later derived for the CKKS scheme in [26] (for one small special prime) and in [24] for the more general case. The hybrid key switching technique [24] is the most efficient one in practice, both in terms of performance and noise growth, as our detailed comparison of the BV, GHS, and Hybrid key switching in the full version shows. Hence we use the hybrid key switching in our implementation.

3 Improved BFV Scheme

One can notice from Eqs. (2), (3) and (4) that the noise of BFV is impacted by the $r_t(Q)$ factor which does not appear in BGV. This factor causes faster noise growth for BFV when using larger plaintext moduli, as compared to BGV. In this section we show that this problem is not inherent to the MSD encoding of BFV, but rather comes from the choice for its instantiation in [16] and prior LWE-based Brakerski scheme [8]. We show that by instantiating the scheme in a more natural way, we can get rid of this $r_t(Q)$ term. We also present a modified homomorphic multiplication procedure that significantly improves the complexity of BFV homomorphic multiplication, as compared to all prior variants of BFV.

In this section, $\mathbf{ct} = (\mathbf{c}_0, \mathbf{c}_1)$ and $\mathbf{ct}' = (\mathbf{c}'_0, \mathbf{c}'_1)$ denote two BFV ciphertexts encrypting, respectively, the messages $[\mathbf{m}]_t$ and $[\mathbf{m}']_t$ with noise \mathbf{v} and \mathbf{v}' .

3.1 Noise Reduction

To fully understand the problem with faster noise growth in BFV for larger plaintext moduli, let us examine more carefully the noise bound after a multiplication in BFV (Eq. (4)). This bound can be simplified by analyzing only the dominant terms, which determine the noise magnitude. More precisely, if we assume that the two ciphertexts have their noise bounded by V , and $B_{\text{key}} = 1$, the size of the noise of their multiplication can be reasonably approximated by

$$\delta_{\mathcal{R}} t \left((5 + \delta_{\mathcal{R}}) V + \frac{r_t(Q)}{2} (\delta_{\mathcal{R}} + 5) \right) + \frac{\delta_{\mathcal{R}}^2}{2} \approx \delta_{\mathcal{R}}^2 t \left(V + \frac{r_t(Q)}{2} \right). \quad (7)$$

Since the noise grows significantly with homomorphic multiplications, V becomes larger than $r_t(Q)/2 < t$ after we perform the first multiplication. However, this is not necessarily true for the first multiplication itself since, like in BGV, the noise of fresh ciphertexts in BFV is bounded by $B_{\text{err}}(2\delta_{\mathcal{R}}B_{\text{key}} + 1) \approx 2\delta_{\mathcal{R}}B_{\text{err}}$. The homomorphic encryption standard [4] recommends using an error distribution with $\sigma_{\text{err}} = 3.2$. Therefore, the noise of a fresh ciphertext can be estimated as $V_{\text{init}} = 2 \times 6 \times 3.2 \times 2\sqrt{N} < 77\sqrt{N}$. Since in practice the dimension N typically does not exceed 2^{16} , a fresh ciphertext always has its noise size not higher than 14 bits, while $r_t(Q)$ can be as large as t . As a consequence, when $r_t(Q)$ is larger than 2^{15} , it becomes responsible for the larger noise growth after

the first multiplication in BFV. For instance, if $t = 2^{32}$ and $r_t(Q) \approx t/2 \approx 2^{31}$, the noise after the first multiplication will be at least 16 bits larger than in the case when $r_t(Q) < V_{\text{init}}$. Note that this difference will not lead to a larger noise growth on the next multiplications since, as shown in Eq. (7), the noise growth after a multiplication is linear in V . However, this difference of 16+ bits will be carried through until the end of the computation. In the case of $t = 2^{60}$, this difference would become at least 44 bits.

The easiest way to circumvent this problem would be to choose the moduli q_i , as in the original BGV, i.e., such that $q_i \equiv 1 \pmod t$, which would lead to $r_t(Q) = 1$. However, for the same reasons as in BGV, this restriction would make the finding of the moduli challenging for large t . Although it is possible to relax this condition by choosing, for instance, $r_t(Q) < \sqrt{N}$, i.e., finding $r_t(Q)$ through trial and error, we believe this would be a patch rather than a real solution. We show next that there is a more natural way to fix this problem.

Indeed, the $r_t(Q)$ term comes from the difference between Δ^{-1} and t/Q since when computing $\Delta t/Q$ (during decryption and homomorphic multiplication), one obtains $1 - r_t(Q)/Q$. Therefore, to solve this issue we propose to modify the original BFV encryption algorithm by encoding $[\mathbf{m}]_t$ in the ciphertext in a more natural way as $\lfloor Q[\mathbf{m}]_t/t \rfloor$ instead of $\Delta[\mathbf{m}]_t$. The first benefit is seen in the decryption bound since now

$$\begin{aligned} \left\lfloor \frac{t}{Q} [c_0 + c_1 \cdot s]_Q \right\rfloor &= \left\lfloor \frac{t}{Q} \left(\frac{Q}{t} [\mathbf{m}]_t + \mathbf{v} + \varepsilon + \mathbf{k}Q \right) \right\rfloor = [\mathbf{m}]_t + \left\lfloor \frac{t}{Q} (\mathbf{v} + \varepsilon) \right\rfloor + t\mathbf{k} \\ &= [\mathbf{m}]_t + \left\lfloor \frac{t}{Q} (\mathbf{v} + \varepsilon) \right\rfloor \pmod t, \end{aligned}$$

where $\mathbf{k} \in \mathcal{R}$ and ε is the rounding error coming from $\lfloor Q[\mathbf{m}]_t/t \rfloor = Q[\mathbf{m}]_t/t + \varepsilon$, such that $\|\varepsilon\|_\infty \leq 1/2$. Therefore the decryption will be correct as long as:

$$\frac{t}{Q} \|\mathbf{v} + \varepsilon\|_\infty < \frac{1}{2},$$

which is satisfied if

$$\|\mathbf{v}\|_\infty < \frac{Q}{2t} - \frac{1}{2}. \tag{8}$$

Remark 2. Note that one can compute $\lfloor Q[\mathbf{m}]_t/t \rfloor \pmod Q$ directly in RNS as long as t and Q are coprime since

$$\left\lfloor \frac{Q[\mathbf{m}]_t}{t} \right\rfloor = \frac{Q[\mathbf{m}]_t - [Q\mathbf{m}]_t}{t} = -\frac{[Q\mathbf{m}]_t}{t} \pmod Q.$$

The second benefit is observed in the addition since now we have

$$\begin{aligned} c_0 + c_1 \cdot s + c'_0 + c'_1 \cdot s &= \frac{Q}{t} ([\mathbf{m}]_t + [\mathbf{m}']_t) + \mathbf{v} + \mathbf{v}' + \varepsilon + \varepsilon' \\ &= \frac{Q}{t} ([\mathbf{m} + \mathbf{m}']_t + t\mathbf{u}) + \mathbf{v} + \mathbf{v}' + \varepsilon + \varepsilon' \\ &= \frac{Q}{t} [\mathbf{m} + \mathbf{m}']_t + \mathbf{v} + \mathbf{v}' + \varepsilon + \varepsilon' \pmod Q. \end{aligned}$$

Hence the noise after a homomorphic addition is bounded by

$$\|\mathbf{v}_{\text{new-add}}\|_\infty \leq \|\mathbf{v}\|_\infty + \|\mathbf{v}'\|_\infty + 1. \tag{9}$$

Note that the decryption bound (8) and the addition bound (9) are now exactly the same as for BGV.

The equations for homomorphic multiplication can be simplified the same way. Denoting $\tilde{\mathbf{v}} = \mathbf{v} + \varepsilon$, $\mathbf{ct}_{\text{tensor}}$ is computed as

$$\begin{aligned} (\mathbf{c}_0 + \mathbf{c}_1 \cdot \mathbf{s}) \cdot (\mathbf{c}'_0 + \mathbf{c}'_1 \cdot \mathbf{s}) &= \left(\frac{Q}{t}[\mathbf{m}]_t + \tilde{\mathbf{v}} + \mathbf{k}Q\right) \cdot \left(\frac{Q}{t}[\mathbf{m}']_t + \tilde{\mathbf{v}}' + \mathbf{k}'Q\right) \\ &= \frac{Q^2}{t^2}[\mathbf{m} \cdot \mathbf{m}']_t + \frac{Q}{t}\mathbf{v}_{\text{new-tensor}} + \frac{Q^2}{t}\mathbf{k}_{\text{new-tensor}}, \end{aligned}$$

where

$$\begin{aligned} \mathbf{v}_{\text{new-tensor}} &= [\mathbf{m}]_t \cdot \tilde{\mathbf{v}}' + [\mathbf{m}']_t \cdot \tilde{\mathbf{v}} + \frac{t}{Q}\tilde{\mathbf{v}} \cdot \tilde{\mathbf{v}}' + t(\tilde{\mathbf{v}} \cdot \mathbf{k}' + \tilde{\mathbf{v}}' \cdot \mathbf{k}), \\ \mathbf{k}_{\text{new-tensor}} &= [\mathbf{m}]_t \cdot \mathbf{k}' + [\mathbf{m}']_t \cdot \mathbf{k} + t\mathbf{k} \cdot \mathbf{k}' + \mathbf{r}_m. \end{aligned}$$

Then after the scaling by t/Q and rounding, similarly to the original BFV scheme, the noise of the multiplication is given by: $\mathbf{v}_{\text{new-mult}} = \mathbf{v}_{\text{new-tensor}} + \mathbf{v}_r$, and is bounded by:

$$\begin{aligned} \|\mathbf{v}_{\text{new-mult}}\|_\infty &\leq \frac{\delta_{\mathcal{R}}t}{2} \left(\frac{2\|\tilde{\mathbf{v}}\|_\infty \|\tilde{\mathbf{v}}'\|_\infty}{Q} + (4 + \delta_{\mathcal{R}}B_{\text{key}})(\|\tilde{\mathbf{v}}\|_\infty + \|\tilde{\mathbf{v}}'\|_\infty) \right) \\ &\quad + \frac{1 + \delta_{\mathcal{R}}B_{\text{key}} + \delta_{\mathcal{R}}^2B_{\text{key}}^2}{2}. \end{aligned} \tag{10}$$

We will see in Sect. 5 that this bound is similar to the bound for BGV.

Similarly to [27], we can derive a bound on the noise after having evaluated a binary tree of depth L from (10). By assuming that the size of the noise of ciphertexts is bounded by V before the first multiplication, the noise of the resulting ciphertext will be bounded by

$$C_1^L V + C_2 \sum_{i=0}^{L-1} C_1^i \leq C_1^L V + LC_2 C_1^{L-1}$$

with $C_1 = \delta_{\mathcal{R}}t(5 + \delta_{\mathcal{R}}B_{\text{key}})$ and $C_2 = (1 + \delta_{\mathcal{R}}B_{\text{key}} + \delta_{\mathcal{R}}^2B_{\text{key}}^2)/2 + V_{\text{ks}}$ where V_{ks} represents the noise added by the key-switching (see the full version).

Last we want to highlight further the similarity between BGV and BFV. Just like in the GHS variant of BGV, one can encrypt a ciphertext \mathbf{ct} in BFV using a slightly larger modulus Qp and then rescale the ciphertext by $1/p$, which will have the same effect as modulus switching in BGV:

$$\mathbf{ct}_{\text{scale}} = \left\lfloor \frac{1}{p} \mathbf{ct} \right\rfloor \bmod Q,$$

so that its noise after scaling is bounded by

$$\|\mathbf{v}_{\text{scale}}\|_{\infty} \leq \frac{\|\mathbf{v}\|_{\infty}}{p} + \frac{1}{2p} + \frac{1 + \delta_{\mathcal{R}} B_{\text{key}}}{2}.$$

Therefore, like in BGV, this allows to reduce the noise of a freshly encrypted to $(1 + \delta_{\mathcal{R}} B_{\text{key}})/2 \approx \delta_{\mathcal{R}}/2$. Note that the noise benefit of the BFV encryption proposed in our work will become more significant as the fresh noise is several bits larger than the modulus switching noise. Moreover, when using GHS or Hybrid key switching, p can be chosen as one of the moduli of the key-switching basis, and therefore this technique will not impact the selection of Ring-LWE security parameters.

3.2 Modified Homomorphic Multiplication

In the previous subsection, we showed how to instantiate BFV in such a way that it is not worse than BGV in terms of noise growth. Now the main difference left between the two schemes is in the complexity of their homomorphic multiplication procedure. In a nutshell, in BGV the tensoring can be done directly modulo Q while in BFV it must be done without any modular reduction. In practice, this requires using a second RNS basis $P = p_1 \cdots p_{k'}$ such that $\text{ct}(\mathbf{s}) \cdot \text{ct}'(\mathbf{s}')$ does not wrap around modulo PQ . More precisely, the critical part in practice is to avoid the wraparound of the dominant term $Qt\mathbf{k} \cdot \mathbf{k}'$ modulo P during the scaling (see Sect. 2.2). This requires to choose $P > t\delta_{\mathcal{R}}^3 Q/4$, which in practice is satisfied by setting $k' = k + 1$ for the RNS instantiation. Algorithm 1 recalls the original homomorphic multiplication procedure of BFV.

Algorithm 1 . Original BFV Multiplication Algorithm

procedure ORIGINALMULT($\text{ct} = (\mathbf{c}_0, \mathbf{c}_1) \in \mathcal{R}_Q^2, \text{ct}' = (\mathbf{c}'_0, \mathbf{c}'_1) \in \mathcal{R}_Q^2$)
 Expand: $\text{ct} \in \mathcal{R}_Q^2$ and $\text{ct}' \in \mathcal{R}_Q^2 \rightarrow \text{ct} \in \mathcal{R}_{QP}^2$ and $\text{ct}' \in \mathcal{R}_{QP}^2$:
 ▷ $\text{ct}(\mathbf{s}) = \Delta[\mathbf{m}]_t + \mathbf{v} + Q\mathbf{k} \pmod{\mathcal{R}_{QP}}$
 ▷ $\text{ct}'(\mathbf{s}') = \Delta[\mathbf{m}']_t + \mathbf{v}' + Q\mathbf{k}' \pmod{\mathcal{R}_{QP}}$
 Tensor: $\text{ct}_{\text{tensor}} = (\mathbf{c}_0 \cdot \mathbf{c}'_0, \mathbf{c}_0 \cdot \mathbf{c}'_1 + \mathbf{c}_1 \cdot \mathbf{c}'_0, \mathbf{c}_1 \cdot \mathbf{c}'_1) \in \mathcal{R}_{QP}^3$:
 ▷ $\text{ct}_{\text{tensor}}(\mathbf{s}) = \frac{Q}{t} \Delta[\mathbf{m} \cdot \mathbf{m}']_t + \frac{Q}{t} \mathbf{v}_{\text{tensor}} + \frac{Q^2}{t} \mathbf{k}_{\text{tensor}} \pmod{\mathcal{R}_{QP}}$
 ScaleDown: $\text{ct}_{\text{scale}} = \lfloor \frac{t}{Q} \text{ct}_{\text{tensor}} \rfloor \in \mathcal{R}_P^3$
 ▷ $\text{ct}_{\text{scale}}(\mathbf{s}) = \Delta[\mathbf{m} \cdot \mathbf{m}']_t + \mathbf{v}_{\text{tensor}} + Q\mathbf{k}_{\text{tensor}} + \mathbf{v}_r \pmod{\mathcal{R}_P}$
 SwitchBasis: $\text{ct}_{\text{scale}} \in \mathcal{R}_P^3 \rightarrow \text{ct}_{\text{scale}} \in \mathcal{R}_Q^3$:
 ▷ $\text{ct}_{\text{scale}}(\mathbf{s}) = \Delta[\mathbf{m} \cdot \mathbf{m}']_t + \mathbf{v}_{\text{tensor}} + \mathbf{v}_r \pmod{\mathcal{R}_Q}$

We propose a new homomorphic multiplication algorithm, with a reduced computational complexity. Instead of multiplying two ciphertexts modulo Q and dealing with a multiple of Q^2 modulo QP , the idea is to switch one of the two ciphertexts to modulus P so that after the tensoring we obtain a multiple of PQ that vanishes modulo PQ . As explained in the above paragraph, this would

allow us to reduce the size of P since the original dominant term would now disappear. More precisely, the procedure starts as usual with two ciphertexts encrypted modulo Q :

$$\text{ct}(s) = \frac{Q}{t}[\mathbf{m}]_t + \tilde{\mathbf{v}} + \mathbf{k}Q \text{ and } \text{ct}'(s) = \frac{Q}{t}[\mathbf{m}']_t + \tilde{\mathbf{v}}' + \mathbf{k}'Q.$$

with $\tilde{\mathbf{v}} = \mathbf{v} + \varepsilon$ and $\tilde{\mathbf{v}}' = \mathbf{v}' + \varepsilon'$, like in Sect. 3.1.

Then one can perform the modulus switching of one of the two ciphertexts, say ct , to convert it to modulus P by computing

$$\hat{\text{ct}} = \left\lfloor \frac{P}{Q} \text{ct} \right\rfloor \text{ mod } P,$$

which satisfies:

$$\hat{\text{ct}}(s) = \frac{P}{t}[\mathbf{m}]_t + \frac{P\tilde{\mathbf{v}}}{Q} + \mathbf{k}P + \varepsilon_{\text{round}} \text{ mod } P,$$

where the rounding error $\varepsilon_{\text{round}} = \hat{\text{ct}}(s) - P\text{ct}(s)/Q$ has its norm bounded by $(1 + \delta_{\mathcal{R}}B_{\text{key}})/2$ as usual. From there, $\hat{\text{ct}}$ is expanded from P to QP , ct' is expanded from Q to QP , and one can perform the tensoring as usual to obtain

$$\begin{aligned} \hat{\text{ct}}_{\text{tensor}}(s) &= \frac{PQ}{t^2}[\mathbf{m} \cdot \mathbf{m}']_t + \frac{P}{t}\hat{\mathbf{v}}_{\text{tensor}} + \frac{QP}{t}\hat{\mathbf{k}}_{\text{tensor}} \\ &\quad + \varepsilon_{\text{round}} \cdot \left(\frac{Q}{t}[\mathbf{m}']_t + \tilde{\mathbf{v}}' + \mathbf{k}'Q \right) \text{ mod } PQ, \end{aligned}$$

with $\hat{\mathbf{v}}_{\text{tensor}} = \mathbf{v}_{\text{new-tensor}}$ from Sect. 3.1 and

$$\hat{\mathbf{k}}_{\text{tensor}} = [\mathbf{m}]_t \cdot \mathbf{k}' + [\mathbf{m}']_t \cdot \mathbf{k} + \mathbf{r}_m \in \mathcal{R}.$$

Note that this time $\hat{\mathbf{k}}_{\text{tensor}}$ does not contain a multiple of $\mathbf{k} \cdot \mathbf{k}'$. Then to get back a valid ciphertext modulo Q , one must scale down the result by t/P , instead of t/Q in the original case, which leads to

$$\frac{t}{P}\hat{\text{ct}}_{\text{tensor}}(s) = \frac{Q}{t}[\mathbf{m} \cdot \mathbf{m}']_t + \hat{\mathbf{v}}_{\text{tensor}} + Q\hat{\mathbf{k}}_{\text{tensor}} + \frac{t\varepsilon_{\text{round}}}{P} \cdot \left(\frac{Q}{t}[\mathbf{m}']_t + \tilde{\mathbf{v}}' + \mathbf{k}'Q \right).$$

After the rounding, the multiple of Q will vanish modulo Q and one will have to take into account the rounding error term \mathbf{v}_r of norm bounded by $(1 + \delta_{\mathcal{R}}B_{\text{key}} + \delta_{\mathcal{R}}^2B_{\text{key}}^2)/2$, which adds to the noise, like in the original BFV scheme. Therefore, the noise of this variant of homomorphic multiplication is bounded by

$$\|\hat{\mathbf{v}}_{\text{new-mult}}\|_{\infty} \leq \|\mathbf{v}_{\text{new-mult}}\|_{\infty} + \frac{t\delta_{\mathcal{R}}(\delta_{\mathcal{R}}B_{\text{key}} + 1)}{2P} \left(\|\tilde{\mathbf{v}}\|_{\infty}' + \frac{Q(\delta_{\mathcal{R}}B_{\text{key}} + 4)}{2} \right). \tag{11}$$

Notice that the only difference in the noise growth between Eq. (10) and Eq. (11) is due to rounding error $\varepsilon_{\text{round}}$ occurring during the first modulus switching.

However, we can control the size of this noise with P . As explained in Sect. 3.1, the norm of $\mathbf{v}_{\text{new-mult}}$ is dominated by $\delta_{\mathcal{R}}^2 tV$, where V is the bound on the size of the noise of $\tilde{\mathbf{v}}$ and $\tilde{\mathbf{v}}'$. If we look carefully at the other term and choose $P \approx Q$, it will be dominated by $\delta_{\mathcal{R}}^3 t/4$. Since the noise of a fresh ciphertext, even scaled-down after encryption, has always a size larger than $\delta_{\mathcal{R}}/2$, this error term will add at most half a bit to the noise of the first multiplication in the worst case, which can be considered as negligible in practice as a larger cushion is typically added to the heuristic expression for δ_R , or, more precisely, δ_R^3 in this case. This means that one can choose $P \approx Q$ and hence $k' = k$, instead of $k' = k + 1$ in the original case, which reduces the size of P and still achieves the same noise growth as in Sect. 3.1. We summarize our new multiplication algorithm in Algorithm 2.

Algorithm 2 . New BFV Multiplication Algorithm

procedure NEWMULT($\mathbf{ct} = (c_0, c_1) \in \mathcal{R}_Q^2, \mathbf{ct}' = (c'_0, c'_1) \in \mathcal{R}_Q^2$)
 ModSwitch: $\mathbf{ct} \in \mathcal{R}_Q^2 \rightarrow \hat{\mathbf{ct}} \in \mathcal{R}_P^2$
 Expand: $\hat{\mathbf{ct}} \in \mathcal{R}_P^2 \rightarrow \hat{\mathbf{ct}} \in \mathcal{R}_{QP}^2$ and $\mathbf{ct}' \in \mathcal{R}_Q^2 \rightarrow \mathbf{ct}' \in \mathcal{R}_{QP}^2$
 $\triangleright \hat{\mathbf{ct}}(s) = \frac{P}{t} [\mathbf{m}]_t + \frac{P}{Q} \tilde{\mathbf{v}} + P\mathbf{k} + \varepsilon_{\text{round}} \pmod{\mathcal{R}_{QP}}$
 $\triangleright \mathbf{ct}'(s) = \frac{Q}{t} [\mathbf{m}']_t + \tilde{\mathbf{v}}' + Q\mathbf{k}' \pmod{\mathcal{R}_{QP}}$
 Tensor: $\hat{\mathbf{ct}}_{\text{tensor}} = (\hat{c}_0 \cdot c'_0, \hat{c}_0 \cdot c'_1 + \hat{c}_1 \cdot c'_0, \hat{c}_1 \cdot c'_1) \in \mathcal{R}_{QP}^3$:
 $\triangleright \hat{\mathbf{ct}}_{\text{tensor}}(s) = \frac{QP}{t^2} [\mathbf{m} \cdot \mathbf{m}']_t + \frac{P}{t} \hat{\mathbf{v}}_{\text{tensor}} + \frac{QP}{t} \hat{\mathbf{k}}_{\text{tensor}} \pmod{\mathcal{R}_{QP}}$
 ScaleDown: $\mathbf{ct}_{\text{scale}} = \lfloor \frac{t}{P} \hat{\mathbf{ct}}_{\text{tensor}} \rfloor \in \mathcal{R}_Q^3$
 $\triangleright \mathbf{ct}_{\text{scale}}(s) = \frac{Q}{t} [\mathbf{m} \cdot \mathbf{m}']_t + \hat{\mathbf{v}}_{\text{tensor}} + \mathbf{v}_r \pmod{\mathcal{R}_Q}$

Remark 3. Notice that since one must scale the tensored ciphertext by t/P instead of t/Q , it can be done directly in the basis Q . Therefore the homomorphic multiplication procedure requires 2 basis extensions from Q to P for $\mathbf{ct} = (c_0, c_1)$ in the beginning, 4 more basis extensions to expand the two ciphertexts $\hat{\mathbf{ct}} = (\hat{c}_0, \hat{c}_1)$ and $\mathbf{ct}' = (c'_0, c'_1)$ from P and Q , respectively, to PQ . Finally one needs 3 additional basis extensions from PQ to Q for $\hat{\mathbf{ct}}_{\text{tensor}} \in \mathcal{R}_{PQ}^3$ to perform the downscaling modulo Q . Thus it requires a total of 9 basis extensions instead of $4 + 3 + 3 = 10$ in the original BFV algorithm. Moreover, since P is slightly smaller, each basis extension will be cheaper to compute.

Leveled BFV Multiplication. If one wants to make BFV even closer to BGV in terms of performance, one could consider a “leveled” approach to BFV by working with ciphertexts modulo $Q_\ell = q_1 \cdots q_\ell$ and performing modulus switching as the computation progresses. However, as in BGV, one would have to manage ciphertexts at different levels and deal with more challenging noise estimation.

To keep the usability of BFV, we propose instead a “leveled” multiplication that pre-scales both ciphertexts by $\frac{Q_\ell}{Q}$ (using internal modulus switching to Q_ℓ), and then multiplies the result by $\frac{Q}{Q_\ell}$ after the multiplication procedure. In this case, the ciphertexts will always stay modulo Q outside the multiplication

procedure, but the multiplication will be done internally modulo $Q_\ell < Q$ and hence will be more efficient.

In this case, the noise of input ciphertexts after internal modulus switching from Q to Q_ℓ will be equal to

$$\hat{v} = \frac{Q_\ell}{Q}v + \varepsilon_{\text{round}} \text{ and } \hat{v}' = \frac{Q_\ell}{Q}v' + \varepsilon'_{\text{round}},$$

where $\varepsilon_{\text{round}}$ and $\varepsilon'_{\text{round}}$ have their norm bounded by $(1 + \delta_{\mathcal{R}}B_{\text{key}})/2 \approx \delta_{\mathcal{R}}/2$. On the one hand, the main consideration here is to choose ℓ such that $\frac{Q_\ell}{Q} \|v\|_\infty$ remains significantly larger than $\|\varepsilon_{\text{round}}\|_\infty \approx \frac{\delta_{\mathcal{R}}}{2}$, so that the noise brought about by the modulus switching procedure will not significantly impact the overall noise growth. This is equivalent to

$$\|v\|_\infty \gg \frac{Q\delta_{\mathcal{R}}}{2Q_\ell} \text{ and } \|v'\|_\infty \gg \frac{Q\delta_{\mathcal{R}}}{2Q_\ell},$$

or in practice

$$\|v\|_\infty > 8\frac{Q\delta_{\mathcal{R}}}{Q_\ell} \text{ and } \|v'\|_\infty > 8\frac{Q\delta_{\mathcal{R}}}{Q_\ell}. \tag{12}$$

On the other hand, in order to gain as much as possible in efficiency, Q_ℓ must be chosen as the smallest modulus satisfying inequalities (12). Theoretically this requires to have a precise estimate of the average (or lower bound within a certain confidence interval) noise size. But in practice it is enough to add a heuristic “cushion” to our worst-case bound. See the full version for details.

Algorithm 3 . New Leveled BFV Multiplication Algorithm

procedure LEVELEDNEWMULT($\text{ct} = (c_0, c_1) \in \mathcal{R}_Q^2, \text{ct}' = (c'_0, c'_1) \in \mathcal{R}_Q^2$)
 ModSwitchDown: $\hat{\text{ct}} = \lfloor \frac{Q_\ell}{Q} \text{ct} \rfloor \in \mathcal{R}_{Q_\ell}^2$ and $\hat{\text{ct}}' = \lfloor \frac{Q_\ell}{Q} \text{ct}' \rfloor \in \mathcal{R}_{Q_\ell}^2$
 $\triangleright \hat{v} = \frac{Q_\ell}{Q}v + \varepsilon_{\text{round}}$ and $\triangleright \hat{v}' = \frac{Q_\ell}{Q}v' + \varepsilon'_{\text{round}}$
 $\hat{\text{ct}}_m = \text{NewMult}(\hat{\text{ct}}, \hat{\text{ct}}') \in \mathcal{R}_{Q_\ell}^3$
 $\triangleright \hat{\text{ct}}_m(s) = \frac{Q_\ell}{t} [m_1 m_2]_t + \hat{v}_m \pmod{\mathcal{R}_{Q_\ell}}$
 ModSwitchUp: $\text{ct}_m = \lfloor \frac{Q}{Q_\ell} \hat{\text{ct}}_m \rfloor \in \mathcal{R}_Q^3$
 $\triangleright \text{ct}_m(s) = \frac{Q}{t} [m_1 m_2]_t + v_m \pmod{\mathcal{R}_Q}$
 $\triangleright v_m = \frac{Q}{Q_\ell} \hat{v}_m + v_r$

Remark 4. Note that this “leveled” optimization can be equally applied to key switching. The only difference in this case would be to ensure that the noise of the scaled ciphertext remains larger than the noise brought about by the key-switching procedure itself.

Table 1. Complexities of different multiplication methods.

	# NTTs	# integer-mult	# floating-point-oper
MultOld	$14k + 7$	$(10k^2 + 26k + 9)n$	$(10k + 3)n$
MultNew	$14k$	$(9k^2 + 15k)n$	$7kn$
MultNewLeveled	14ℓ	$(4k\ell + 5\ell^2 + 2k + 18\ell)n$	$(2k + 5\ell)n$

Remark 5. Modulus switching from Q to Q_ℓ and then from Q_ℓ to P_ℓ in Algorithm 3 can be combined into a single modulus switching from Q directly to P_ℓ . This reduces the number of integer multiplications from $(k\ell + \ell^2 + 2\ell)n$ to $(k\ell + \ell)n$. Note that an approximate modulus switching (instead of an exact one with a floating-point correction technique from [21]) can be employed by adding extra $\log k$ bits to the noise estimate used for reducing the number of levels inside homomorphic multiplication. Both exact and approximate procedures for switching the moduli of ciphertexts between arbitrary RNS bases are described in the full version.

Table 1 summarizes the computational complexities of different multiplication algorithms by assuming that the extension from Q to QP is performed using the technique from [21] with floating-point operations.

Lazy Scaling in BFV Multiplication. An additional optimization can be implemented by noticing that tensoring and scaling can be separated to optimize some evaluation circuits. For example, consider the inner product circuit of two vectors of ciphertexts. We evaluate it by multiplying (tensoring and scaling) the pairs of ciphertexts and then adding the results (mod \mathcal{R}_Q). It is more efficient to do this in a different way: first we apply the tensoring subroutine to the pairs of ciphertexts, then add the results (mod \mathcal{R}_{QP}), and finally perform the expensive scaling subroutine only once. Indeed, after tensoring we already have the information about the multiplicative noise v_{tensor} , thus changing the order of scaling and additions does not affect the v_{tensor} noise. Moreover, as we perform the scaling down only once instead of doing it for each pair of ciphertexts, the total noise from the inner product is actually reduced. We call this technique *lazy scaling* and describe the pseudocode in the full version. The experimental results in the full version suggest that this optimization can speed up inner products by more than 2x in practice.

3.3 Improved Decryption in the HPS RNS Variant

A significant practical limitation of the HPS approach is that high-precision (“long double” or even quad-precision) floating-point arithmetic is required to support larger CRT moduli [21]. We introduce a general-purpose digit decomposition technique (inspired by digit decomposition in key switching) and apply it to the HPS decryption procedure to add support for arbitrary CRT moduli

using only regular double-precision floating-point arithmetic, hence overcoming this limitation of the HPS variant.

The idea of HPS scaling [21] for decryption can be briefly explained as follows: for $x \in \mathbb{Z}_Q$ with CRT representation (x_1, \dots, x_k) we want to compute an integer $y = \lceil t/Qx \rceil \in \mathbb{Z}_t$, and use a CRT composition formula to derive the following expression:

$$\begin{aligned}
 y &:= \left\lceil \frac{t}{Q}x \right\rceil = \left\lceil \left(\sum_{i=1}^k x_i \cdot \left[\left(\frac{Q}{q_i} \right)^{-1} \right]_{q_i} \cdot \frac{Q}{q_i} \cdot \frac{t}{Q} \right) - u \cdot Q \cdot \frac{t}{Q} \right\rceil \\
 &= \left\lceil \left[\left(\sum_{i=1}^k x_i \cdot \left(\left[\left(\frac{Q}{q_i} \right)^{-1} \right]_{q_i} \cdot \frac{t}{q_i} \right) \right) \right] \right\rceil_t = \left[\left(\sum_{i=1}^k x_i \cdot \omega_i \right) + \left[\sum_{i=1}^k x_i \cdot \theta_i \right] \right]_t, \\
 &\text{where } \left[\left(\frac{Q}{q_i} \right)^{-1} \right]_{q_i} \cdot \frac{t}{q_i} = \omega_i + \theta_i \text{ with } \omega_i \in \mathbb{Z}_t \text{ and } \theta_i \in \left[-\frac{1}{2}, \frac{1}{2} \right).
 \end{aligned}$$

As we can only store approximate values $\tilde{\theta}_i = \theta_i + \epsilon_i$, the magnitude of the error term $|e'| = |\sum_i x_i \epsilon_i|$ in the fractional part is limited by $kq_m \epsilon_m$, where $q_m = \max_i(q_i)$ and $\epsilon_m = \max_i(\epsilon_i)$. If we restrict the floating-point precision to “doubles”, which are natively supported by modern CPUs, we have to introduce a constraint $kq_m < 2^{51}$. To support larger CRT moduli, we need “long doubles” or even quad-precision arithmetic: long doubles are needed for CRT moduli from 47 to 58 bits, and quad-precision floats are needed for higher CRT moduli [21].

Our main idea is to perform digit decomposition, somewhat similar to how digit decomposition is done in BV key-switching, to replace the factor q_m with a smaller digit of it. For base $B_s \in \mathbb{Z}$, $B_s \geq 2$, let $d_s = \lceil \log(q_m)/\log(B_s) \rceil$. Let $x_i = \sum_{j=0}^{d_s-1} x_{i,j} \cdot B_s^j$ be the B_s base decomposition of x_i . Then the expression for y can be rewritten as

$$\begin{aligned}
 y &= \left\lceil \left[\left(\sum_{i=1}^k \sum_{j=0}^{d_s-1} x_{i,j} \cdot \left(\frac{t}{q_i} \cdot \left[\left(\frac{Q}{q_i} \right)^{-1} \right]_{q_i} \cdot B_s^j \right) \right) \right] \right\rceil_t \\
 &= \left[\left(\sum_{i=1}^k \sum_{j=0}^{d_s-1} [x_{i,j} \cdot \omega_{i,j}]_t \right) + \left[\sum_{i=1}^k \sum_{j=0}^{d_s-1} x_{i,j} \cdot \theta_{i,j} \right] \right]_t,
 \end{aligned}$$

where $\frac{t}{q_i} \cdot \left[\left(\frac{Q}{q_i} \right)^{-1} \right]_{q_i} \cdot B_s^j = \omega_{i,j} + \theta_{i,j}$, with $\omega_{i,j} \in \mathbb{Z}_t$ and $\theta_{i,j} \in \left[-\frac{1}{2}, \frac{1}{2} \right)$.

Note that $\omega_{i,j}$ and $\theta_{i,j}$ are the new precomputation factors instead of ω_i and θ_i .

Error Analysis. The magnitude of the error term $|e'| = |\sum_{i,j} x_{i,j} \epsilon_{i,j}|$ is now limited by $|\sum_{i,j} x_{i,j} \epsilon_{i,j}| < kd_s B_s \epsilon_m$. In practice, our moduli q_i are normally

bounded by 64 (or often by 60) bits. We have already considered the case of $kq_m < 2^{51}$. If $kq_m > 2^{51}$, we can take $d_s = 2$, $B_s = 2^{\lceil \log_2 q_m / 2 \rceil}$. Then $|\epsilon'| < 2^{-19}k < 1/4$ for $k < 2^{17}$. Hence the floating-point error will have no effect on the result for any practically reasonable value of k .

Complexity Analysis. The procedure takes kd_s floating-point multiplications, kd_s modular integer multiplications, some modular additions, and one rounding to compute $\lceil u \rceil$. However, if $tkd_s B_s < 2^{64}$, then we can replace modular multiplications and modular additions by plain integer multiplications and additions respectively, and do one modular reduction at the end.

Remark 6. Note that this digit decomposition technique can be applied to other mixed integer/floating-point RNS operations to reduce precision requirements for floating-point arithmetic or avoid extra noise due to floating-point rounding. For instance, it can be used in the scaling for homomorphic multiplication.

4 More Usable BGV Scheme

The practical use of the BGV scheme requires accurate dynamic noise estimation to decide when the modulus operation should be executed, and what scaling factor should be chosen for modulus switching [20]. Each modulus switching decision may significantly impact the noise not only for the current operation, but also for all subsequent operations. An error in a noise estimate may eventually lead even to a decryption failure. Therefore, fine-tuned noise estimation techniques are used to estimate the noise for various operations (see [23] for a more detailed discussion). In contrast, the BFV scheme is much more robust to inaccuracies in noise estimation and typically only requires an upper bound on the error for the desired multiplicative depth. This robustness of BFV is related to the use of the MSD encoding and scaling down by a large factor Q/t during BFV homomorphic multiplication, and the “fragility” of BGV is caused by the LSD encoding and scaling down by a small factor, comparable in magnitude to the noise incurred in operations after the previous modulus switching. For this reason, many modern homomorphic encryption libraries implement BFV as the scheme for finite fields, while only HELib and PALISADE (since quite recently) provide efficient implementations of BGV.

We present an alternative approach for instantiating BGV, which does not require dynamic noise estimation. For this instantiation, one only needs to specify the multiplicative depth of the computation, maximum number of additions per multiplicative level, and the number of additions and automorphisms before first multiplication. Then all moduli $Q_L, Q_{L-1}, \dots, Q_1, Q_0$ are chosen so that a small, constant level of noise can be maintained throughout the computation. All modulus switching operations are automatically performed right before a homomorphic multiplication. Conceptually, this BGV instantiation is similar in usability to BFV, where only the multiplicative depth needs to be specified upfront and all “modulus switching” operations are performed automatically.

The logic for choosing the moduli is as follows. We start with a fresh encryption that has a noise $\|\mathbf{v}_{\text{fresh}}\|_{\infty} = B_{\text{err}}(2\delta_{\mathcal{R}}B_{\text{key}} + 1)$. Then we perform automorphism operations and additions, and apply modulus switching right before the first multiplication. This additional modulus switching before first multiplication allows us to reset the noise to a value comparable to the modulus switching noise, which will be the constant noise level $\|\mathbf{v}_{\text{c}}\|_{\infty}$ we will maintain throughout the computation. This can be expressed as

$$\frac{Q_L}{Q_{L+1}} ((n_{\text{add}} + 1) \|\mathbf{v}_{\text{fresh}}\|_{\infty} + n_{\text{ks}} \|\mathbf{v}_{\text{ks}}\|_{\infty}) + \frac{1 + \delta_{\mathcal{R}}B_{\text{key}}}{2} \leq \|\mathbf{v}_{\text{c}}\|_{\infty},$$

where n_{add} and n_{ks} are the numbers of additions and automorphism operations, respectively, that are performed before first multiplication, and $\|\mathbf{v}_{\text{ks}}\|_{\infty}$ is the bound on key switching noise. Note that here we introduced a new level and corresponding new modulus Q_{L+1} to account for an extra level we added before first multiplication. It is best to choose Q_{L+1}/Q_L such that $\|\mathbf{v}_{\text{c}}\|_{\infty} \approx 1 + \delta_{\mathcal{R}}B_{\text{key}}$ to achieve the smallest constant error because this error will allow us to minimize the subsequent values of Q_{i+1}/Q_i for most practical scenarios, hence resulting in the minimum value of ciphertext modulus Q_{L+1} (see the full version for more details, and more general expression for optimal $\|\mathbf{v}_{\text{c}}\|_{\infty}$).

Then for multiplication levels (from L to 1), we have to satisfy

$$\begin{aligned} \frac{Q_i}{Q_{i+1}} \left((n'_{\text{add}} + 1) \frac{\delta_{\mathcal{R}}t}{2} (2\|\mathbf{v}_{\text{c}}\|_{\infty}^2 + 2\|\mathbf{v}_{\text{c}}\|_{\infty} + 1) + (n'_{\text{ks}} + 1) \|\mathbf{v}_{\text{ks}}\|_{\infty} \right) \\ + \frac{1 + \delta_{\mathcal{R}}B_{\text{key}}}{2} \leq \|\mathbf{v}_{\text{c}}\|_{\infty}, \end{aligned}$$

where n'_{add} and n'_{ks} are the maximum numbers of additions and key switching operations, respectively, allowed per any multiplication level (going from L down to 1). For simplicity we use these maximum values across all levels so that Q_{i+1}/Q_i could have roughly same value for all $i \in \{1, \dots, L-1\}$. Note that for Hybrid key switching and relatively large plaintext moduli, such as $t = 2^{16} + 1$, which is often used for CRT packing, the multiplication noise is always much higher than $\|\mathbf{v}_{\text{ks}}\|_{\infty}$ (see derivations in the full version). Hence for this case we can rewrite the expression as

$$\frac{Q_i}{Q_{i+1}} \left((n'_{\text{add}} + 1) \frac{\delta_{\mathcal{R}}t}{2} (2\|\mathbf{v}_{\text{c}}\|_{\infty}^2 + 2\|\mathbf{v}_{\text{c}}\|_{\infty} + 1) \right) + \frac{1 + \delta_{\mathcal{R}}B_{\text{key}}}{2} \leq \|\mathbf{v}_{\text{c}}\|_{\infty}.$$

The last modulus Q_0 is chosen such that decryption is correct for a ciphertext with noise bounded by $\|\mathbf{v}_{\text{c}}\|_{\infty}$. This implies that $Q_0 > 2t\|\mathbf{v}_{\text{c}}\|_{\infty} - t$.

It is easy to show that once L , n_{add} , n_{ks} , n'_{add} , and n'_{ks} (only needed for small t) are given, all moduli Q_i can be derived. First we compute Q_0 , then all values Q_1, \dots, Q_L , and finally we can find Q_{L+1} .

This logic is simple to implement and avoids any dynamic noise estimation during the computation. It is also robust to inaccurate estimates as long as the upper bound for $\delta_{\mathcal{R}}$ is chosen appropriately, which is very similar to what is

done for BFV. There is a cost for this simplicity and robustness. The moduli Q_{L+1} may be larger than the values obtained using the more granular approach with dynamic noise estimation [23] because we use the maximum values of n'_{add} and n'_{ks} over all intermediate levels. However, our experimental results show that this instantiation of BGV can be significantly faster than the improved BFV implementation described in Sect. 3.

Remarks on the RNS Instantiation. Recall that for original BGV, we choose $Q = q_0 \cdots q_L$ and denote $Q_i = q_0 \cdots q_i$ for $0 \leq i \leq L$, where all $q_i = 1 \bmod 2N$ and co-prime to each other. In the case of our BGV variant, an extra q_{L+1} is introduced to reset the “fresh” noise to modulus switching noise. It is easy to show that for this setup, $Q_0 = q_0$, $Q_{i+1}/Q_i = q_{i+1}$, and $Q_{L+1}/Q_L = q_{L+1}$.

The expressions for finding q_0, q_i, q_{L+1} , where $i \in \{1, \dots, L\}$, can be written as follows:

$$q_0 > 2t \|\mathbf{v}_c\|_\infty - t, \quad (13)$$

$$q_i > 2 \left((n'_{\text{add}} + 1) \frac{\delta_{\mathcal{R}} t}{2} (2 \|\mathbf{v}_c\|_\infty + 2 + \frac{1}{\|\mathbf{v}_c\|_\infty}) + (n'_{\text{ks}} + 1) \frac{\|\mathbf{v}_{\text{ks}}\|_\infty}{\|\mathbf{v}_c\|_\infty} \right), \quad (14)$$

$$q_{L+1} > 2 \left((n_{\text{add}} + 1) \frac{\|\mathbf{v}_{\text{fresh}}\|_\infty}{\|\mathbf{v}_c\|_\infty} + n_{\text{ks}} \frac{\|\mathbf{v}_{\text{ks}}\|_\infty}{\|\mathbf{v}_c\|_\infty} \right),$$

where we take $\|\mathbf{v}_c\|_\infty = 1 + \delta_{\mathcal{R}} B_{\text{key}}$.

Handling Crosslevel Operations and Scaling Factors. The GHS variant implemented in HELib uses ciphertext-specific scaling factors, which introduces some complications that may affect the usability and may require additional scalar multiplications to bring two ciphertexts to the same scaling factor. In our BGV variant, we chose a simpler approach where the same scaling factor is used for all ciphertexts at a specific level, which reduces the number of scalar multiplications. This approach was originally introduced for the CKKS scheme in [25], and in our work we adapt it to BGV.

5 Comparison of BFV and BGV

5.1 Noise Growth

When comparing BGV and BFV, it is convenient to use the leveled approach of BGV, first comparing Q_0 , then Q_i , and finally Q_{L+1} .

For Q_0 , our modified variant of BFV has identical noise as BGV, i.e., Eq. (8) is exactly the same as Eq. (13).

For Q_{i+1}/Q_i , which corresponds to each multiplicative level, the dominant term in the BFV expression given by Eq. (11) is $\delta_R^2 t B_{\text{key}} V$, where V is the largest of the errors in two multiplied ciphertexts. For BGV, Eq. (14) suggests that the dominant term is $2\delta_R^2 t B_{\text{key}} V$. In other words, the expressions for BFV and BGV are identical except for the extra multiplicative factor of 2. This factor appears in BGV because we ensure that at each level the downscaled noise matches the added

modulus switching noise, keeping the noise level constant at twice the modulus switching noise (see the full version for details). In the case of BFV, the quadratic noise (product of prior noises for each ciphertext) is negligible as we downscale the ciphertext by a large factor Q/t , and we only observe the pure modulus switching noise. In other words, BFV has a small benefit of using one bit less per multiplication level.

There is also an extra advantage of BFV for small plaintext moduli, e.g., $t = 2$. As the analysis in the full version shows, the key switching noise in this case becomes comparable to multiplication noise for BGV, which implies higher values of Q_{i+1}/Q_i . In contrast, the key switching noise may only affect the initial level in BFV, as afterwards the accumulated noise from prior multiplications will be much higher than additive key switching noise, which is independent of current ciphertext noise. When we switch to larger plaintext moduli, this BFV advantage disappears as the key-switching noise in BGV becomes negligible compared to multiplication noise (as shown in the full version).

Using the $(L + 1)$ -th level (q_{L+1} in the RNS version) is preferred in BGV to achieve the smallest constant noise (twice the modulus switching noise). If $(L + 1)$ -th level is not used, then the fresh noise will make each Q_{i+1}/Q_i larger by a factor $\|\mathbf{v}_{\text{fresh}}\|_{\infty} / \|\mathbf{v}_c\|_{\infty} \approx 2B_{\text{err}} \approx 37$. Although one could use an auxiliary modulus in hybrid key switching during encryption instead (see the end of Sect. 3.1), extra noise can be accumulated from additions and/or key switching operations performed before first multiplication, which would increase all subsequent Q_{i+1}/Q_i . So the least level of constant noise in BGV, and hence smallest Q_{i+1}/Q_i , can be guaranteed only by introducing a relatively small extra “noise budget” for pseudo-level $L + 1$. Note that in BFV it is best to use an auxiliary modulus to reset the fresh noise to smaller modulus switching noise, without increasing the ciphertext modulus (see Sect. 3.1). Hence no pseudo-level $L + 1$ is needed in BFV, which is another small advantage of BFV over BGV.

In summary, the improved variants of BFV and BGV presented in this work have very similar noise growth, but BFV has some minor advantages over BGV, resulting in somewhat reduced ciphertext moduli needed to support the same computations.

5.2 Computational Complexity

The main difference between BFV and BGV in terms of computational complexity is due to the scaling method used in the multiplication operation. As was previously mentioned, BFV uses the MSD encoding and scales down the tensor product by a large Q/t factor, while BGV uses the LSD encoding technique to scale the tensor product only by a relatively small factor, comparable to the noise of previous multiplication. Considering that the noise growth is very similar in both schemes, one can expect that the computational complexity of BFV multiplication will be significantly higher. The purpose of this section is to quantify this difference, and examine the effect of plaintext moduli on this difference. Note that all other operations, such as addition and automorphism, use the same

approach in both schemes, and do not have any significant difference in terms of theoretical complexity. Hence we focus on the operation of multiplication.

The analysis in Sect. 3.2 shows that our leveled BFV multiplication takes 14ℓ NTTs and $(4k\ell + 5\ell^2 + 2k + 18\ell)n$ integer multiplications (we ignore for simplicity a much smaller contribution of floating-point operations). We also add the computational cost of hybrid key switching used for relinearization as there is a difference in its cost between BFV and BGV. For key-switching we assume that the ciphertext element is decomposed into $d_{\text{num}} = \ell/\alpha$ digits, i.e. each digit is considered modulo α moduli. The cost of relinearization for BFV is $4\ell + 2\alpha$ NTTs and $n(3\alpha\ell + 2d_{\text{num}}\ell + 2\alpha + 5\ell)$ integer multiplications (see the full version for details). Here, for simplicity of analysis we assume that ℓ is the same for leveled BFV multiplication and subsequent relinearization. The total cost of multiplication and relinearization in BFV is $17\ell + 2\alpha$ NTTs and $n(4k\ell + 5\ell^2 + 2k + 23\ell + 3\alpha\ell + 2d_{\text{num}}\ell + 2\alpha)$ integer multiplications.

In the case of our BGV variant, the total cost of multiplication includes two modulus switching operations for input ciphertexts, the tensor product, and, finally, the relinearization. The cost of modulus switching is $4(\ell' + 1)$ NTTs and $4n(\ell' + 2)$ integer multiplications, where ℓ' is the number of CRT moduli after modulus switching. The cost of tensor product is $4n\ell'$ integer multiplications. The cost of relinearization in the case of BGV is: $4\ell' + 2\alpha'$ NTTs and $n(3\alpha'\ell' + 2d_{\text{num}}\ell' + 4\alpha' + 7\ell')$ integer multiplications. Hence the total cost of multiplication and relinearization is $8\ell' + 4 + 2\alpha'$ NTTs and $n(3\alpha'\ell' + 2d_{\text{num}}\ell' + 4\alpha' + 15\ell' + 8)$ integer multiplications.

One can observe that the number of NTTs needed for BFV multiplication appears to be 2x or even higher than for BGV. But we should keep in mind that typically $\ell' > \ell$. For example, when $t = 2$, we can even have $\ell' > 3\ell$ since in BFV we work with large (60-bit) moduli vs the moduli of size $\delta_R^2 t$ (less than 20 bits) in BGV. On the other hand, the cost of integer multiplications in BFV appears to be significantly higher due to multiple basis extension operations. The above may suggest that the complexity of BFV could be lower than for BGV at small t , while more significant benefits of BGV are expected as t is increased, when the ratio of ℓ'/ℓ becomes smaller than 2, which corresponds to the typical value of $t = 2^{16} + 1$ used for CRT packing. One could argue that this is essentially due to the assumption that the computations modulo each CRT moduli are implemented on different machine words, which is typically true for practical implementations of homomorphic encryption. As a consequence, while BGV might be practically slower than BFV at small t for classical implementations, we stress that this is only due to the way the CRT representation is usually implemented and that BGV still has a lower theoretical complexity than BFV even for small plaintext moduli.

Remark 7. To reduce even further the computational cost of BGV, one could trunk some CRT moduli together in the same 64-bit machine word. This would allow one to divide the number of moduli ℓ' , and thus of NTTs, by a factor of 2 when the moduli are smaller than 30 bits ($t \approx 2^{11}$) and by a factor of 3 when they are smaller than 20 bits ($t \approx 2$).

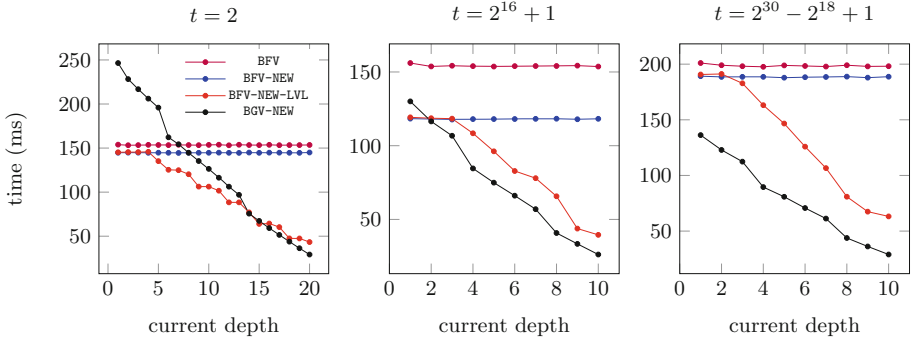


Fig. 1. Comparison of homomorphic multiplication runtimes for BFV and BGV variants at various depths as a function of plaintext modulus t . Hybrid key switching with 3 digits, i.e., $d_{\text{num}} = 3$, was used, and N was set to 2^{15} .

5.3 Software Implementation and Experimentation Setup

We implemented all variants of BFV and BGV in PALISADE v1.10.4. The evaluation environment was a commodity desktop computer system with an Intel(R) Core(TM) i7-9700 CPU @ 3.00 GHz and 64 GB of RAM, running Ubuntu 18.04 LTS. The compiler was g++ 9.3.0. All experiments were executed in the single-threaded mode.

PALISADE includes the implementation of both BEHZ and HPS variants of BFV. The runtime results and noise growth for both variants are roughly the same (as shown in Sect. 5.4). We chose the HPS variant as the main RNS variant for our BFV modifications due to its relative simplicity. We denote our modified BFV variant as BFV-NEW, our modified BFV variant with leveled multiplication as BFV-NEW-LVL, and our BGV variant as BGV-NEW. Note that our implementation does not trunk small CRT moduli in BGV for small values of t , i.e., it does not include the optimization suggested in Remark 7.

5.4 Performance Comparison

Figure 1 illustrates the comparison of homomorphic multiplication runtimes for the BFV and BGV variants developed in this work to the baseline for the prior state-of-the-art BFV implementation of the HPS variant [21]. The first major observation is that BFV-NEW-LVL outperforms BGV-NEW for small plaintext moduli (at least up to depth 20), while BGV-NEW runs significantly faster than BFV-NEW-LVL for intermediate and large plaintext moduli, i.e., $t = 2^{16} + 1$ and $t = 2^{30} - 2^{18} + 1$. This observation is in agreement with our theoretical complexity analysis in Sect. 5.2 since our implementation does not include the optimization suggested in Remark 7, i.e., small moduli are not trunked together. The second significant observation is that our best BFV variant, labeled as BFV-NEW-LVL, speeds up the runtime of deeper multiplications (depth-20 for $t = 2$ and depth-10 for higher t) by 3x-4x, as compared to the BFV baseline.

Table 2 shows the comparison of noise growth and runtimes for a binary tree computation ranging in multiplicative depth from 1 to 7. First, we want to point out that the noise growth and runtimes for the BEHZ and HPS variants are very close, with HPS having somewhat better runtime efficiency, which agrees well with the noise analysis in [6] and runtime comparison in [3]. In view of this, we chose HPS as the main variant for our BFV improvements (but similar gains can be expected for the BEHZ variant). Our next observation is that BGV has a slightly faster noise growth, as compared to all BFV variants, with the difference in noise increasing with depth (as predicted in Sect. 5.1). Note that the original BFV variants have somewhat higher noise (by almost constant number of bits) as compared to our BFV variants because they do not use the technique of encrypting with a slightly larger modulus Qp , followed with scaling by p . Our final observation is that BGV-NEW has a minor speed-up over the best BFV variant for the chosen plaintext modulus $t = 2^{16} + 1$. Note that the speed-up is observed only for this or higher plaintext moduli, with BFV-NEW-LVL becoming faster for $t = 2$ (see the full version for details). Tables in the appendix of the full version also show the more significant effect of $r_t(Q)$ on noise magnitude at larger plaintext moduli for the original BFV, as theoretically predicted in Sect. 3.

Table 2. Comparison of noise growth and runtimes of BFV and BGV variants for a benchmark computation $\prod_{i=1}^{2^k} x_i$. Hybrid key switching with 3 digits, i.e., $d_{\text{num}} = 3$, was used, t was set to $2^{16} + 1$, and $\lambda \geq 128$. Here, e denotes the current noise magnitude, $\log Q$, the size of the BFV ciphertext modulus, and $\log Q_L$, the equivalent ciphertext modulus in BGV without the last CRT modulus q_{L+1} .

k	Original BFV						Our BFV						Our BGV											
	params			BEHZ			HPS			params			BFV-NEW			BFV-NEW-LVL			params			BGV-NEW		
	$\log N$	$\log q_i$	$\log Q$	$\log e$	Time (s)	$\log e$	Time (s)	$\log N$	$\log q_i$	$\log Q$	$\log e$	Time (s)	$\log e$	Time (s)	$\log N$	$\log q_i$	$\log Q_L$	$\log e$	Time (s)	$\log N$	$\log q_i$	$\log Q_L$	$\log e$	Time (s)
1	13	31	62	45	0.011	44	0.01	13	59	59	36	0.004	35	0.004	13	33	58	34	0.005					
2	13	47	94	66	0.034	66	0.03	13	45	90	63	0.025	63	0.025	13	33	91	67	0.02					
3	14	43	129	102	0.24	103	0.21	14	41	123	95	0.19	96	0.18	13	33	124	100	0.063					
4	14	53	159	131	0.52	132	0.45	14	52	156	125	0.4	125	0.39	13	33	157	133	0.17					
5	14	48	192	158	1.41	161	1.2	14	47	188	155	1.07	155	1.04	14	34	196	171	0.8					
6	14	56	224	189	2.85	189	2.44	14	55	220	184	2.18	184	2.13	14	34	230	205	2.03					
7	14	51	255	221	7.61	220	6.51	14	50	250	214	5.98	214	5.73	14	34	264	239	4.86					

Table 3 illustrates the comparison of noise growth and runtimes for a polynomial evaluation benchmark. Our first observation is that BGV-NEW has a significantly higher noise than all BFV variants because the moduli q_i in this case require extra room for the additions at each level (the deepest level has the most significant effect on all q_i 's). BGV-NEW again has a minor advantage in terms of runtime as compared to our best BFV variant for $t = 2^{16} + 1$, but BFV-NEW-LVL becomes faster when we decrease t to smaller values (see the full version for details). Note that for $k = 8$, BGV-NEW has a smaller ring dimension than all BFV variants, which is an effect of the automated logic for hybrid key switching used in the implementation, rather than a result of better noise growth in BGV (since $\log Q$ in BFV is significantly smaller than $\log Q_L$ in BGV).

Table 3. Comparison of noise growth and runtimes of BFV and BGV variants for a benchmark computation $\prod_{i=0}^k a_i x^i$: $|a_i| < 16$. Hybrid key switching with 3 digits, i.e., $d_{\text{num}} = 3$, was used, t was set to $2^{16} + 1$, and $\lambda \geq 128$. Here, e denotes the current noise magnitude, $\log Q$, the BFV ciphertext modulus, and $\log Q_L$, the equivalent ciphertext modulus in BGV without the last CRT modulus q_{L+1} .

k	Original BFV						Our BFV						Our BGV						
	params			BEHZ		HPS		params			BFV-NEW		BFV-NEW-LVL		params			BGV-NEW	
	$\log N$	$\log q_t$	$\log Q$	$\log e$	Time (s)	$\log e$	Time (s)	$\log N$	$\log q_t$	$\log Q$	$\log e$	Time (s)	$\log e$	Time (s)	$\log N$	$\log q_t$	$\log Q_L$	$\log e$	Time (s)
2	13	34	68	41	0.012	40	0.01	13	32	64	35	0.009	36	0.009 s	13	38	68	38	0.007
4	13	50	100	76	0.034	76	0.03	13	48	96	67	0.026	67	0.025 s	13	38	107	74	0.024
8	14	45	135	106	0.25	107	0.22	14	43	129	100	0.19	100	0.18 s	13	39	148	116	0.061
16	14	56	168	138	0.53	138	0.46	14	54	162	130	0.4	130	0.33 s	14	41	197	163	0.28
32	14	50	200	166	1.43	167	1.22	14	49	196	161	1.1	161	0.78 s	14	42	244	208	0.61
48	14	58	232	197	2.16	198	1.85	14	57	228	191	1.66	190	1.22 s	14	42	286	251	1.07
64	14	58	232	199	2.89	199	2.48	14	57	228	191	2.22	191	1.54 s	14	43	293	256	1.27

6 Concluding Remarks

Our theoretical analysis and experimental results show that the modified BFV variant has somewhat better noise growth than BGV for all plaintext moduli, though previous results suggested that BGV has a better noise growth than BFV for larger plaintext moduli [13, 14]. This result is mainly due to our modification of the BFV encryption procedure. The other major conclusion is that, when the moduli of BGV are not trunked together, BFV is significantly faster for small plaintext moduli, e.g., $t = 2$, with BGV becoming faster as the plaintext modulus is increased.

The variant of BGV presented in this paper was mainly motivated by improving the usability of the scheme, which is known to be more challenging for use than BFV. From this perspective, this BGV variant is as easy to use as the implementation of BFV in PALISADE. However, the usability also has some performance cost, e.g., we have to choose the size of CRT moduli more conservatively. It would be interesting to examine how the performance of our BGV variant compares to the BGV design with dynamic noise estimation, which is implemented in HELib. It would not be fair to compare the PALISADE implementation directly with the HELib implementation as one would mostly observe the effect of differences in the efficiency of primitive ring operations, such as NTTs, rather than the differences between the BGV variants. For a fair comparison, a PALISADE implementation of the dynamic-noise BGV variant would be needed. Another potential improvement for BGV is to consider the idea of trunking multiple small CRT moduli mentioned in Remark 7. We plan to examine both ideas in our future work.

Acknowledgments. Andrey Kim and Yuriy Polyakov’s NJIT work was supported in part by the Defense Advanced Research Projects Agency (DARPA) and the US Navy SPAWAR Systems Center Pacific (SSCPAC) under Contract Number N66001-17-1-4043 and the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via Contract No. 2019-1902070006. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied,

of the Department of Defense, ODNI, IARPA, or the U.S. Government. Vincent Zucca's KU Leuven work was supported in part by the Research Council KU Leuven grant C14/18/067, CyberSecurity Research Flanders with reference number VR20192203, and the IARPA HECTOR project under the solicitation number IARPA-BAA-17-05. We also thank Charlotte Bonte for a careful review of the first version of the paper, her feedback, and fruitful discussions that helped us to improve the paper.

References

1. Lattigo v2.1.1, December 2020. <http://github.com/ldsec/lattigo>. ePFL-LDS
2. PALISADE Lattice Cryptography Library (release 1.10.6), December 2020. <https://palisade-crypto.org/>
3. Al Badawi, A., Polyakov, Y., Aung, K.M.M., Veeravalli, B., Rohloff, K.: Implementation and performance evaluation of RNS variants of the BFV homomorphic encryption scheme. *IEEE Trans. Emerg. Top. Comput.* **9**(2), 941–956 (2021). <https://doi.org/10.1109/TETC.2019.2902799>
4. Albrecht, M., Chase, M., Chen, H., et al.: Homomorphic encryption security standard. Technical report, HomomorphicEncryption.org, Toronto, Canada, November 2018
5. Bajard, J.-C., Eynard, J., Hasan, M.A., Zucca, V.: A full RNS variant of FV like somewhat homomorphic encryption schemes. In: Avanzi, R., Heys, H. (eds.) SAC 2016. LNCS, vol. 10532, pp. 423–442. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-69453-5_23
6. Bajard, J.C., Eynard, J., Martins, P., Sousa, L., Zucca, V.: Note on the noise growth of the RNS variants of the BFV scheme. *Cryptology ePrint Archive*, Report 2019/1266 (2019). <https://eprint.iacr.org/2019/1266>
7. Bos, J., et al.: CRYSTALS - kyber: a CCA-secure module-lattice-based KEM. In: 2018 IEEE European Symposium on Security and Privacy (EuroS P), pp. 353–367 (2018). <https://doi.org/10.1109/EuroSP.2018.00032>
8. Brakerski, Z.: Fully homomorphic encryption without modulus switching from classical GapSVP. In: Safavi-Naini, R., Canetti, R. (eds.) CRYPTO 2012. LNCS, vol. 7417, pp. 868–886. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-32009-5_50
9. Brakerski, Z., Gentry, C., Vaikuntanathan, V.: (leveled) fully homomorphic encryption without bootstrapping. *ACM Trans. Comput. Theory (TOCT)* **6**(3), 1–36 (2014)
10. Brakerski, Z., Vaikuntanathan, V.: Fully homomorphic encryption from ring-LWE and security for key dependent messages. In: Rogaway, P. (ed.) CRYPTO 2011. LNCS, vol. 6841, pp. 505–524. Springer, Heidelberg (2011). https://doi.org/10.1007/978-3-642-22792-9_29
11. Cheon, J.H., Kim, A., Kim, M., Song, Y.: Homomorphic encryption for arithmetic of approximate numbers. In: Takagi, T., Peyrin, T. (eds.) ASIACRYPT 2017. LNCS, vol. 10624, pp. 409–437. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-70694-8_15
12. Chillotti, I., Gama, N., Georgieva, M., Izabachène, M.: Faster fully homomorphic encryption: bootstrapping in less than 0.1 seconds. In: Cheon, J.H., Takagi, T. (eds.) ASIACRYPT 2016. LNCS, vol. 10031, pp. 3–33. Springer, Heidelberg (2016). https://doi.org/10.1007/978-3-662-53887-6_1

13. Costache, A., Smart, N.P.: Which ring based somewhat homomorphic encryption scheme is best? In: Sako, K. (ed.) CT-RSA 2016. LNCS, vol. 9610, pp. 325–340. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-29485-8_19
14. Costache, A., Laine, K., Player, R.: Evaluating the effectiveness of heuristic worst-case noise analysis in FHE. In: Chen, L., Li, N., Liang, K., Schneider, S. (eds.) ESORICS 2020. LNCS, vol. 12309, pp. 546–565. Springer, Cham (2020). https://doi.org/10.1007/978-3-030-59013-0_27
15. Ducas, L., Micciancio, D.: FHEW: bootstrapping homomorphic encryption in less than a second. In: Oswald, E., Fischlin, M. (eds.) EUROCRYPT 2015. LNCS, vol. 9056, pp. 617–640. Springer, Heidelberg (2015). https://doi.org/10.1007/978-3-662-46800-5_24
16. Fan, J., Vercauteren, F.: Somewhat practical fully homomorphic encryption. IACR Cryptol. ePrint Arch. **2012**, 144 (2012)
17. Gentry, C.: Fully homomorphic encryption using ideal lattices. In: Proceedings of the Forty-First Annual ACM Symposium on Theory of Computing, pp. 169–178 (2009)
18. Gentry, C., Halevi, S.: Implementing gentry’s fully-homomorphic encryption scheme. In: Paterson, K.G. (ed.) EUROCRYPT 2011. LNCS, vol. 6632, pp. 129–148. Springer, Heidelberg (2011). https://doi.org/10.1007/978-3-642-20465-4_9
19. Gentry, C., Halevi, S., Smart, N.P.: Fully homomorphic encryption with polylog overhead. In: Pointcheval, D., Johansson, T. (eds.) EUROCRYPT 2012. LNCS, vol. 7237, pp. 465–482. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-29011-4_28
20. Gentry, C., Halevi, S., Smart, N.P.: Homomorphic evaluation of the AES circuit. In: Safavi-Naini, R., Canetti, R. (eds.) CRYPTO 2012. LNCS, vol. 7417, pp. 850–867. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-32009-5_49
21. Halevi, S., Polyakov, Y., Shoup, V.: An improved RNS variant of the BFV homomorphic encryption scheme. In: Matsui, M. (ed.) CT-RSA 2019. LNCS, vol. 11405, pp. 83–105. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-12612-4_5
22. Halevi, S., Shoup, V.: Bootstrapping for HELib. Cryptology ePrint Archive, Report 2014/873 (2014). <https://eprint.iacr.org/2014/873>
23. Halevi, S., Shoup, V.: Design and implementation of HELib: a homomorphic encryption library. Cryptology ePrint Archive, Report 2020/1481 (2020)
24. Han, K., Ki, D.: Better bootstrapping for approximate homomorphic encryption. In: Jarecki, S. (ed.) CT-RSA 2020. LNCS, vol. 12006, pp. 364–390. Springer, Cham (2020). https://doi.org/10.1007/978-3-030-40186-3_16
25. Kim, A., Papadimitriou, A., Polyakov, Y.: Approximate homomorphic encryption with reduced approximation error. Cryptology ePrint Archive, Report 2020/1118 (2020). <https://eprint.iacr.org/2020/1118>
26. Kim, M., Song, Y., Li, B., Micciancio, D.: Semi-parallel logistic regression for GWAS on encrypted data. BMC Med. Genomics **13**(7), 1–13 (2020)
27. Lepoint, T., Naehrig, M.: A comparison of the homomorphic encryption schemes FV and YASHE. In: Pointcheval, D., Vergnaud, D. (eds.) AFRICACRYPT 2014. LNCS, vol. 8469, pp. 318–335. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-06734-6_20
28. Micciancio, D., Polyakov, Y.: Bootstrapping in FHEW-like cryptosystems. Cryptology ePrint Archive, Report 2020/086 (2020). <https://eprint.iacr.org/2020/086>
29. Microsoft SEAL (2020). <https://github.com/Microsoft/SEAL>
30. Smart, N.P., Vercauteren, F.: Fully homomorphic SIMD operations. Des. Codes Crypt. **71**(1), 57–81 (2012). <https://doi.org/10.1007/s10623-012-9720-4>