

# Revisiting Key Decomposition Techniques for FHE: Simpler, Faster and More Generic

M.G. Belorgey<sup>1</sup>, S. Carpov<sup>1</sup>, N. Gama<sup>2</sup>, S. Guasch<sup>2</sup>, D. Jetchev<sup>1</sup>

<sup>1</sup> Inpher  
<sup>2</sup> SandboxAQ

**Abstract.** Ring-LWE based homomorphic encryption computations in large depth use a combination of two techniques: 1) decomposition of big numbers into small limbs/digits, and 2) efficient cyclotomic multiplications modulo  $X^N + 1$ . It was long believed that the two mechanisms had to be strongly related, like in the full-RNS setting that uses a CRT decomposition of big numbers over an NTT-friendly family of prime numbers, and NTT over the same primes for multiplications. However, in this setting NTT was the bottleneck of all large-depth FHE computations. A breakthrough result from Crypto'2023 [22] managed to overcome this limitation by introducing a second gadget decomposition and by showing that it indeed shifts the bottleneck and renders the cost of NTT computations negligible compared to the rest of the computation. In this paper, we extend this result (far) beyond the Full-RNS settings and show that we can completely decouple the big number decomposition from the cyclotomic arithmetic aspects. As a result, we get modulus switching/rescaling for free, and the memory footprint for storing re-linearization keys across different levels is considerably lower compared to the CRT-based counterparts, by typically a factor  $\ell/3$  where  $\ell$  is the deepest level of multiplication depth supported. We verify both in theory and in practice that the performance of key-switching, external and internal products and automorphisms using our representation are similar or faster than the one achieved by [22], and we discuss the high impact of these results for people who work on low-level or hardware optimizations as well as the benefits of the new parametrizations for people currently working on compilers for FHE. We even manage to lower the running time of the gate bootstrapping of TFHE by eliminating 12.5% of its FFTs.

## 1 Introduction

Homomorphic encryption allows computations on encrypted data without decrypting it. Since the first fully homomorphic encryption (FHE) scheme introduced by Gentry in [17], several improvements, implementations and new designs have been proposed. Most of them are based on the ring version of the Learning-With-Error problem (RingLWE) defined in [27]. Among the most popular schemes today are BFV [16], BGV [6], CKKS [9], FHEW [15] and TFHE [10,11]. Subsequent work [5] has been done on the interoperability of TFHE with the other two schemes (BFV and CKKS).

Most of the basic building blocks of the practical homomorphic encryption schemes (e.g., homomorphic multiplication and relinearization, key switching, automorphisms or bootstrapping) reduce to efficient *homomorphic external products*. The latter are basic operations that are combinations of decompositions of one of the inputs into higher-dimensional tensor with entries of small norms (*gadget decompositions*) and polynomial multiplications in order to control the accumulation of noise in the ciphertexts and ensure correct decryption (see, e.g., [5, Defn.1] and [11, Defn.3.12] for typical examples of external product).

A major line of research has been undertaken on using Residue Number Systems (RNS) in the acceleration of the basic underlying polynomial multiplication operations used for computing external products [18]. The full-RNS approach of [4], originally introduced to speedup BFV operations (see also [21]) and subsequently applied to other schemes too such as CKKS [8], allows for manipulating large numbers using smaller machine-size moduli and expressing the homomorphic products (internal products) as well as the gadget decompositions of the BFV and CKKS in the NTT domain. Subsequent hardware acceleration efforts have been made to accelerate RNS-variants of NTT for GPU architectures [26,31] and FPGA architectures [14,25,28,30].

One can thus conceptually differentiate two important aspects of optimization of FHE operations: 1) efficient large integer arithmetic (frontend aspect) that is addressed via mathematical techniques such as gadget decompositions; 2) cyclotomic arithmetic (backend aspects) - the need for efficient arithmetic on the backend to support the basic cyclotomic operations - the choice of floating-point arithmetic or integer arithmetic on the backend. Unfortunately, in most practical implementations, the two aspects are often coupled together, thus, limiting the flexibility for optimization techniques. For instance, by nature, the full-RNS approach to gadget decomposition requires modular arithmetic operations on the backend and thus, is bound to NTT leaving little room for really fast FFT-favored operations (such as the optical FFT approach undertaken in [23]).

Note that there are essentially two classical approaches to 1): decomposition in base  $2^K$  and CRT decomposition over primes of  $K$  bits. Assuming that single elementary operations operate natively over  $K$ -bit numbers (e.g., 32-bit or 64-bit integers or floating point numbers with 52-bits of mantissa), the fact that CRT arithmetic is exempt from any carry propagation yields very efficient parallel additions and multiplications of  $\ell$   $K$ -bit numbers in  $O(\ell)$  operations, whereas the base- $2^K$  counterpart would struggle between  $O(\ell \log \ell)$  and  $O(\ell^2)$  multiplications and force these elementary operations in sequential mode. For this reason, a natural choice for homomorphic encryption was to prefer CRT representations for efficient implementations of CKKS and BFV.

A recent breakthrough has been made by Kim et al. towards decoupling 1) and 2) by an approach based on an auxiliary gadget decomposition that replaces expensive modular arithmetic in the computation of the external product by pure arithmetic with small integers via gadget decompositions based on RNS [22].

Note that the RNS-type approach to 1) often has the following drawbacks:

1. *Arithmetic modulo different primes.* The arithmetic has to be carried out modulo different prime numbers of same size, even if they are selected to be as *friendly* as possible, numbers modulo which arithmetic is really efficient are rare (e.g., Mersenne primes, Fermat primes). There is also a small loss from hardware aspects such as the necessity to have one hardware multiplier per prime.
2. *Unfeasibility of bit extraction/coarse noise granularity.* Bit extraction is not possible in CRT domain modulo a product of small primes, the best granularity of HE noise levels is those of the primes: e.g., 16 to 40 bits. As explained in the analysis of the Hecate compiler [24] this impacts negatively the optimal use of homomorphic budget for CRT-based CKKS - one often needs to re-scale the noise to the nearest available level before attempting a homomorphic product after an operation producing less noise (such as a sum, a small linear combination or a trace).
3. *Sub-optimal key material for relinearization, bootstrapping or key switching.* The necessity to store a different key for each multiplicative level yields a  $O(\ell)$  overhead in key material.
4. *Sub-optimal memory usage with respect to truncations.* Restrictions from high-precision representations to low-precision representations are not simple truncations of limbs (as in the natural base- $2^K$  representation), but require a separate RNS representation. This is a significant memory overhead.

The above-mentioned carry-less nature and efficiency of parallel CRT multiplications would have been perfect if most homomorphic operations were actual ring operations. However, internal ring multiplications over large numbers never come standalone: first of all, in a BFV or CKKS internal product they always come in packs of  $N$  (giving a  $O(N\ell)$  complexity to multiply two polynomials modulo  $X^N + 1$ ), and second, they are always followed by an external product that comprises at least  $\ell$  NTTs (so at least  $O(\ell N(\ell + \log N))$  if we use the latest double gadget decomposition of [22] and  $O(\ell^2 N \log N)$  before). In other words, since  $N \gg \ell$ : typically,  $N$  is between 1024 and 65536 and  $\ell$  is lower than 100, the complexity of an internal BFV or CKKS product remains bounded both asymptotically and practically by the complexity of its underlying external product. This means that as long as we make the external product more efficient, we have some margin to degrade the complexity of the internal multiplication of polynomials up to  $O(N\ell^2)$  without witnessing any negative effect on the running time of CKKS and BFV operations.

It is thus desirable to extend the auxiliary gadget decomposition approach of 1) to the natural  $2^K$ -base large integer representation to overcome the above drawbacks.

### **Our contributions.**

In this work, this is exactly what we achieve - we show how to adapt the *auxiliary gadget decomposition* of [22] to the natural base- $2^K$  representation of numbers and thus, obtain an external product that is simpler and faster.

Our starting point is a common plaintext representation space for all the practical RingLWE-based schemes (BFV, CKKS and TFHE) described in [5], namely the  $\mathcal{R}$ -module  $\mathbb{T}_N[X] := \mathbb{R}[X]/\langle X^N + 1 \rangle / \mathcal{R}$  where  $\mathcal{R} = \mathbb{Z}[X]/\langle X^N + 1 \rangle$  (also thought of as *formal* polynomials of degree  $N - 1$  whose coefficients are in the torus  $\mathbb{T} = \mathbb{R}/\mathbb{Z}$ ). Inspired by the classical Schoenhage–Strassen algorithm for large integer multiplication based on FFT [29], in section 3, we propose a bivariate integral polynomial representation with small coefficients of these formal polynomials where the evaluation on  $Y = 2^{-K}$  ( $K$  is the limb size) yields a sufficiently accurate approximation of the toric coefficients. Introducing the auxiliary variable implicitly corresponds to having the second gadget decomposition: on one hand, it yields the natural decomposition in base  $2^K$  of toric elements (evaluation at  $Y = 2^{-K}$ ); on the other hand, it provides a plenitude of choice of bivariate polynomials to approximate the given toric polynomial. A careful normalization and reduction yields a representative with small coefficients - the bivariate polynomial representation that is used for efficient external product (see subsection 3.2), thus making it specifically tuned for TRLWE operations.

Since the space of TRLWE ciphertexts is  $\mathbb{T}_N[X]^2$ , we express the external product in terms of multiplications of bivariate polynomials with small integer coefficients. It is this reduction that decouples 1) and 2) above - the multiplication can be performed using either FFT or NTT backend arithmetic which we explain in section 6, thus, enabling optimal use of different architectures for hardware optimization.

A major advantage of our representation is that it addresses challenges 1–4. Indeed, our representation supports efficient sums, normalizations, reductions modulo  $\mathbb{Z}$ , left and right shifts in  $O(\ell)$  and products in  $O(\ell \log \ell)$  (only bilinear expressions are needed, we do not need higher degree terms in FHE operations!), bit-decompositions and bit-extraction. In addition, a key advantage is the *prefix property* of the natural base- $2^K$  representation yields a memory-efficient way to pass from a higher-precision representation to a lower precision representation. It provides a single-bit noise granularity, thus, optimizing the key material overhead and enabling efficient compositions as proposed in [5]. The asymptotic complexity matches the best available full-RNS algorithms for both external products of TFHE, and also internal products CKKS and BFV.

The fact that carry propagation must be treated sequentially is mitigated by the parallelism induced by the fact that the operations in our external product occur in a SIMD manner over  $N$  elements each time. We increase the complexities of the internal multiplications from  $O(N\ell)$  to  $O(N\ell \log \ell)$ , keeping them negligible compared to an external product, and hence, due to the speed-up of the latter ones, also accelerate internal products of CKKS and BFV, so that at the end of the day, by not using any of the CRT of full-RNS techniques, and instead, applying the double-gadget decomposition to old-school base- $2^K$  representations directly, we end up speeding up all the existing homomorphic operations of TFHE, CKKS and BFV.

Our practical experiments and benchmarks support the theoretical evidence: first, Table 1 provides evidence for the number of elementary operations (SIMD

products and DFTs) in an external product computation in the context of the various FHE building blocks (or more specifically, in the more efficient approach via `HalfRGSW` introduced in subsection 3.2). The experiments demonstrate that for FHE parameters  $N = 65536$ ,  $\ell \sim 100$  and  $K \sim 20$ , the cost of SIMD products largely dominates the cost of DFTs, an already strong indication of the parallelization-friendly nature of our approach. Second, we compare our bivariate representation (using both FFT and NTT) to the approach from [22] to demonstrate that, prior to any parallelization or hardware acceleration efforts, we get a comparable performance (see Table 4). Finally, in section 4 we even manage to lower the running time of the gate bootstrapping of TFHE by eliminating 12.5% of its FFTs.

### Future directions.

One of the most important contributions of our work is the flexibility it provides for future hardware acceleration efforts. Everything that we have done in this work is single-core optimizations.

First of all, the normalization and reduction Algorithm 1 (and its extension presented in the Appendix) allows for fast AVX and GPU-friendly implementations. The bivariate representation yields several advantages making the parallelization very natural. First, as mentioned, it is directly inspired by the Schoenage–Strassen algorithm for large integer multiplication based on NTT with complexity  $O(n \log n \log \log n)$ , where  $n$  is the number of binary digits of the inputs [29]. In fact, the case of large integer multiplication is recovered directly from our bivariate representation by setting  $N = 1$  via the evaluation map<sup>3</sup>  $\mathcal{R}[Y] \rightarrow \mathcal{R}$  given by  $f(Y) \mapsto f(2^K)$ . Note that the Schoenage–Strassen is an NTT version of the DFT-based fast polynomial multiplication. Even if the original Schoenage–Strassen algorithm is not completely easy to parallelize, the large degree  $N$  on the variable  $X$  makes our bivariate polynomial multiplication friendlier for parallel architectures.

Second, a recursive version of DFT or more generally, a Cooley–Tukey transforms [12] allows for parallelization on the  $Y$ -aspect. The naïve multiplication, single iteration of Karatsuba or DFT algorithms discussed here can be performed in-place. The well-known *in-place* and *in-order* variants of the Cooley–Tukey transform together with the prefix property opens the door for highly optimized memory usage as well.

As already mentioned, the trade-off of using FFT and NTT can benefit from the various hardware acceleration initiatives such as GPU acceleration [26], FPGA acceleration [25,30], optical computing [23] (particularly favorable for the FFT approach), ASICs [13] as well as the outcome of the DPRIVE program [1], especially the recent work [7]. Of particular interest considering the trade-off of compute cost vs. performance is the FPGA approach.

---

<sup>3</sup> Note the sign difference in the exponent; it is due to the fact that multiplication is not well-defined on the torus but only on the integers.

Finally, our work opens up the possibility for an efficient implementation of the scheme switching framework **Chimera** [5] that allows to mix **TFHE**, **CKKS**, **BFV** arithmetic using continuous homomorphic levels. The original proposal was suggesting the use of large numbers fixed-point arithmetic in order to perform efficient additions, multiplication and bit-extraction. Unfortunately, none of the available representations of big numbers in **GMP**, **MPFR** are sufficiently efficient in practice, and a CRT representation would not allow efficient bit-extraction either, thus leading to the same setbacks as Full-RNS. For these reasons, the early attempts to implement **Chimera** were, thus, slow in practice for levels  $\geq 3$ .

In conclusion, our approach supports the following general principle: instead of one trying to find a numerical representation (e.g., RNS representation) with a gadget decomposition that is not too convoluted, it is more natural to start from the gadget decomposition (e.g., base- $2^K$  representation) and take the output of the gadget decomposition directly as the numerical representation.

## 2 Preliminaries

### Notation

Let  $\mathbb{T} = \mathbb{R}/\mathbb{Z}$  be the real torus. For a ring  $A$  (e.g.,  $A = \mathbb{Z}, \mathbb{R}, \mathbb{C}$ ), let  $A_N[X] := A[X]/\langle X^N + 1 \rangle$  be the ring of polynomials modulo  $X^N + 1$  with coefficients in  $A$  and  $N$  a power of 2. In particular, let  $\mathcal{R} = \mathbb{Z}_N[X]$ , which is also the ring of integers of the cyclotomic field  $\mathbb{Q}(\zeta_{2N})$ .

Let  $\mathbb{T}_N[X] = \mathbb{R}_N[X]/\mathcal{R}$  (a.k.a  $\mathbb{R}[X] \bmod X^N + 1 \bmod \mathbb{Z}$ ) which we view as an  $\mathcal{R}$ -module (it has no ring structure). Elements of this module can formally be represented using notations borrowed from polynomials, i.e.,  $a(X) = \sum_{i=0}^{N-1} a_i X^i$  where  $a_i \in \mathbb{T}$ . Since the coefficient space  $\mathbb{T}$  is not a ring, we cannot evaluate these polynomials over any non-integer value, nor multiply two such elements together. However, the notation allows coefficient-wise addition/subtraction of polynomials, projection from and lifts to the continuous ring  $\mathbb{R}_N[X]$ , and most importantly, we often refer to the module action of  $\mathcal{R}$  on  $\mathbb{T}_N[X]$  as an *external multiplication* of  $a$  by integer polynomial  $u(x) = \sum_{i=0}^{N-1} u_i X^i \in \mathcal{R}$  where  $u_i \in \mathbb{Z}$  via the natural Cauchy product formula:  $(u \cdot a)(X) = \sum_{i=0}^{2N-2} \left( \sum_{j=0}^i u_j \cdot a_{i-j} \right) X^i \bmod X^N + 1$ .

We denote by  $\|\cdot\|_p$  and  $\|\cdot\|_\infty$  the standard norms for scalars and vectors over real field or over the integers. By extension, the norms  $\|P(X)\|_p$  and  $\|P(X)\|_\infty$  of a real or integer polynomial  $P$  are the norms of their coefficient vectors. The norm of an element of  $\mathbb{T}_N[X]$  is the norm of its centered lift in  $\mathbb{R}_N[X]$  with coefficients in  $[-1/2, 1/2)$ .

### 2.1 TRLWE

TRLWE encrypts elements of a subset of the  $\mathcal{R}$ -module  $\mathbb{T}_N[X]$  called the plaintext space, and which can be finite (**TFHE**), a discrete subgroup (**BFV**) or a continuous set of small elements (**CKKS**). The ciphertext of  $\mu$  are of the form  $c = (a, b =$

$s \cdot a + \mu + e$ ), where  $s \in \mathcal{B}$  is the secret key,  $a \in_R \mathbb{T}_N[X]$  is uniformly random and the  $e \in \mathbb{T}_N[X]$  is chosen randomly from an error distribution with mean zero and suitably chosen standard deviation. The decryption procedure starts by evaluating the phase function on the ciphertext, that is, by computing  $\varphi_s(a, b) := b - s \cdot a$ . Since  $\varphi_s(a, b) = \mu + e$ , the plaintext  $\mu$  is the mean of the phase (the mean being computed over the random coins in the generation of the noise  $e$  during encryption).

- **Parameters:** A security level  $\lambda$  and a minimal noise parameter  $\alpha$ .
- **KeyGenTRLWE**( $\lambda, \alpha$ ): A uniformly random binary key  $s \in \mathcal{B} \subset \mathcal{R}$ .
- **EncTRLWE**( $\mu, s, \alpha$ ): Choose a uniform random element  $a \in \mathbb{T}_N[X]$  and a small Gaussian error  $e \in \mathbb{T}_N[X]$  and return  $c = (a, b = s \cdot a + \mu + e)$
- **DecTRLWE**( $c, s$ ): Compute the phase  $\varphi_s(c)$  and round  $\varphi_s(c)$  to the nearest point in the plaintext space (when it is discrete), or returns a close approximation (CKKS).

## 2.2 Approximate gadget decompositions and TRGSW ciphertexts

A common challenge with various FHE constructions (e.g., GSW) is the control on the accumulation of noise after homomorphic operations. A general technique from lattice-based cryptography to address this problem is a concept known as *flattening gadget* or *gadget decomposition*, that is, a map that transforms a ciphertext into a higher dimensional vector with small  $\|\cdot\|_\infty$  while preserving some linear-algebraic properties.

Two classical examples of gadget decompositions are numerical base representations and RNS representations. Since most of the FHE literature uses gadget decompositions on integral values, the decomposition themselves are exact. A notable exception to this principle is TFHE where one uses floating point arithmetic and therefore, an approximate gadget decomposition. Below we recall the basic definition (see also [11, Defn.3.6]):

### Definition 2.1 (approximate gadget decomposition for TRLWE samples).

We say that an algorithm  $Decomp_{H, \beta, \epsilon}$  is a valid gadget decomposition on gadget  $H \in (\mathbb{T}_N[X]^2)^{2\ell}$ , quality (or  $\|\cdot\|_\infty$ -bound),  $\beta$  and precision  $\epsilon > 0$  if for any input  $v \in \mathbb{T}_N[X]^2$ , it outputs an element  $u \in \mathcal{R}^{2\ell}$ ,  $\|u\|_\infty \leq \beta$  such that  $\|u \cdot H - v\|_\infty < \epsilon$ .

As already mentioned in [11], there is a canonical approximate gadget decompositions coming from numerical base representations for any base  $B_g$ :

$$H^T = \begin{pmatrix} B_g^{-1} & \dots & B_g^{-\ell} & 0 & \dots & 0 \\ 0 & \dots & 0 & B_g^{-1} & \dots & B_g^{-\ell} \end{pmatrix}$$

The above notion is important for two main reasons: 1) allows us to define TRGSW ciphertexts; 2) allows us to define an external product between TRGSW ciphertexts and TRLWE ciphertexts.

We first recall the TRGSW ciphertexts: the TRGSW encrypts elements of the ring  $\mathcal{R}$  with bounded infinity norm. Intuitively, the idea of the original GSW scheme

[19] is to encrypt a plaintext  $\mu$  into a ciphertext that is a matrix  $C_\mu$  such that the secret key  $\mathbf{s}$  is an approximate eigenvector of  $C_\mu$  with eigenvalue  $\mu$ . Such an encryption scheme leads to a natural homomorphic addition and multiplication operations that are simply matrix additions and matrix multiplications. The naïve idea does not quite work without gadget decompositions since one cannot control the propagation of the errors in the approximate eigenvalues under homomorphic multiplications. The schemes FHEW and TFHE use a ring-variant of the original scheme, the last one with an approximate gadget decomposition:

**Definition 2.2 (TRGSW ciphertexts).** *Let  $H \in (\mathbb{T}_N[X]^2)^{2\ell}$  be a gadget and let  $\mu \in \mathcal{R}$  be a plaintext with bounded  $\ell_\infty$ -norm. The space of valid TRGSW ciphertexts for  $\mu$  as  $\text{TRGSW}(\mu) := \{Z + \mu H\}$ , where each element of  $Z \in (\mathbb{T}_N[X]^2)^{2\ell}$  is a valid TRLWE ciphertext of zero.*

### 2.3 External products, relinearization keys and internal products

For the purposes of the current work, we will need a relatively uniform treatment of the various RLWE-based FHE frameworks, most notably, BFV, CKKS and TFHE, as well as their internal products.

It is explained in [5, §2.5] (as well as the particular discussion of BFV and CKKS in Sections 3 and 4, respectively, of *loc.cit.*) how to uniformize the plaintext spaces for all these schemes and view them as various subgroups of the  $\mathcal{R}$ -module  $\mathbb{T}_N[X]$ . Once this is done, the homomorphic internal products of BFV and CKKS are expressed in terms of the following basic primitive, the *external product*, a major operation of interest in the current work:

**Definition 2.3 (TFHE external product).** *Given a flattening gadget  $H$  and an approximate gadget decomposition  $\text{Decomp}_{H,\beta,\epsilon}$ , the external product in TFHE is a map*

$$\square: \text{TRGSW} \times \text{TRLWE} \rightarrow \text{TRLWE}$$

defined by

$$C \square c := \text{Decomp}_{H,\beta,\epsilon}(c) \cdot C,$$

where  $C \in \text{TRGSW}(\mu_1)$  and  $c \in \text{TRLWE}(\mu_2)$ .

One can show [11, Thm.3.13] that under certain specific noise conditions, the above external product  $C \square c$  is a valid ciphertext for  $\mu_1\mu_2$ .

Once we have defined this primitive, we use it to express the various internal products, the major idea being the concept of *relinearization* and *relinearization keys* - we refer the reader to [5, pp.325–326] for a more detailed explanation. The important point is that by using extra key material (relinearization key), one can express the internal product of the BFV scheme in terms of the above external product. More specifically, if we define<sup>4</sup> the relinearization key

<sup>4</sup> There are many equivalent ways to define the relinearization key. Here, we use an external product and a zero term in the TRLWE ciphertext, which can be propagated in the algorithm and simplified, the original references of Fan–Vercauteren [16] present a key-switching that substitute a key  $s^2$  by  $s$ .

as  $\text{RK} = \text{TRGSW}(s)$ , then for  $\star \in \{\text{BFV}, \text{CKKS}\}$  the homomorphic internal product  $\boxtimes: \mathbb{T}_N[X]^2 \times \mathbb{T}_N[X]^2 \rightarrow \mathbb{T}_N[X]^2$  of two ciphertexts  $(a_1, b_1)$  and  $(a_2, b_2)$  of plaintexts  $\mu_1$  and  $\mu_2$ , respectively becomes

$$(a_1, b_1) \boxtimes (a_2, b_2) = (C_1, C_0) - \text{RK} \boxdot (C_2, 0),$$

where  $C_0 = b_1 \otimes_\star b_2$ ,  $C_1 = a_1 \otimes_\star b_2 + a_2 \otimes_\star b_1$ ,  $C_2 = a_1 \otimes_\star a_2$ , the important point being that this ciphertext is a valid encryption of the plaintext  $\mu_1 \mu_2$ . Here,  $\otimes_\star: \mathbb{T}_N[X] \times \mathbb{T}_N[X] \rightarrow \mathbb{T}_N[X]$  indicates a certain product map (depending on the scheme) whose restriction to the plaintext subgroup of  $\mathbb{T}_N[X]$  yields the plaintext product.

For instance, if for  $\bullet \in \mathbb{T}_N[X]$ ,  $\tilde{\bullet} \in \mathbb{R}[X]/\langle X^N + 1 \rangle$  denotes the unique lift with coefficients in the interval  $[-1/2, 1/2)$  then the BFV product  $\otimes_{\text{BFV}}$  with plaintext modulo  $p$  (Montgomery product) is:

$$\otimes_{\text{BFV}}: \mathbb{T}_N[X] \times \mathbb{T}_N[X] \rightarrow \mathbb{T}_N[X], \quad (u, v) \mapsto p \cdot \tilde{u} * \tilde{v}. \quad (1)$$

Similarly, CKKS plaintext products at input level  $L_{\text{in}}$  are:

$$\begin{aligned} \otimes_{\text{CKKS}}: \mathbb{T}_N[X] \times \mathbb{T}_N[X] &\rightarrow \mathbb{T}_N[X], \\ (u, v) &\mapsto 2^{L_{\text{in}}} \cdot \text{Round}_{2^{-L_{\text{in}}}}(\tilde{u}) * \text{Round}_{2^{-L_{\text{in}}}}(\tilde{v}) \end{aligned} \quad (2)$$

In both formulas for (BFV and CKKS product), the symbol  $\cdot$  is the external product by an integer, and  $*$  represents multiplication in  $\mathbb{R}_N[X]$ : the input coefficients are first lifted to the real interval  $[-1/2, 1/2)$  and in the second case, also rounded to the nearest exact multiple of  $2^{-L}$ . None of these functions is an actual products over the entire space - one needs to restrict the inputs to the subgroups  $p^{-1}\mathcal{R}/\mathcal{R}$  and  $2^{-L}\mathcal{R}/\mathcal{R}$  of  $\mathbb{T}_N[X]$ , respectively to obtain products. Yet, the function  $\otimes_{\text{BFV}}$  has the property that if  $u$  and  $v$  are close from  $i/p$  and  $j/p$  where  $i$  and  $j$  are integers, their product is close to  $(ij \bmod p)/p$ , which is handy to encode plaintext arithmetic modulo  $p$ . The product  $\otimes_{\text{CKKS}}$  has the property that if  $u$  and  $v$  are at distance  $\leq 2^{L+1}$  away from  $i/2^L$  and  $j/2^L$  where  $i$  and  $j$  are smaller than  $2^\rho$ , then their product is at distance  $2^{L-\rho}$  away from  $ij/2^L$ , which is good to encode fixed-point number arithmetic on  $\rho$ -bits numbers.

## 2.4 Table of symbols, orders of magnitudes

When reading the paper, different complexities shall intervene in the different theorems. It is important to keep in mind the difference in order of magnitudes, as for instance: having an arithmetic in  $O(N^2)$  is prohibitive, and we must absolutely stick to  $O(N \log N)$ , however  $O(\ell^2)$  is perfectly realistic in some scenario.

The parameters  $N, K$  and  $\tilde{K}$  are set once and for all during parameter and key generation, while  $L, \ell$ , and  $\tilde{\ell}$  evolve during a homomorphic evaluation computation, following the noise rate variations: They decrease across homomorphic operations, and are reset to a large value after a bootstrapping.

Variable	Range	Meaning
$N$	$[2^{10}, 2^{16}]$	Power-of-two polynomial modulus: $X^N + 1$
$K, \tilde{K}$	$[10 - 60]$	Multi-precision representation of Torus elements consist of $K$ -bits limbs (so machine words have at least $2K$ bits to handle products), and gadget decompositions produce $\tilde{K}$ -bit outputs. $\tilde{K}$ is in general equal to $K$ , but can be chosen smaller in rare circumstances, since $\tilde{K}$ intervenes in the noise propagation of external products.
$L$	$[20 - 2000]$	Targeted number of bits of precision of the multiprecision arithmetic (to handle a RLWE ciphertext whose noise rate is $\alpha \approx 2^{-L}$ ). For a given cryptographic security parameter, each key dimension is associated to a maximal noise level $L$ : for 128-bit security, count $L_{\max} = 20, 880, 1761$ for respectively $N = 2^{10}, 2^{15}, 2^{16}$
$\ell, \tilde{\ell}, \tilde{\ell}_A$	$[1-100]$	Number of limbs per coefficient in a RLWE ciphertext (without tilde) or in a gadget decomposition (with tilde, and/or subscript depending on the context). These $\ell$ can be thought as $\ell \approx L/K$ , it is the main asymptotic parameter in all the complexities, directly related to the number of elementary vector operations. In practice, for the key sizes and precision we consider above, all these $\ell$ 's fall within the range $[1, 100]$

### 3 A Bivariate Polynomial Representation

We retain the notation and the basic setting from section 2. Let  $K$  be a limb size. We represent (rational) approximations of elements of  $\mathbb{R}[X]/\langle X^N + 1 \rangle$  by elements of the (discrete) quotient ring  $\mathbb{Z}[X, Y]/\langle X^N + 1 \rangle$  (also equal to  $\mathcal{R}[Y]$ , the polynomials in  $Y$  over  $\mathcal{R}$ ) of the bivariate polynomial ring  $\mathbb{Z}[X, Y]$ . The representations are obtained via the evaluation map (ring homomorphism) on the  $Y$  variable

$$\phi_K: \mathcal{R}[Y] \rightarrow \mathbb{R}_N[X], \quad P(X, Y) \mapsto P(X, 2^{-K}). \quad (3)$$

More explicitly, the elements of  $\mathcal{R}[Y]$  are represented by bivariate integer polynomials  $P(X, Y) = \sum a_{i,j} X^i Y^j$  whose degrees in  $X$  are at most  $N - 1$ . Since the ring homomorphism  $\phi_K$  is clearly not injective, an element of  $\mathbb{R}_N[X]$  can have multiple pre-images in  $\mathcal{R}[Y]$ . We will use this property in a crucial way and will be particularly interested in representatives with small coefficients.

**Definition 3.1.** A bivariate polynomial  $P(X, Y) = \sum_{i=0}^{N-1} \sum_{j \geq 0} a_{i,j} X^i Y^j$  is called  $K$ -normalized if  $a_{i,j} \in [-2^{K-1}, 2^{K-1}]$  for all  $i = 0, \dots, N - 1$  and  $j \geq 1$ .

As we are mainly interested in representing elements of  $\mathbb{T}_N[X]$ , that is, formal polynomials over the torus, we often reduce the coefficients of real polynomials to the real interval  $[-1/2, 1/2]$  and hence, we say that

**Definition 3.2.** A bivariate polynomial  $P(X, Y) = \sum_{i=0}^{N-1} \sum_{j \geq 0} a_{i,j} X^i Y^j$  is  $K$ -normalized and reduced if, in addition to being  $K$ -normalized, it satisfies  $a_{i,0} = 0$  for all  $i = 0, \dots, N-1$ .

We now approximate any element of the  $\mathcal{R}$ -module  $\mathbb{T}_N[X]$  up to an arbitrary precision with  $K$ -normalized and reduced bivariate polynomials from  $\mathcal{R}[Y]$ . The proofs of the lemma and the corollary below are simple consequences of the decomposition of the coefficients in base  $2^K$ .

**Lemma 3.1.** For every integer  $L > 0$  and every polynomial  $Q \in \mathbb{R}_N[X]$ , there exists a  $K$ -normalized polynomial  $P(X, Y) \in \mathcal{R}[Y]$  of degree  $\leq \lceil L/K \rceil$  in  $Y$  such that  $\|\phi_K(P) - Q\|_\infty \leq 2^{-L}$ .

**Corollary 3.1.** For all elements  $Q \in \mathbb{T}_N[X]$ , there exists a normalized and reduced polynomial  $P \in \mathcal{R}[Y]$  such that

$$\|\phi_K(P) \bmod \mathcal{R} - Q\|_\infty \leq 2^{-L}.$$

Note that the above approximation of the elements  $Q \in \mathbb{T}_N[X]$  via bivariate polynomials  $P(X, Y) \in \mathcal{R}[Y]$  is reminiscent to a *gadget decomposition* in classical lattice-based cryptography. Intuitively, the presence of the redundant variable  $Y$  corresponds to representing a vector in higher-dimensional space with lower  $\|\cdot\|_\infty$ -norm - a key technique necessary for the design of various FHE schemes (e.g., GSW [20] and TFHE [10]).

The following lemma whose proof is rather formal shows that for any limb size  $K$  and any polynomial in  $P \in \mathcal{R}[Y]$ , there is a unique  $K$ -normalized and reduced polynomial  $Q$  with the same image under  $\phi_K$ .

**Lemma 3.2.** For any limb size  $K$  and any polynomial  $P(X, Y) \in \mathcal{R}[Y]$ , there exists a unique normalized polynomial  $Q(X, Y)$  of the same degrees in both  $X$  and  $Y$  such that  $\phi_K(P) = \phi_K(Q)$ . Additionally, there exists a unique  $K$ -normalized and reduced polynomial  $Q(X, Y)$  such that  $\phi_K(P) = \phi_K(Q) \bmod \mathcal{R}$ .

We are now interested in an efficient algorithm for normalization and reduction of polynomials.

**Lemma 3.3.** Let  $K$  be a limb size and let  $L > 0$  be a specified precision. Given a polynomial  $A(X, Y) \in \mathcal{R}[Y]$  satisfying  $\|A\|_\infty < 2^M$  for some  $M$ , Algorithm 1 outputs a  $K$ -normalized and reduced polynomial  $R(X, Y) \in \mathcal{R}[Y]$  such that  $\|(\phi_K(A) - \phi_K(R)) \bmod \mathcal{R}\| \leq 2^{-L}$  in  $O((L+M)/K)$  element-wise operations (additions/subtractions or binary shifts) on  $N$ -length vectors of integers in  $(-2^{M+1}, 2^{M+1})$ .

A proof of this lemma is given in appendix section A. Using a parameter value  $M$  such that  $M+1 < 64$ , Algorithm 1 allows for efficient AVX and GPU-friendly implementations.

---

**Algorithm 1** Normalization and Reduction

---

**Input:** A target precision of  $L$  bits (that represents  $2^{-L}$ ) and a limb size  $K$

**Input:** An input polynomial  $A(X, Y) = \sum_{k \geq 0} A_k(X)Y^k$  satisfying  $\|A\|_\infty < 2^M$

**Output:** A  $K$ -normalized and reduced polynomial  $R(X, Y) = \sum_{k=1}^{\ell} R_k(X)Y^k$  of degree  $\ell = \lceil L/K \rceil$  in  $Y$  such that  $\|\phi_K(R) - \phi_K(A)\|_\infty \leq 2^{-L}$ .

- 1:  $\text{acc}(X) = 0$
  - 2: **for**  $k = \lceil (L + M)/K \rceil$  **downto** 1 **do**
  - 3:      $\text{acc}(X) \leftarrow \text{acc}(X) + A_k(X)$
  - 4:      $R_k(X) \leftarrow \text{centermod}_{2^K}(\text{acc}(X))$
  - 5:      $\text{acc}(X) \leftarrow (\text{acc}(X) - R_k(X))/2^K$
  - 6: **end for**
  - 7: Return  $R(X, Y) = \sum_{k=1}^{\ell} R_k(X)Y^k$
- 

Since the representations that occur throughout a FHE computation can be parameterized with different limb sizes, we are going to present the general theorems using two limb sizes  $K$  and  $\tilde{K}$  that are not necessarily the same. This is why in appendix section B we present a slightly more general algorithm that converts a  $K$ -representation into a normalized  $\tilde{K}$ -representation. Even if slightly more complex than Algorithm 1, it has the same performance.

### 3.1 Evaluation of external products

Most homomorphic operations reduce to efficient evaluations of Lipschitz functions. Recall that a function  $f: \mathbb{T}_N[X] \rightarrow \mathbb{T}_N[X]$  is called  $\kappa$ -Lipschitz for some parameter  $\kappa > 0$  if  $\|f(x) - f(y)\|_\infty \leq \kappa \|x - y\|_\infty$  for all  $x, y \in \mathbb{T}_N[X]$ . Most of the time, the Lipschitz functions used in homomorphic evaluation are, in addition,  $\mathcal{R}$ -module homomorphisms. Since  $\text{Hom}_{\mathcal{R}}(\mathbb{T}_N[X], \mathbb{T}_N[X]) \simeq \mathcal{R}$ , all Lipschitz functions that are  $\mathcal{R}$ -module homomorphisms are external products by a fixed element of  $\mathcal{R}$ , that is,  $f_u: \mathbb{T}_N[X] \rightarrow \mathbb{T}_N[X]$ ,  $x \mapsto u \cdot x$ . We thus explain how to evaluate efficiently external products. The lemma below formalizes the following simple fact: evaluating  $f_u$  up to a certain target precision  $2^{-L}$  amounts to evaluating the reduced and  $K$ -normalized representations of the external products  $u \cdot 2^{-K}, \dots, u \cdot 2^{-\tilde{K}\tilde{\ell}}$  for sufficiently large  $\tilde{\ell}$ .

The following technical lemma enables us to compute a sufficiently good  $K$ -normalized and reduced representation of the external product  $u \cdot \bullet$  with a prescribed target precision of  $L$  bits by providing (in a precomputation) sufficiently good approximations of the external product of  $u$  with the negative powers  $2^{-Kj}$ .

**Lemma 3.4.** *Let  $u \in \mathcal{R}$  be an integer polynomial and consider a target precision of  $L > 0$  bits. Let  $K$  and  $\tilde{K}$  be two positive integers (limb sizes). Let  $B_1(X, Y), \dots, B_{\tilde{\ell}}(X, Y)$  be  $K$ -normalized and reduced representations of  $u \cdot 2^{-\tilde{K}}, u \cdot 2^{-2\tilde{K}}, \dots, u \cdot 2^{-\tilde{K}\tilde{\ell}}$  with precision  $(\tilde{\ell}N)^{-1}2^{-(\tilde{K}+L-1)}$ . For any  $\tilde{K}$ -normalized and*

reduced bivariate polynomial  $A(X, Y) = \sum_{i=1}^{\ell} A_i(X)Y^i \in \mathcal{R}[Y]$ ,

$$C(X, Y) = \sum_{i=1}^{\tilde{\ell}} A_i(X)B_i(X, Y) \quad (4)$$

is a (non-reduced)  $K$ -representation of  $u \cdot \phi_{\tilde{K}}(A)$  of precision  $2^{-L}$ , i.e.,

$$\|\phi_K(C) - u \cdot \phi_{\tilde{K}}(A)\|_{\infty} \leq 2^{-L} \quad \text{and} \quad \|C\|_{\infty} \leq \tilde{\ell}N2^{K+\tilde{K}-2}.$$

*Proof.* The second bound on  $\|C\|$  comes from the fact that  $A$  and  $B$  are resp  $K$  and  $\tilde{K}$ -normalized, thus each  $\|A_i\|_{\infty}$  and  $\|B_i\|_{\infty}$  are bounded by resp.  $2^{K-1}$  and  $2^{\tilde{K}-1}$ .  $\phi_K(C) = \sum_{i=1}^{\tilde{\ell}} A_i\phi_K(B_i) = \sum_{i=1}^{\tilde{\ell}} A_i(u \cdot 2^{-\tilde{K}i} + e_i)$  where by definition, each  $\|e_i\|_{\infty} \leq (\tilde{\ell}N)^{-1}2^{-(K+L-1)}$ . Therefore,  $\phi_K(C) = u \cdot \phi_{\tilde{K}}(A) + \sum_{i=1}^{\tilde{\ell}} A_i e_i$ , so the first inequality becomes  $\left\| \sum_{i=1}^{\tilde{\ell}} A_i e_i \right\|$  which is  $\leq \sum_{i=1}^{\tilde{\ell}} N \|A_i\|_{\infty} \|e_i\|_{\infty} \leq 2^{-L}$ .  $\square$

If we expand further the right-hand side of (4) via  $B_i(X, Y) = \sum_{j=1}^{\ell} B_{i,j}(X)Y^j$  for  $\ell \leq L + 1 + \tilde{K} + \log_2 \tilde{\ell}$ , we obtain  $C(X, Y) = \sum_{j=1}^{\ell} \sum_{i=1}^{\tilde{\ell}} A_i(X)B_i(X)Y^j$  that can be computed using

- $\ell\tilde{\ell}$  internal polynomial products over  $\mathcal{R}$  with inputs of infinity norm bounded by  $2^{\tilde{K}-1}$  and  $2^{K-1}$ , respectively and output of norm bounded by  $N \cdot 2^{\tilde{K}+K-2}$ .
- $\ell\tilde{\ell}$  additions of polynomials whose norm is bounded by  $N\tilde{\ell} \cdot 2^{\tilde{K}+K-2}$ .

In most FHE operations with GSW ciphertexts such as key switching, relinearization, automorphisms and bootstrapping,  $u \in \mathcal{R}$  depends only on the secret key and is thus known in advance of the homomorphic operation. Anything that depends only on  $u$  can thus be precomputed in an offline phase (in general during the key generation) and only the cross-terms must be evaluated in an online phase (HE evaluation). With this in mind, should we decide to use DFT over only  $X$ , all the multiplications between  $A_i(X)$  and  $B_{i,j}(X)$  become element-wise products on the DFT space and are performed separately for each power of  $Y$ . Therefor, we obtain a nice offline/online phase separation:

**Offline Phase (most often during keygen)**

- $\ell\tilde{\ell}$  bounded DFT's of  $B_{i,j}(X)$  precomputed and given as input

**Online Phase**

- $\ell$  bounded DFT's of  $A_i(X)$
- $\ell\tilde{\ell}$  element-wise multiplications in DFT domain
- $\ell\tilde{\ell}$  element-wise additions in DFT domain
- $\ell$  bounded DFT's for the results  $C_j(X)$
- $O(\ell)$  element-wise additions, shifts or masks to normalize the result (if needed)

Even if we have  $\ell^2$  products to evaluate, the online phase requires only a linear number of DFT's instead of a quadratic number. Fundamentally, it is the exact same root cause that lead [22] to its new asymptotic speedup compared to full-RNS. The main advantage here is that we obtain the representation in a natural way from the normalized gadget decomposition on  $\mathbb{T}_N[X]$ .

### 3.2 External products by secret polynomials over TRLWE.

TRGSW ciphertexts are used to encrypt Lipschitz functions and to evaluate them homomorphically on TRLWE ciphertexts. As explained before,  $\mathcal{R}$ -module homomorphisms are of the form  $f_u$  for a polynomial  $u \in \mathcal{R}$ . When  $u$  is secret, evaluating homomorphically the external product  $u \cdot a$  over an encrypted input  $a \in \mathbb{T}_N[X]$  is the TRGSW-TRLWE external product from [11].

As noted in [5, p.326], the homomorphic evaluation of the external product  $v \rightarrow s \cdot v$  where  $s$  is a small TRLWE secret key could be used as an alternative to the traditional relinearization of the quadratic term  $s^2$  originally used in the CKKS and BFV products. This is done via a relinearization key  $\text{RK} = \text{TRGSW}_s(s)$  using an external product of the form  $\text{RK} \square(a, 0)$  (see [5, Defn.3]). The latter requires only half of a TRGSW ciphertext and half of the running time as we formally explain below via a concept that we call **HalfRGSW** ciphertext. The latter can be used to multiply a secret  $u \in \mathcal{R}$  with a public  $v \in \mathbb{T}_N[X]$  via the  $\odot$  secret-public external product operation

$$\text{HalfRGSW}(u) \odot v \rightarrow \text{TRLWE}(u \cdot v)$$

rather than via the  $\square$  and the full TRGSW ciphertext of  $s^2$ .

For a secret  $u \in \mathcal{R}$  and  $\ell$  TRLWE ciphertexts of zero  $Z := (A|B) \in \mathbb{T}_N[X]^{\ell \times 2}$ , the **HalfRGSW** is defined as

$$\text{HalfRGSW}(u) := (A, B + u(B_g^{-1}, \dots, B_g^{-\ell})) \in \mathbb{T}_N[X]^{2 \times \ell},$$

Note that for all valid TRLWE encryption  $(a, b) \in \mathbb{T}_N[X]^2$  of  $v$  the element  $\text{HalfRGSW}(u) \odot b - \text{HalfRGSW}(us) \odot a \in \mathbb{T}_N[X]^2$  is a valid TRLWE ciphertext (under  $s$ ) of the plaintext  $u \cdot v$ . That yields a more flexible computation of a ciphertext of  $u \cdot v$  instead of computing the traditional external product  $\text{TRGSW}(u) \square(a, b)$ , thus, justifying the principle that *two halves make a whole*. As a bonus, a speed-up can be achieved by using different parameters for the two halves, especially in small levels as in the bootstrapping of TFHE.

**Definition 3.3 (bivariate RLWE ciphertexts).** *Let  $K$  be a limb size, a **bivRLWE** ciphertext  $C$  under a small key  $S \in \mathbb{Z}_N[X]$  of the message  $m \in \mathbb{T}_N[X]$  is materialized by a tuple  $(A, B) \in \mathcal{R}[Y]^2$  of  $K$ -normalized and reduced representations whose phase  $\varphi_{S,K}(A, B) \stackrel{\text{def}}{=} \phi_K(B) - S \cdot \phi_K(A)$  is equal to  $m + e$  where  $e \in \mathbb{T}_N[X]$  is a small Gaussian error. We will note  $\text{bivRLWE}_{S,K,2^{-L}}(m)$  such bivariate RLWE encryption of  $m$  with error bound  $2^{-L}$*

Note that **bivRLWE** ciphertexts satisfy the following *prefix property*: a higher precision ciphertext with small error yields a lower precision one by simply restricting the representation to the first few limbs (prefix). For instance, when restricting a high precision ciphertext  $(A, B)$  of error norm  $\leq 2^{-L+1}$  to degrees  $\ell_B = \lceil (L+2)/K \rceil$  and  $\ell_A = \lceil (L+2 + \log_2 \|S\|_1)/K \rceil$  in  $B$  and  $A$ , respectively, we obtain an encryption of the same plaintext with error  $\leq 2^{-L}$ . The importance of the *prefix property* is that it is a computation-free version of modulus switching. We assume throughout that all **bivRLWE** ciphertexts of precision  $2^{-L}$  are

instantiated with degrees  $\ell_A, \ell_B$  in  $Y$ . We also consider the bivariate HalfRGSW counterpart:

**Definition 3.4 (bivariate HalfRGSW encryptions).** Let  $K, \tilde{K}$  be limb sizes, let  $u \in \mathcal{R}$  be a small polynomial of norm  $\|u\|_1 \leq \kappa$  and let  $S \in \mathcal{R}$  be a small key. We define  $\text{bivHalfRGSW}_{S,K,\tilde{K},2^{-L}}(u)$  (bivariate half RingGSW encryption of  $u$  under  $S$  with precision  $L$ ) to be a family of  $\text{bivRLWE}_{S,K,2^{-L}}$  ciphertexts of  $u \cdot 2^{-\tilde{K}}, u \cdot 2^{-2\tilde{K}}, \dots, u \cdot 2^{-\tilde{\ell}\tilde{K}}$ . The family can be restricted to its first  $\tilde{\ell} = \lceil (L + 2 + \log_2 \kappa) / \tilde{K} \rceil$  ciphertexts. If the  $\text{bivHalfRGSW}$  encryption is given in DFT basis, we denote it by  $\text{bivHalfRGSW}^{\text{DFT}}$ .

These ciphertexts also satisfy an even stronger prefix property: from any  $\text{bivHalfRGSW}$  ciphertext  $C$  of high precision  $\leq 2^{-(L+1)}$ , the truncations to degrees  $\ell_A, \ell_B$  above of its first  $\tilde{\ell} = \lceil (L + 2 + \log_2(\kappa)) / \tilde{K} \rceil$   $\text{bivRLWE}$  ciphertexts form a  $\text{bivHalfRGSW}$  of the same message with lower precision  $\leq 2^{-L}$ .

Because of the prefix property, these ciphertexts can be passed to any function that require a lower precision level  $L' < L$ : in this case, the function will only access the terms of degree  $\ell'_A \leq \ell_A$  and  $\ell'_B \leq \ell_B$  from the first  $\tilde{\ell}' \leq \tilde{\ell}$  elements. This property holds both on paper and also in efficient implementations, where ciphertexts are passed by pointers.

**Theorem 3.1 (half external product (secret  $\times$  public)).** Let  $u \in \mathcal{R}$  be a small polynomial of norm  $\|u\|_1 \leq \kappa$  for some  $\kappa > 0$ , let  $a \in \mathbb{T}_N[X]$  and let  $L$  an output precision parameter. Let  $L_\alpha, L_\beta \geq L$  be parameters satisfying  $2^{-L_\alpha} + 2^{-L_\beta} \leq 2^{-L}$  and let  $L_1 = L_\alpha + \log_2 \kappa$ ,  $\tilde{\ell} = \lceil L_1 / \tilde{K} \rceil$  and  $L_2 = L_\beta + \tilde{K} + \log_2(N\tilde{\ell})$ . If  $C_f = (c_1, \dots, c_{\tilde{\ell}})$  is a  $\text{bivHalfRGSW}_{S,K,\tilde{K},2^{-L_2}}^{\text{DFT}}$  encryption of  $u$  with precision  $\leq 2^{-L_2}$  and  $A(X, Y) = \sum_{i=1}^{\tilde{\ell}} A_i \cdot Y^i$  is a  $\tilde{K}$ -normalized and reduced representation of  $a \in \mathbb{T}_N[X]$  up to  $2^{-L_1}$  then the ciphertext

$$C_f \boxplus a = \text{normalizeReduce} \left( \text{iDFT} \left( \sum_{i=1}^{\tilde{\ell}} \text{DFT}(A_i) \cdot c_i \right) \right)$$

is a  $\text{bivRLWE}_{S,K,2^{-L}}$  encryption of  $u \cdot a$ . Homomorphic evaluation of such an external product with  $\tilde{\ell} = \lceil L_1 / \tilde{K} \rceil$  and  $\ell = \lceil (L_2 + \log_2 N + 2) / \tilde{K} \rceil$  requires

- $\tilde{\ell}$  bounded DFT's of the  $A_i$ 's of norm  $\leq 2^{\tilde{K}-1}$ ,
- $2\tilde{\ell}$  element-wise **admul**'s in DFT domain for polynomials of norm  $\leq N\tilde{\ell}2^{\tilde{K}+K-2}$ ,
- $2\ell$  bounded DFT's for the results  $C_j(X)$  with norm bound  $N\tilde{\ell}2^{\tilde{K}+K-2}$ ,
- $2\ell$  element-wise additions/shifts/masks to normalize and truncate the result.

*Proof.* Letting  $e = a - \sum_{i=1}^{\tilde{\ell}} A_i \cdot 2^{\tilde{K}i}$ , we have  $\|e\|_\infty \leq 2^{-L_1}$ . If  $c'_i = \text{iDFT}(c_i)$  for  $i = 1, \dots, \tilde{\ell}$  then  $\varphi_{S,K}(c'_i) = u \cdot 2^{-\tilde{K}i} + e_i$  where  $\|e_i\|_\infty \leq 2^{-L_2}$ . Since  $C_f \boxplus a = \text{normalizeRed}(\sum_{i=1}^{\tilde{\ell}} A_i \cdot c'_i)$ , it follows that  $\varphi_{S,K}(C_f \boxplus a) = \sum_{i=1}^{\tilde{\ell}} A_i \cdot \varphi_{S,K}(c'_i) = u \cdot \left( \sum_{i=1}^{\tilde{\ell}} A_i \cdot 2^{-\tilde{K}i} \right) + \sum_{i=1}^{\tilde{\ell}} A_i \cdot e_i$ , and therefore,  $\|\varphi_{S,K}(C_f \boxplus a) - u \cdot a\|_\infty \leq \|u \cdot e\|_\infty + \sum_{i=1}^{\tilde{\ell}} \|A_i \cdot e_i\|_\infty \leq 2^{-L_\alpha} + 2^{-L_\beta} \leq 2^{-L}$ . The rest of the theorem is a simple count of operations, the degrees of  $c_i$  being bounded by  $\ell$ .  $\square$

The above theorem applies to the traditional key switching, automorphisms and relinearization operations. To recover the traditional homomorphic product of a secret  $u$  by a  $\text{bivRLWE}$  encrypted ciphertext  $(a, b) \in \mathcal{R}[Y]^2$  under a small key  $s \in \mathbb{Z}_N[X]$ , we can use one  $\text{bivHalfRGSW}$  ciphertext  $C_u$  of  $u$  and one  $\text{bivHalfRGSW}$  ciphertext  $C_{su}$  of  $su$ , and call the pair  $C = (C_u, C_{su})$  a full  $\text{bivRGSW}(u)$ . The full external product is then

$$C \boxtimes (a, b) = C_u \boxtimes b - C_{su} \boxtimes a.$$

Note that since the norms  $\|u\|_1$  and  $\|su\|_1$  in the two halves are distinct, Theorem 3.1 suggests the use of distinct parameters for each half. This reflects the natural property that in a  $\text{RingLWE}$  ciphertext  $(a, b)$ , it suffices to provide the term  $b$  up to a lower precision compared to the term  $a$ . In the next section, we show that this observation speeds up the bootstrapping of  $\text{TFHE}$  by removing 12.5% of the FFTs.

We also merge the final normalizations of the two half-products together to obtain the following

**Corollary 3.2 (full external product (secret $\times$ secret)).** *Let  $K, \tilde{K}, \tilde{K}'$  be limb sizes, let  $u \in \mathcal{R}$  be a small polynomial, let  $v \in \mathbb{T}_N[X]$  be a message, let  $s \in \mathcal{R}$  a small key and  $2^{-L} > 0$  be a target output precision. For all  $L_\alpha, L_\beta, L_\gamma \geq L$  satisfying  $2^{-L_\alpha} + 2^{-L_\beta} + 2^{-L_\gamma} \leq 2^{-L}$ , let  $L_1 = L_\alpha + \log_2 \|u\|_1$ ,  $\tilde{\ell}_B = \lceil (L_1 + 2) / \tilde{K} \rceil$ ,  $\tilde{\ell}_A = \lceil (L_1 + 2 + \log_2 \|s\|_1) / \tilde{K}' \rceil$ ,  $L_2 = L_\beta + \tilde{K} + \log_2(N\tilde{\ell}_B)$  and  $L'_2 = L_\gamma + \tilde{K}' + \log_2(N\tilde{\ell}_A)$ . If  $(A, B) \in \mathcal{R}[Y]^2$  is a  $\text{bivRLWE}_{s, K}$  of  $v \in \mathbb{T}_N[X]$  with noise  $\leq 2^{-L_1}$ , if  $C_u = (c_1, \dots, c_{\tilde{\ell}})$  is a  $\text{bivHalfRGSW}_{s, K, \tilde{K}}$  encryption of  $u$  with precision  $2^{-L_2}$  and if  $C_{su} = (d_1, \dots, d_{\tilde{\ell}'})$  is a  $\text{bivHalfRGSW}_{s, \tilde{K}'}$  encryption of  $su$  with precision  $2^{-L'_2}$  then*

$$(C_u, C_{su}) \boxtimes (a, b) = \text{normalizeRed} \left( \text{idFT} \left( \sum_{i=1}^{\tilde{\ell}_B} \text{DFT}(B_i) \cdot c_i - \sum_{i=1}^{\tilde{\ell}_A} \text{DFT}(A_i) \cdot d_i \right) \right),$$

is a  $\text{bivRLWE}_{s, K, 2^{-L}}$  encryption of  $u \cdot v$ . Here,  $A$  and  $B$  have been  $\tilde{K}'$  and  $\tilde{K}$ -normalized, respectively. In addition, computing this encryption with  $\ell = \lceil (L'_2 + \log_2 N + 2) / K \rceil$  requires at most

- $\tilde{\ell}_A + \tilde{\ell}_B$  bounded DFT's of the  $A_i, B_i$ 's of norm  $\leq 2^{\tilde{K}-1}$ ,
- $2\ell(\tilde{\ell}_A + \tilde{\ell}_B)$  element-wise  $\text{addmul}$ 's in DFT domain for polynomials of norm  $\leq 2N\tilde{\ell}2^{\tilde{K}+K-2}$ ,
- $2\ell$  bounded DFT's for the results  $C_j(X)$  with norm bound  $2N\tilde{\ell}2^{\tilde{K}+K-2}$ ,
- $\tilde{\ell}_A + \tilde{\ell}_B + 2\ell$  element-wise additions/shifts or masks to normalize and truncate the input ciphertext and the final result.

*Proof.* Let  $d'_i = \text{idFT}(d_i)$  and  $c'_i = \text{idFT}(c_i)$ , by the same proof as in Theorem 3.1,  $\varphi_{S, K}(\sum_{i=1}^{\tilde{\ell}_B} B_i \cdot c'_i) = u \cdot \phi_{\tilde{K}}(B) + e_1$  where  $\|e_1\|_\infty \leq 2^{-L_\beta}$  and  $\varphi_{S, K}(\sum_{i=1}^{\tilde{\ell}_A} A_i \cdot d'_i) = su \cdot \phi_{\tilde{K}'}(A) + e_2$  where  $\|e_2\|_\infty \leq 2^{-L_\gamma}$ . Therefore,  $\varphi_{S, K}(C \boxtimes (A, B))$  is the difference  $u \cdot (\phi_{\tilde{K}}(B) - s \cdot \phi_{\tilde{K}'}(A)) + e_1 - e_2 = u \cdot \varphi_{S, K}(A, B) + e_1 - e_2$ . Since  $\varphi_{S, K}(A, B) = v + e$  where  $e \leq 2^{-L_1}$ , we have  $\|\varphi_{S, K}(C \boxtimes (A, B)) - u \cdot v\|_\infty \leq 2^{-L_\alpha} + 2^{-L_\beta} + 2^{-L_\gamma} \leq 2^{-L}$ .  $\square$

### 3.3 Automorphisms in BFV and CKKS

Unlike external products that can operate with very short keys and large noise, BFV and CKKS arithmetic usually operate on much larger parameters, where a single limb is in general the optimal choice. For the rest of the section, we will therefore consider that there is a unique limb size  $K$  (i.e.  $K = \tilde{K}$ ).

Automorphisms of  $\mathbb{T}_N[X]$  are  $\mathcal{R}$ -module homomorphisms. These are the *rotation/conjugation* functions  $\sigma_k: \mathbb{T}_N[X] \rightarrow \mathbb{T}_N[X]$  that substitutes the variable  $X$  with  $X^k$  where  $k$  is odd.  $\sigma_k(a)$  can be computed efficiently in coefficient space over the bivariate representations  $\mathbb{Z}[X, Y]/\langle X^N + 1 \rangle$  by treating each power of  $Y$  independently and mapping  $\sum_{i=1}^{\ell} A_i(X)Y^i$  to  $\sum_{i=1}^{\ell} \sigma_k(A_i(X))Y^i$ .

Over ciphertexts, we just need to observe that  $\sigma(b - sa) = \sigma(b) - \sigma(s) \cdot \sigma(a)$ , so given a `bivHalfRGSW` encryption of  $\sigma(s)$ , we recover the well known homomorphic evaluation:

**Lemma 3.5 (Automorphism in BFV and CKKS).** *Let  $K$  be a limb size, let  $v \in \mathbb{T}_N[X]$  be a message, let  $\sigma$  be an automorphism of  $\mathbb{T}_N[X]$  and let  $L$  be an output precision parameter. For all  $L_\alpha, L_\beta \geq L$  that satisfy  $2^{-L_\alpha} + 2^{-L_\beta} \leq 2^{-L}$ , if  $(a, b) \in \mathcal{R}[Y]^2$  is a `bivRLWE` <sub>$s, K$</sub>  encryption of  $v \in \mathbb{T}_N[X]$  with noise  $2^{-L_\alpha}$  and  $C_{\sigma(s)}$  is a `bivHalfRGSW` <sub>$s, K, K$</sub>  of  $\sigma(s)$  with noise  $2^{-L_\beta}$ , then*

$$\boxed{\sigma}((a, b)) = (0, \sigma(b)) - C_{\sigma(s)} \boxed{\sigma}(a)$$

is a `bivRLWE` <sub>$s, K$</sub>  encryption of  $\sigma(v)$  with noise  $\leq 2^{-L}$ . Note that we can omit the `normalizeReduce` operation of the  $\boxed{\sigma}$  operation in Theorem 3.1 and compute it at the end, after the subtraction. Computing this operation with  $\tilde{\ell} = \lceil L_\alpha/K \rceil$  and  $\ell = \lceil (L_\beta + \log_2 N + 2)/K \rceil$  requires, as in Theorem 3.1,

- $\tilde{\ell}$  bounded DFT's of the  $A_i$ 's, of norm  $\leq 2^{K-1}$ ,
- $2\tilde{\ell}$  element-wise `admmul`'s in DFT domain for polynomials of norm  $\leq N\tilde{\ell}2^{2K-2}$ ,
- $2\ell$  bounded DFT's for the results  $C_{\sigma(s)}(X)$  with norm bound  $N\tilde{\ell}2^{2K-2}$ ,
- $2\ell$  element-wise additions/shifts/masks to normalize and truncate the final result,
- (in addition to the computations in Theorem 3.1),  $2\ell$  evaluations of  $\sigma$  over the input  $A_i$  and  $B_i$ , of norm bound  $2^K$ .

*Proof.* We first verify that  $\varphi_{s, K}(\boxed{\sigma}((a, b)))$  is close to  $\sigma(v)$ . We compute the noise level of the output as  $\|\varphi_{s, K}(\boxed{\sigma}((a, b)) - \sigma(v))\|_\infty = \|\phi_K(\sigma(b)) - \varphi_{s, K}(C_{\sigma(s)} \boxed{\sigma}(a)) - \sigma(v)\|_\infty = \|\sigma(\phi_K(b)) - \sigma(v) - \sigma(s\phi_K(a)) + \sigma(s\phi_K(a)) - \varphi_{s, K}(C_{\sigma(s)} \boxed{\sigma}(a))\|_\infty \leq \|\sigma(\phi_K(b) - s\phi_K(a) - v)\|_\infty + \|\sigma(s\phi_K(a)) - \varphi_{s, K}(C_{\sigma(s)} \boxed{\sigma}(a))\|_\infty$ . Per the LWE definition, and because  $\sigma$  is an isometry,  $\|\sigma(\phi_K(b) - s\phi_K(a) - v)\|_\infty = \|\sigma(e)\|_\infty = \|e\|_\infty \leq 2^{-L_\alpha}$ . Per Theorem 3.1, the second term is a (half) external product noise bounded by  $\|\sigma(s)\sigma(\phi_K(a)) - \varphi_{s, K}(C_{\sigma(s)} \boxed{\sigma}(a))\|_\infty \leq 2^{-L_\beta}$ . Summing the two, the output noise is bounded by  $2^{-L_\alpha} + 2^{-L_\beta} \leq 2^{-L}$ .  $\square$

### 3.4 Internal products in BFV and CKKS

In the previous section, we showed how leveled-FFT allows to compute efficient homomorphic external products and automorphisms. We now show that we also obtain efficient homomorphic internal products between two TRLWE ciphertexts as defined in section 2 and in [5].

Now that the `bivHalfRGSW` can act as a relinearization, we can simplify the second term  $\text{TRGSW}(s) \boxtimes (C_2, 0) = \text{bivHalfRGSW}(s) \boxtimes 0 - \text{bivHalfRGSW}(s^2) \boxtimes C_2$  of the internal product from section 2 back to  $-\text{bivHalfRGSW}(s^2) \boxtimes C_2$ , which is closer to the original keyswitch formulation. To evaluate a BFV or CKKS product between two RingLWE ciphertexts, all we need is to compute :

$$(a_1, b_1) \boxtimes (a_2, b_2) = (C_1, C_0) + \text{bivHalfRGSW}(s^2) \boxtimes C_2,$$

where  $C_0 = b_1 \otimes_* b_2$ ,  $C_2 = a_1 \otimes_* a_2$  and  $C_1 = a_1 \otimes_* b_2 + b_1 \otimes_* a_2$  and where  $\otimes_*$  is one of those two functions, depending on whether the product is a BFV (1) or a CKKS product (2).

We recall that the  $*$  operator is the multiplication over  $\mathbb{R}_N[X]$ , and that the coefficients of normalized and reduced bivariate representations of  $(u, v)$  come automatically lifted in the real interval  $[-1/2, 1/2)$  as  $(\tilde{u}, \tilde{v})$ . Unlike for CRT representations, the rounding to the nearest multiple of  $2^{-L_{in}}$  at the start of the CKKS product is very easy to do in the bivariate base- $2^K$  representation of the coefficients: just truncate the representation to degree  $\ell = \lceil L_{in}/K \rceil$ , and round the last limb to the nearest multiple of  $2^{L_{in} \bmod K}$ . At the end of the product, the multiplication by the small integers  $p$  or the large power of two  $2^{L_{in}}$  are also very straightforward to express on limbs. The major part of the computation remains the multiplication over  $\mathbb{R}_N[X]$  in the middle of the formula. Fortunately, the bivariate representation is a ring isomorphism between  $\mathcal{R}[Y]$  and  $\mathbb{R}_N[X]$ , so the product of two bivariate representations represents the product over the reals.

In the full-RNS, that is based on a CRT-representation, products are carry-less and had only a  $O(\ell N)$  complexity. Unfortunately, our bivariate representation do not have such equivalent, we have to multiply two representations over  $\mathbb{Z}[X, Y]/\langle X^N + 1 \rangle$  of degree  $\leq \ell - 1$  and obtain a result of degree  $\leq 2\ell - 2$ . The fastest multiplication algorithms for such polynomials involve a double DFT on the variables  $X$  and  $Y$  and run asymptotically in time  $O(\ell N \log_2(\ell N)) = O(\ell N(\log_2(\ell) + \log_2(N)))$ , which is of course a logarithmic factor above the carry-less multiplication. However, remember that this operation is in any case not the bottleneck: an internal product over ciphertexts require at least one external product in  $O(\ell^2 N)$  for the relinearization, so asymptotically, an internal product over `bivRLWE` ciphertexts keeps the same cost as its underlying external relinearization product.

To compute a multiplication over two  $K$ -normalized and reduced representations  $a$  and  $b$  of degree  $\ell$ : we first do the DFT on the  $X$  variable to obtain  $N/2$  univariate polynomials  $(A_1(Y), \dots, A_{N/2}(Y))$  and  $(B_1(Y), \dots, B_{N/2}(Y))$ , each in  $\mathbb{C}_N[X]$  of degree  $\leq \ell - 1$ . Then, we simultaneously multiply each  $A_i(Y) \times B_i(Y)$  to form  $N/2$  univariate polynomials  $(C_1(Y), \dots, C_{N/2}(Y))$  of degree at most  $2\ell - 2$ , and to that end, we use the fastest algorithm in practice between the

naive multiplication, Karatsuba multiplication, or DFT over  $Y$ . This takes the form of an SIMD computation, where all the elementary steps necessary to compute the first product  $C_1 = A_1 \times B_1$  are replicated coordinate-wise on a vector of  $N/2$  complex numbers for the other products. Finally, we apply the inverse DFT on  $X$  to recover the final result.

The SIMD part of the computation can be sped-up by noticing that in the BFV product, we actually only need to compute the  $\ell$  terms of lower degree since, higher degree terms in  $Y$  represent a negligible quantity, whereas in CKKS, the final multiplication by  $2^{L_{in}}$  eliminates all the low-degree terms after reduction modulo 1, so we need only the  $\ell$  terms of higher degree. The initial DFTs only need to be computed once per input, and two of the final inverse DFTs and normalizations can be postponed to after the external product. Overall, the following lemma summarizes the total complexity of a CKKS or BFV product between two `bivRLWE` ciphertexts.

**Theorem 3.2 (CKKS or BFV product).** *The homomorphic CKKS or BFV product between two `bivRLWE` ciphertexts  $(a_1, b_1)$  and  $(a_2, b_2) \in \mathcal{R}[Y]^2$  with a relinearization key  $RK = \text{bivHalfRGSw}(s^2)$ , where  $s$  is the secret key requires:*

- $8\ell$  DFT/IDFTs on  $X$  (incl.  $3\ell$  of them for the external product)
- $2\ell^2 + \min(7\ell \log_2(\ell) + 9\ell, \frac{3}{2}\ell(\ell - 1) + 4\ell)$  SIMD operations (add, mul, addmul or twiddle factors) on vectors of size  $N$  (incl.  $2\ell^2$  for the external product).
- $3\ell$  SIMD rounding/normalization operations on vectors of size  $N$  (incl.  $2\ell$  of them in the external product).

The count of operations in this lemma shows that asymptotically, when  $\ell$  grows, the complexity of an internal product is exactly the same as the underlying external product (i.e.  $O(\ell^2 N)$ ) and the overhead induced by the rest of the operations in an internal product is negligible. Also, in practice for small values of  $\ell$ , since the SIMD operations dominate, the running time is at most 1.75 times the running time of the underlying external product.

*Proof.* All our inputs  $a_1, b_1, a_2, b_2$  are given in coefficient space, so  $4\ell$  DFT's on the variable  $X$  are consumed to bring them in DFT-space. Then, to compute  $C_0, C_1$  and  $C_2$ , we study three strategies:

*Naive product:* for one plaintext internal product we have exactly  $\ell(\ell - 1)/2$  `addmul` operations to obtain either the  $\ell$  coefficients of lowest degree (BFV) or the  $\ell$  coefficients of highest degree (CKKS). In the internal homomorphic encryption we have 4 products (to compute  $C_0, C_1$  and  $C_2$ ). This yields a baseline of  $4\ell(\ell - 1)/2$  products.

*Karatsuba product:* The first iteration of Karatsuba is quite interesting, as it notices that  $C_2$  is in fact (close to)  $(a_1 + a_2) \otimes_\star (b_1 + b_2) - C_1 - C_2$  so three products plus four additions/subtractions are needed instead of four products in the naive algorithm, so  $3\ell(\ell - 1)/2$  products and  $4\ell$  additions. Pursuing the

recursion would end-up executing less multiplications than the naive approach (the number of multiplications would be in  $O(\ell^{\log_2 3})$ ), however, one machine word addition often has the same throughput as one multiplication (which is the case of floating point and int64 operations on a CPU), and if we take additions into account, the number of binary operations follows the recursion  $U(\ell) = 1$  if  $\ell = 1$ , otherwise  $3U(\ell/2) + 6\ell$ , which is super-quadratic. More over the subsequent recursive calls of Karatsuba algorithm cannot be done in-place anymore, which require expensive memory allocation. In the end, Karatsuba has to be limited to a single iteration to achieve  $3\ell(\ell - 1)/2 + 4\ell$  operations.

*FFT product:* we evaluate each polynomial on at least  $2\ell - 1$  points, which is a DFT in  $Y$ , multiply those evaluations, and compute the inverse DFT. The most efficient choice is to take the  $2^{\ell'}$ -th roots of unity, where  $\ell'$  is the smallest power of two larger or equal to  $\ell$ . Because  $A_i$  have degrees lower or equal to  $\ell'$ , each of the four direct DFT has  $\ell' \log_2 \ell'$  twiddle factors, and each of the three output DFT that have degree  $2\ell$  have  $\ell'(\log_2 \ell' + 1)$  twiddle factors. That is a total of  $7\ell' \log_2 \ell' + 9\ell'$  operations. This count is exact when  $\ell = \ell'$  is a power of two, however, in the other cases, we would in theory have to consider that  $\ell'$  can grow up to  $2\ell$ . In practice, since both the input and output polynomials are padded with zero coefficients, we can smooth the resulting complexity by propagating these known zero positions across the twiddle factors and eliminate those that are trivial.

Once the products  $C_0, C_1, C_2$  are computed, we need to inverse the DFT in  $X$  and normalize the result: for  $C_2$ , this normalization needs to happen before the external product, however for  $C_0$  and  $C_1$ , this iDFT and normalization can be delayed and merged with the ones at the output of the external product, on the very final result. Therefore we count only  $4\ell$  additional DFT's  $\square$

*Remark 1.* We can decrease the number of DFTs in  $X$  from 8 to 6 for the internal product and from 3 to 2 for the external product, if we take the convention that the default representation for a `bivRLWE` ciphertext is in DFT space (`bivRLWEDFT`) rather than in coefficient space. In the present case, since DFTs are not the limiting factor, it in fact only has a small effect on the overall complexity, and it becomes less natural if the external product uses multiple limb sizes  $K, \bar{K}$ . Also, unlike the Full-RNS case where reduction mod  $Q$  can be done in NTT space as well, all the normalizations in this work and also in [22] require a round-trip in coefficient space.

We conclude this section with Table 1, which recalls the number of  $N$ -dimensional DFT's and the number of SIMD operations over vectors of  $N/2$  complexes, for one half external product (so either one keyswitch or one automorphism), as well as the number of such operations in one CKKS or BFV product, at the maximal levels  $L = 1761$  and 880 that are secure for respectively  $N = 65536$  and 32768. These numbers can be directly comparable to the ones provided in [22]. The table confirms that the larger we can choose the limb size  $K$ , the less elementary operations are required. We will show in the following

sections that each choice of backend naturally comes with one maximal value of  $K$  it can process: for instance, the double floating points FFT is limited to  $K = 19$ , but other 128-bit constructions can afford larger limb sizes, which we explain in section 6 together with some benchmarks.

**Table 1.** Number of operations in one half-external product (i.e. one relinearization, one keyswitch, or one automorphism) and in one BFV/CKKS internal product (IP).

$N = 65536, L = 1761$						$N = 32768, L = 880$				
$K$	$\ell$	# DFTs	# SIMD prods	# IP DFTs	# IP SIMD prods	$\ell$	# DFTs	# SIMD prods	# IP DFTs	# IP SIMD prods
13	136	408	36992	952	50345	68	204	9248	476	15155
16	111	333	24642	777	31128	55	165	6050	385	8900
<b>19</b>	<b>93</b>	<b>279</b>	<b>17298</b>	<b>651</b>	<b>23541</b>	<b>47</b>	<b>141</b>	<b>4418</b>	<b>329</b>	<b>7160</b>
22	81	243	13122	567	19203	40	120	3200	280	5540
25	71	213	10082	497	16028	36	108	2592	252	4482
...										
43	41	123	3362	287	5822	21	63	882	147	1512
46	39	117	3042	273	5265	20	60	800	140	1370
49	36	108	2592	252	4482	18	54	648	126	1107
<b>52</b>	<b>34</b>	<b>102</b>	<b>2312</b>	<b>238</b>	<b>3995</b>	<b>17</b>	<b>51</b>	<b>578</b>	<b>119</b>	<b>986</b>

The dominant parts of the computation are spent in the SIMD products (or assimilated), the cost of the DFTs remain very small.  $K = 19$  is the largest limb-size achieved via `float64` FFT whereas  $K = 52$  is achieved via 120-bit NTT. The two give similar final running time: SIMD products for  $N = 65536$  are microbenchmarked at  $21\mu s$  in the `float64` scenario, whereas the equivalent counterpart for 120-bit NTT take  $88\mu s$ . The number of SIMD products per internal product is an upper-bound that considers the minimum time between the naive multiplication in  $3\ell(\ell - 1)/2$  and the  $O(\ell \log_2(\ell))$ . Because  $\ell$  has to be rounded to the next power of two, some nodes in the resulting Tom-Cook FFT algorithms would compute zeros: these nodes have been eliminated from the count in this table.

## 4 Accelerating TFHE Gate Bootstrapping

In this section, we show how the concept of `halfTRGSW` can speed-up the TFHE library [2]. One TFHE gate bootstrapping is computing  $n = 635$  successive TRGSW-TRLWE external products. We defer to the noise propagation analysis in [10] and the lattice estimator [3] for the explanations on lattice security but in summary, there are essentially two constraints on the external products:

- *Lattice security constraint.* Any TRLWE or TRGSW ciphertext encrypted with a  $N = 1024$  dimensional key (binary or ternary) must have a noise variance

parameter of at least  $V_{\text{GSW}} = V_{\text{min}} = 2^{-50}$  (stdev  $2^{-25}$ ) to provide 128-bits of security.

- *Correctness of bootstrapping constraint.* Each individual external product output shall make the noise variance grow by at most  $\Delta_{\text{max}} = 8.4961 \cdot 10^{-8}$  to ensure that the final ciphertext is decryptable.

The gap between  $V_{\text{min}}$  and  $\Delta_{\text{max}}$  is very small, so we have to pick the external product parameters with extreme care. For instance, it is unrealistic to try to take  $K = \tilde{K}$ , instead, the choice in **TFHE** is to decrease  $\tilde{K}$  as much as needed to contain the noise growth of **TRGSW** ciphertexts, and maintain  $K$  to a fixed value 32, so that **TRLWE** use a single limb of 32-bits, and correspond to our **bivRLWE** with  $\ell = 1$ . Also, the worst-case bounds on noise amplitude increase too fast compared to the reality: as it is shown in **LWE**, we can use the average-case noise propagation formula, which essentially transcriptions of the worst-case theorems where all noises are assimilated to independent Gaussian samples and the theorem operates on their variance instead of their infinity norm.

With this in mind, in **TFHE**, the noise propagation of a **TRGSW-TRLWE** external product (adapted to use the notations of this paper), between a **TRGSW** ciphertext of error variance  $V_{\text{TRGSW}}$  for a decomposition in  $\tilde{\ell}$  limbs of  $\tilde{K}$  bits, and an input **TRLWE** ciphertext of noise variance  $V_{\text{in}}$ , is:

$$V_{\text{out}} - V_{\text{in}} \leq 2\tilde{\ell}N4^{\tilde{K}-1}V_{\text{GSW}} + (1 + \underbrace{N}_{\|s\|_2^2})4^{-\tilde{K}\tilde{\ell}-1}. \quad (5)$$

The choices of  $\tilde{K}, \tilde{\ell}$  in **TFHE** library are 7 and 3, which correspond to a variance growth of  $2.2410 \cdot 10^{-8} < \Delta_{\text{max}}$  per external product. As we know that the running time of **TFHE** is proportional to the number of FFTs per external products, this number is  $2\tilde{\ell} + 2 = 8$ .

We can also see that any attempt to reduce the number of FFTs by setting  $\tilde{\ell} = 2$  for instance fails, as no more limb size  $\tilde{K}$  exists that makes  $V_{\text{out}} - V_{\text{in}}$  in Eq. (5) smaller than  $\Delta_{\text{max}}$ . That's precisely where the concept of **HalfRGSW** helps: if the input **TRLWE** is  $(A, B)$ , and we are allowed to pick different parameters and dimensions  $\tilde{K}_A, \tilde{\ell}_A$  and  $\tilde{K}_B, \tilde{\ell}_B$ , the variance growth bound of Eq. (5) becomes, by analogy with our Corollary 3.2:

$$\left(\tilde{\ell}_A N 4^{\tilde{K}_A-1} V_{\text{GSW}} + N \cdot 4^{-\tilde{K}_A \tilde{\ell}_A - 1}\right) + \left(\tilde{\ell}_B N 4^{\tilde{K}_B-1} V_{\text{GSW}} + 1.4^{-\tilde{K}_B \tilde{\ell}_B - 1}\right)$$

Not unsurprisingly, the term in  $\tilde{K}_A, \tilde{\ell}_A$  is already tight, so the value (7, 3) remains optimal. However the absence of the factor  $N$  in the second term (analogue of  $\|s\|_1$  in worst-case formulae, which only affects the half external product term in  $A$ ) gives us much more flexibility on the choices of  $\tilde{K}_B, \tilde{\ell}_B$ . It turns out that we can finally reduce  $\tilde{\ell}_B$  to 2 and use  $\tilde{K}_B = 8$ . We indeed obtain the following variance growth:  $2.986053 \cdot 10^{-8}$  for the half external product in  $A$  and  $1.1234 \cdot 10^{-8}$  for the one in  $B$ , and we verify that the sum  $4.10946 \cdot 10^{-8}$  stays  $\leq \Delta_{\text{max}}$ , and is thus suitable for bootstrapping.

Because of that improvement, the number of FFTs in TFHE decreases to  $\tilde{\ell}_A + \tilde{\ell}_B + 2 = 7$  instead of 8 per external product, which eventually translates into a total running time speedup of 12.5%, while maintaining the same security level.

## 5 Frontend and Large Number Arithmetic: Bivariate vs. CRT Representations

The main idea in both [22] and the present work is the decoupling of the cyclotomic arithmetic from the large number arithmetic in BFV, CKKS, TFHE and Chimera which, in turn, allows to fully benefit from the small polynomial coefficients produced by the gadget decomposition. In our work, the large number arithmetic corresponds to the bivariate representation we presented in section 3, which, in essence, is a base- $2^K$  decomposition into  $\ell$  limbs with complexity  $O(\ell \log_2 \ell)$ . In [22], the large number arithmetic is performed via CRT representation of big numbers. The parametrization is a bit more complex to grasp since the various parameters  $d, r, p_i, q_i, Q, \tilde{Q}$  appear for historical reasons and for comparison to the former full-RNS algorithm. Yet, the parameters  $\ell$  and  $B'$  have the same meaning as the ones in our work - large numbers of  $L$  bits are represented via their remainders modulo  $\ell$  primes of size  $K$  bits each (which do not need to be NTT-friendly) and the arithmetic has complexity  $O(\ell)$ . Therefore, following the same steps as in section 3, We could define the CRT representation of  $\mathbb{T}_N[X]$  as:

- A sufficiently large family of  $K$ -bit prime numbers  $(q_i)_{i \in I}$ ,
- The space of CRT-representations  $(u_i)_{i \in I} \in \mathcal{R}^I$  (instead of  $\mathcal{R}[Y]$ ), a representation of *degree*  $\ell$  has support over  $u_1, \dots, u_\ell$  only, the rest is zero. A normalized representation satisfies  $u_i \in [-(q_i - 1)/2, (q_i - 1)/2]$  (analogue to Lemma 3.1).
- An  $\mathcal{R}$ -module homomorphism  $\phi^{\text{crt}} : \mathcal{R}^I \rightarrow \mathbb{R}_N[X]$  (the evaluation homomorphism) given by  $\phi^{\text{crt}}((u_i)_{i \in I}) := \sum_{i \in I} u_i / q_i$ . Any element of  $\mathbb{T}_N[X]$  can be approximated at distance  $\leq 2^{-K\ell}$  by a degree- $\ell$  representation (analogue to Lemma 3.3).
- `crtRLWE` of precision  $2^{-K\ell}$  as RingLWE ciphertexts whose representation of coefficients have degree  $\ell$ . (analogue to Definition 3.3)
- `crtHalfRGSW` of precision  $2^{-K\ell}$  as a family of  $\ell$  RingLWE ciphertexts that encrypt  $1/q_i$  with precision  $2^{-K\ell}$  (analogue to Definition 3.4).
- Half external products, full external products, automorphisms work the same on `crtHalfRGSW` and `crtRLWE` as in their bivariate counterparts and follow the same noise propagation formula as in Theorem 3.1, Corollary 3.2, Lemma 3.5, if we replace `biv` by `crt`. However, a few workarounds are needed to define the internal products because  $\varphi$  is not a ring morphism.

The three main differences between the bivariate representation and the CRT-representation are:

- The CRT-representation does not have the prefix property: a last digit removal procedure allows to recompute a weaker precision from a higher one, this computation is relatively cheap, but not free.
- The CRT representation only has proper subrings of precision  $2^{-Ki}$  where  $i$  is integer: the CKKS product must round to an exact multiple of  $1/\prod_{j=1}^i q_j$ , whereas the bivariate representation can afford any rounding to any  $2^{-L}$  bits of precision.
- There is one different product per degree for which  $\phi^{crt}$  is a ring morphism: namely  $x, y \rightarrow q_1, \dots, q_\ell \cdot xy$  in degree  $\ell$ . It happens to be the product we need for CKKS.
- If all half and full external products, keyswitches, relinearizations and automorphisms have exactly the same cost in both bivariate and CRT-cases at a given precision, internal CKKS or BFV products would be a little bit faster in the CRT representation, with a second order logarithmic difference in complexity: of  $2\ell^2 N + 5\ell N + o(\ell N)$  in the CRT-case instead of  $2\ell^2 N + 7\ell \log_2(\ell)N + o(\ell \log(\ell)N)$  in the bivariate case.

In both cases, bivariate and CRT, the cost of an homomorphic operation is dominated by the cost of a half external product in  $O(\ell^2 N)$ , that is included in keyswitches, automorphisms and in BFV/CKKS products. All the big-number arithmetic operations on  $K$ -bit limbs (whether these are base- $2^K$  digits or residue modulo a  $K$ -bit prime) are carried out in parallel in DFT space on  $N$ -dimensions, which makes the two approaches equally parallelizable or vectorizable even if the big number arithmetic needs to propagate carries. Both approaches operate directly on normalized representations that come from a theoretical decomposition of  $\mathbb{T}_N[X]$ , and only require a linear number of DFTs per homomorphic operation (except the keygen that is done offline), which improve upon the quadratic number in Full-RNS. In the next section about the backend, we will see that for both the bivariate and the CRT approaches, these DFTs can be instantiated with similar performance, as either approximated FFT over the complex field, or NTT. And finally, both approaches require an additional normalization in coefficient space after every homomorphic operation to avoid unwanted wraparounds or overflows on  $B'$  bits.

We now point out some very important advantages of choosing the bivariate representation over the CRT representation: namely how the prefix property leads to a greatly reduced memory footprint for multi-level homomorphic circuits, and how the freedom of picking any input level  $L_{in}$  in the CKKS product allows to pack a larger number of homomorphic operations for a given level. We also show that the split of the full external product in two halves lead to an acceleration of TFHE bootstrapping, due to a better handling of the parameters.

### 5.1 Prefix property: huge memory footprint reduction

A typical leveled HE computation in large depth  $\ell$  consumes more memory than the same plaintext counterpart. One of the reasons is because the ciphertext overhead of the encrypted data is proportional to  $\ell$ , so the encrypted data has a

larger memory footprint; this can in general be mitigated with a pipeline from the network or the hard drive that loads in memory only the variables that is currently operated on. Another item must however remain loaded in memory for efficiency reasons: the relinearization key, evaluation keys or bootstrapping key materials. These often consist in `HalfRGSW` encryptions of the secret key, or of rotations of the secret key. Such ciphertexts (let’s call it  $H_\ell$ ) occupy naturally  $2\ell^2 \cdot N$  limbs in memory. For instance for  $N = 65536$  for precision bit level  $L = 40\ell \approx 1760$ , can take between 1 and 10GB of RAM, depending on the limb size  $K$ , and it is important to note that it is exactly the same size whether we opt for the bivariate encoding like in our work, or a CRT encoding like in [22]. The interesting part is that to be able to operate with ciphertexts at all the intermediate precision levels 1 to  $\ell$ , we need not only the highest precision key, but also all the weaker precision versions  $H_1, \dots, H_\ell$  of all these keys present in memory.

CRT-representations offer a “cheap” way of deducing  $H_{i-1}$  from  $H_i$ , that runs in time  $O(\ell^2 N)$ , which is the CRT-modulus rescaling algorithm. It rounds each big-integer element to the nearest multiple of the last prime of the CRT basis, and eliminates that coordinate. After a preprocessing of  $O(\ell^3 N)$  steps, we obtain a vector  $H_1, \dots, H_\ell$  of  $2N\ell^3/3$  limbs. If  $\ell$  is between 30 and 100, this is between 10x and 30 times larger than the highest precision keysize, and it requires to operate on machines with a very large amount of memory (typically servers of 80GB to 200GB of RAM, or swap space).

In the bivariate approach, weaker precision keys are exact prefixes of higher precision keys  $H_i \subset H_{i+1} \subset \dots \subset H_\ell$ . Thus, the last key  $H_\ell$  can be passed directly by pointer as argument to any arithmetic operation that requires a low precision version. This completely removes the need of computing the lower precision keys, which saves both time, and a huge amount of memory, since only a 1GB to 10GB laptop is required to carry out the same efficient computation. More generally, all modulus rescaling operations are completely for free with the bivariate encoding.

## 5.2 Discrete vs. continuous levels in CKKS

The CKKS internal product formula (2) applies to all input bit precision parameters  $L_{\text{in}}$ : for any fixed plaintext mantissa precision  $\rho$ , feeding `TRLWE` ciphertexts with input noise rate  $2^{-L_{\text{in}}}$  to (2) yields a fixed-point product of output noise rate  $2^{-L_{\text{in}} + \rho + \log_2 N}$  that represents the  $\rho$  most-significant bits of the plaintext product. By varying  $L_{\text{in}}$ , we obtain any desired output precision  $L_{\text{out}}$ .

Since it is unrealistic to go back-and-forth between the original large integer coefficients and their full-RNS / CRT representations of CKKS, the number  $2^{L_{\text{in}}}$  in (2) has to be replaced by one of the discrete products of  $Q_i = q_1, \dots, q_i$  of the CRT basis that is close to  $2^{L_{\text{in}}}$ . In full-RNS, each  $q_i$  must be at least 18 bits to be NTT friendly for  $N = 2^{16}$ , and more generally, even though the CRT-representation in this section and in [22] raised the NTT-friendliness restriction, using a smaller bitsizes  $K$  make the scheme less efficient; in practice, it is usual to have  $K$  between 30-40 bits. If the noise of an input ciphertext is not close

to an exact multiple of  $K$ -bits, (usually referred as half-multiplicative level), it has to be mod-rescaled to the next available one before a CKKS product can be attempted. E.g. after a cheap homomorphic operation like the sum of 1000 ciphertexts, or after computing a trace, or just after a small linear combination with coefficients between -30 and 30, the noise has only increased by 5 to 8 bits, and to be able to start the next CKKS product, we need to multiply/rescale the ciphertext by a factor  $2^{K-8}$  to  $2^{K-5}$  bits to come back to a (half-)level boundary. This represents a considerable loss of homomorphic budget.

Our bivariate approach does not have this limitation: any input noise level can be fed-in to Equation 2, or equivalently, any desired output noise level can be attained without any penalty on the noise overhead, that stays  $\rho + \log_2(N)$  bits. A `bivRLWE` ciphertext at one given level of noise will be able to support more homomorphic operations before the noise is saturated. This statement is corroborated by the current progresses in FHE compilers: whereas the first CKKS compilers were operating on "full" 40-bit levels, the recent Hecate compiler by [24] was able to generate 27% more efficient CKKS circuits by considering "half-levels" of 20-bits. This line of work should continue to give more impressive results using the "continuous" single-bit levels and the simple noise propagation rule  $L_{\text{out}} = L_{\text{in}} - \rho - \log_2 N$  during a product, and matches the BFV product.

While we are waiting for the next generation of CKKS compilers on continuous levels, our scheme is trivially backwards-compatible with CKKS compilers that support discrete full-levels or half-levels only: we just need to set  $L_{\text{in}} = 40.\text{level}$ , and as a bonus, the prefix-property of our representation implicitly relinearizes the input ciphertext for free!

### 5.3 Scheme switching: across schemes and representations

Chimera [5] introduced a few years ago a scheme switching concept that would allow to switch back and forth between TFHE, BFV and CKKS representations, however it did not provide a solution to efficiently deal with big numbers. At that time, the construction remained under the assumption that it would be efficient to do efficient multiprecision fixed-point FFT, and none of the experiments with GMP and MPFR provided fast enough results for levels  $> 2$  compared to the full-RNS approach.

In this paper, the bivariate representation (which we see as a direct consequence of [22] second gadget decomposition), is in some sense the missing piece that allows to instantiate Chimera's scheme switching efficiently in practice.

In addition, the `HalfRGSW` product can efficiently convert not only between schemes, but also back and forth between the `crTRLWE` and `bivRLWE` ciphertexts, preserving the noise level and the plaintext: in the CRT to bivariate direction, all we need is a family of  $K$ -normalized representations of each  $1/q_i$ , and in the bivariate to CRT direction, a family of normalized CRT-representations of  $1/2^{-Ki}$ , both with precision  $2^{-K\ell}$ . In both cases, Theorem 3.1 fulfills such conversion in time  $O(\ell^2 N)$ .

This makes very efficient bridges between not only TFHE, CKKS and BFV ciphertexts, but for each scheme, between their CRT and bivariate representations.

## 6 Backend and Cyclotomic Multiplications: Approximate FFT or NTT

At a low level, we need a efficient arithmetic (additions and multiplications) in the cyclotomic ring  $\mathcal{R}$  where the coefficients of the polynomials are bounded in the worst case by  $B'_{\text{worst}} = 2\ell N 2^{K+\tilde{K}-2}$ . We can get a tighter average case bound if we pay a closer attention to the expressions present in the previous section. We then notice that it is sufficient for such arithmetic to be able to successfully evaluate, with overwhelming probability, expressions of the form  $\sum_{i=1}^{2\ell} a_i(X) \cdot b_i(X)$  where  $a, b$  have their coefficients computationally indistinguishable from uniformly distributed in respectively  $[-2^{K-1}, -2^{K-1})$  and  $[-2^{\tilde{K}-1}, -2^{\tilde{K}-1})$ . Indeed, these inputs are base- $2^K$  decompositions of LWE ciphertexts or fresh GSW ciphertexts, and any computational bias against the uniform distribution would form an attack on these schemes. In other words, if we randomize `bivRLWE` ciphertexts during normalization (e.g. we can always mask them with random ciphertexts of zero), with an overwhelming probability  $1 - \varepsilon$  the arithmetic just needs to handle elements of size  $B'_{\text{avg}} = C_\varepsilon \sqrt{2\ell N} 2^{K+K'-2}$  instead of the worst case  $B'_{\text{worst}} = 2\ell N 2^{K+K'-2}$ . We will use  $C_\varepsilon = 17$  in this section, that corresponds to an error probability  $\varepsilon < 2^{-40}$ .

We present two equivalently good ways of handling such arithmetic: floating point or fixed-point approximations of the continuous FFT on backends whose mantissa can store  $B'$ , or NTT over a friendly modulus larger than  $B'$ . The underlying arithmetic must be sufficiently atomic to be considered as native operations, therefore we will limit ourselves to 64-bit or 128-bit  $B'$ .

In a nutshell, the main result is that each choice of backend is bound to a precision  $B'$  and lead to a maximal limb size  $K$ : for all practical HE dimensions, the `float64` backend corresponds to  $K \leq 19$ , the `float128` backend (or fixed-point backends via 104-bit arithmetic in `AVX_IFMA`) would correspond to  $K \leq 49$ , doing NTT over a 60-bits modulus corresponds to  $K \leq 22$ , and  $K \leq 52$  for a 120-bit modulus. And the key takeaway is that the most important success factor for a backend to be FHE friendly is to support the largest machine-word arithmetic; it is much less important if that arithmetic is modulo a power of two (fixed-point FFT), modulo a user-friendly prime number (NTT), or floating point (FFT).

### 6.1 Floating point backends

The first choice when it comes to FFT on bounded numbers is floating point backends. This algorithm is well studied worldwide, and has been successfully powering the `TFHE` library since its origins. A naive application of the formula  $B'_{\text{avg}} = C_\varepsilon \sqrt{2\ell N} 2^{K+\tilde{K}-2}$  for  $K = \tilde{K}$  and provided that we can use the 52 bits of mantissa without loss would bound  $K$  around 19. However, a lot of intermediate floating point operations occur between the start of the FFT, the products, the inverse FFT and we have to guarantee that the final error remains bounded by  $1/2$  to be recoverable by rounding. We decided to treat the problem experimentally, by sampling uniformly random polynomials  $A(X), B(X) \in \mathcal{R}$

with coefficients  $\in [-2^{K-1}, 2^{K-1}]$ , multiplying these two polynomials via  $C = iFFT(FFT(A) * FFT(B))$  using 64-bit double-precision floats, and measuring the amplitude and the standard deviation of the error  $C - AB$ , as a function of  $N$  and  $K$ . From these measurements, we estimate the probability that a sum of such terms satisfy  $\left\| \sum_{i=1}^{2^\ell} C_i - A_i B_i \right\|_\infty < 1/2$ , which is sufficient to recover the actual result by rounding.

**Lemma 6.1.** *Let  $\ell, N \in \mathbb{N}$  and  $\varepsilon > 0$ . The maximal value of  $\sigma$  such that with probability  $\geq 1 - \varepsilon$ , the sum  $v = \sum_{i=1}^{2^\ell} v_i$  of  $2^\ell$  independent real vectors  $v_i \in \mathbb{R}^N$  with independent Gaussian coordinates of mean 0 and stdev  $\sigma$  is bounded by  $\|v_\infty\| < 1/2$  is  $\sigma \leq 1 / \left( \sqrt{16\ell} \operatorname{erfinv} \left( (1 - \varepsilon)^{1/N} \right) \right)$ .*

For a target probability error  $\varepsilon = 2^{-40}$ ,  $N = 65536$  and  $\ell = 100$ , the maximum  $\sigma$  is 0.00414. Experimentally we have observed that for those values the largest limb size we can choose is  $K = 19$  for 64-bit words. We obtain the same result with  $N = 32768$ . For completeness, we have run the same experiment with 128-bit words (`float128`) that have a mantissa of 112 bits. In this case we see that we can choose  $K = 48$ , however the obtained  $\sigma$  is very close and we could even choose  $K = 49$  given that  $\ell$  is going to be much smaller than 100 due to the constraints outlined in Table 1.

Table 2 shows the results of these experiments for different values of  $K$  and  $N$  for 64-bit (`double`) and 128-bit (`float128`) words. The objective is to find the largest value of  $K$  for a given error boundary.

Unfortunately, to this date, `float128` is not a primitive type on x86 architectures and it's not available in all targets either. GCC supports `float128` operations via the quad-precision math library `libquadmath`, which is shipped along since version 4.6 of GCC, however, `float128` operations provide a poor performance compared to AVX accelerated doubles. Our experiment shows that due to the increased limb size  $K$ , large precision floats have potential, but without any dedicated hardware support, quad floats are not performant enough for our homomorphic backends. As an opening, we could however investigate the newer AVX-IFMA extensions set, whose instructions are already available on most recent commodity CPUs. These instructions allow to easily emulate 104-bit fixed-point arithmetic, and seems to be a perfect hardware-accelerated candidate for FFT computations, on a large limb sizes.

## 6.2 NTT backends over a fixed modulus

As an alternative to the complex FFT, it is well known that cyclotomic multiplications can also be carried out by NTT, which has been the default choice in Full-RNS representation and later in [22]. All we need is one (or a product of) modulus larger than  $B'$ , which must be NTT-friendly unlike the frontend-ones in the CRT-representation, and replace all DFTs with NTT. NTT over moduli that are products of 30 or 60-bit primes can be AVX2 accelerated, however, for a given limb size, because of the overhead necessary modulus reductions, and also

**Table 2.** Experimental standard deviation of floating point errors after a sum of FFT products for given  $N$  and varying  $K$  with 64 and 128-bit floats.

$N$	64-bit representation					128-bit representation					
$64k$	17	18	<b>19</b>	20	21	47	48	<b>49</b>	50	51	$K$
	0.0002	0.0008	<b>0.0031</b>	0.0124	0.0498	0.0003	0.0011	<b>0.00416</b>	0.0169	0.067	$\sigma$
$32k$	17	18	<b>19</b>	20	21	47	48	<b>49</b>	50	51	$K$
	0.0001	0.0005	<b>0.0021</b>	0.0085	0.034	0.0002	0.0007	<b>0.0028</b>	0.0112	0.0448	$\sigma$

because CPU’s and GPU’s better support 64-bit floating point operations rather than their integer counterpart, one NTT on a 30-bit modulus costs roughly the same as one FFT on 64-bits floating points. To compensate this setback, currently the best trade-off we have is to target 120-bit modulus NTT, because the larger supported limb size  $K = 52$  reduces the number of elementary operations in an external product: 120-bit NTT with  $K = 52$  that has approximately the same running time as the 64-bit floating point FFT counterpart with  $K = 19$ .

Comparative micro-benchmarks between floating point and NTT elementary operations are given in Table 3: one DFT/iDFT operation is either an FFT or an NTT on a consecutive array of  $N$  elements. `float64` and `float128` are self-descriptive, whereas in 60-bit and 120-bit NTT, one element  $x$  is represented as respectively two or four (lazily-reduced) 64-bit integers equal to  $x$  modulo  $(q_1, q_2)$  or  $(q_1, q_2, q_3, q_4)$ . One SIMD `addmul/twiddle` consists in either one operation  $r = r + ab$  over vectors in  $\mathbb{C}^{N/2}$  or  $(\mathbb{Z}/Q\mathbb{Z})^N$ , or one twiddle-factor  $(a, b) \rightarrow (a + \omega.b, a - \omega.b)$  where  $\omega$  is a fixed (general) root of unity, whichever is slower. For  $N = 32k$  or  $64k$ , the twiddle factor is in general 10% faster than the `addmul`, despite the fact that the it contains one more subtraction, which indicates that these operations are memory-bound. `float64`, and the two NTTs use AVX2 instructions, whereas `float128` is powered by `libquadmath` and does not benefit from any particular acceleration, except for automorphisms that use only copy and sign-bit flipping.

Additionally, benchmarks on elementary homomorphic operations are shown in Table 4. We expect that the performance numbers in these tables will be in constant evolution, as new hardware tend to support larger precision arithmetic. However, unlike what was commonly believed so far, any chip or device that can: either efficiently approximate the complex FFT with reasonable precision, or execute NTT on even just one single NTT-friendly modulus of reasonable size, is suitable for fast homomorphic computations at any depth. The efficiency of such device is directly related to the number of bits of mantissa or modulus it natively supports.

## Conclusion

In this paper we extend the key decomposition techniques from [22] by using a simpler and more natural base- $2^K$  representation. It allows a better understand-

**Table 3.** DFT and SIMD arithmetic operations.

Backend	Float64 FFT		Float128 FFT		60-bit NTT		120-bit NTT	
K	19		48		22		52	
N	64k	32k	64k	32k	64k	32k	64k	32k
DFT/iDFT	179 $\mu$ s	58 $\mu$ s	45.6ms	21.6ms	528 $\mu$ s	183 $\mu$ s	1332 $\mu$ s	528 $\mu$ s
SIMD addmul/twiddle	23 $\mu$ s	8 $\mu$ s	1.96ms	1ms	42 $\mu$ s	22 $\mu$ s	88 $\mu$ s	42 $\mu$ s
Automorphism	58 $\mu$ s	29 $\mu$ s	63 $\mu$ s	31 $\mu$ s	63 $\mu$ s	31 $\mu$ s	89 $\mu$ s	32 $\mu$ s

**Table 4.** Total running time per homomorphic operation over RLWE ciphertexts: CRT representations for full-RNS and [22], bivariate representations in our case.

Operation	Keyswitch		Automorphism	CKKS product
Size	$N = 64k$	$N = 32k$	$N = 64k$	$N = 64k$
	L=1761	L=880	L=1761	L=1761
- Full-RNS (best r)	1.143s	0.151s	1.261s	1.324s
- [22] (best r)	0.562s	0.101s	0.681s	0.729s
- ours: biv + fft-f64 ( $K = 19$ )	0.558s	0.064s	0.570s	0.816s
- ours: biv + ntt120 ( $K = 52$ )	0.511s	0.069s	0.518s	0.735s

ing of the parametrization of TFHE, CKKS and BFV schemes. We decrease the memory footprint of the relinearization keys, which in turn speeds-up not only the main operations in CKKS and BFV schemes, but also low-depth computations in TFHE.

## References

1. DARPA:data protection in virtual environments (DPRIVE).
2. TFHE: Fast Fully Homomorphic Encryption over the Torus. <https://tfhe.github.io/tfhe/>.
3. M. R. Albrecht, R. Player, and S. Scott. On the concrete hardness of learning with errors. *Journal of Mathematical Cryptology*, 9(3):169–203, 2015.
4. J.-C. Bajard, J. Eynard, M. A. Hasan, and V. Zucca. A full RNS variant of FV like somewhat homomorphic encryption schemes. In R. Avanzi and H. Heys, editors, *Selected Areas in Cryptography – SAC 2016*, pages 423–442, Cham, 2017. Springer International Publishing.
5. C. Boura, N. Gama, M. Georgieva, and D. Jetchev. CHIMERA: combining ring-lwe-based fully homomorphic encryption schemes. *J. Math. Cryptol.*, 14(1):316–338, 2020.
6. Z. Brakerski, C. Gentry, and V. Vaikuntanathan. (leveled) fully homomorphic encryption without bootstrapping. In S. Goldwasser, editor, *Innovations in Theoretical Computer Science 2012, Cambridge, MA, USA, January 8-10, 2012*, pages 309–325. ACM, 2012.
7. R. Cammarota. Intel HERACLES: homomorphic encryption revolutionary accelerator with correctness for learning-oriented end-to-end solutions. In F. Regazzoni and M. van Dijk, editors, *Proceedings of the 2022 on Cloud Computing Security Workshop, CCSW 2022, Los Angeles, CA, USA, 7 November 2022*, page 3. ACM, 2022.
8. J. H. Cheon, K. Han, A. Kim, M. Kim, and Y. Song. A full rns variant of approximate homomorphic encryption. *Selected areas in cryptography : ... annual international workshop, SAC ... proceedings. SAC*, 11349:347–368, 2018.
9. J. H. Cheon, A. Kim, M. Kim, and Y. S. Song. Homomorphic encryption for arithmetic of approximate numbers. In T. Takagi and T. Peyrin, editors, *Advances in Cryptology - ASIACRYPT 2017 - 23rd International Conference on the Theory and Applications of Cryptology and Information Security, Hong Kong, China, December 3-7, 2017, Proceedings, Part I*, volume 10624 of *Lecture Notes in Computer Science*, pages 409–437. Springer, 2017.
10. I. Chillotti, N. Gama, M. Georgieva, and M. Izabachène. Faster fully homomorphic encryption: Bootstrapping in less than 0.1 seconds. In *ASIACRYPT 2016, Proceedings, Part I*, volume 10031 of *LNCS*, pages 3–33. Springer, 2016.
11. I. Chillotti, N. Gama, M. Georgieva, and M. Izabachène. TFHE: fast fully homomorphic encryption over the torus. *J. Cryptol.*, 33(1):34–91, 2020.
12. J. Cooley and J. Tukey. An algorithm for the machine calculation of complex fourier series. *Mathematics of Computation*, 19(90):297–301, 1965.
13. M. Creeger. The rise of fully homomorphic encryption: Often called the holy grail of cryptography, commercial FHE is near. *ACM Queue*, 20(4):39–60, 2022.
14. K. Derya, A. C. Mert, E. Öztürk, and E. Savas. Coha-ntt: A configurable hardware accelerator for ntt-based polynomial multiplication. *Microprocess. Microsystems*, 89:104451, 2022.
15. L. Ducas and D. Micciancio. FHEW: bootstrapping homomorphic encryption in less than a second. In *EUROCRYPT 2015, Proceedings, Part I*, volume 9056 of *LNCS*, pages 617–640. Springer, 2015.
16. J. Fan and F. Vercauteren. Somewhat practical fully homomorphic encryption. *IACR Cryptol. ePrint Arch.*, page 144, 2012.

17. C. Gentry. Fully homomorphic encryption using ideal lattices. In M. Mitzenmacher, editor, *Proceedings of the 41st Annual ACM Symposium on Theory of Computing, STOC 2009, Bethesda, MD, USA, May 31 - June 2, 2009*, pages 169–178. ACM, 2009.
18. C. Gentry, S. Halevi, and N. P. Smart. Homomorphic evaluation of the AES circuit. *IACR Cryptol. ePrint Arch.*, page 99, 2012.
19. C. Gentry, A. Sahai, and B. Waters. Homomorphic encryption from learning with errors: Conceptually-simpler, asymptotically-faster, attribute-based. In R. Canetti and J. A. Garay, editors, *Advances in Cryptology - CRYPTO 2013 - 33rd Annual Cryptology Conference, Santa Barbara, CA, USA, August 18-22, 2013. Proceedings, Part I*, volume 8042 of *Lecture Notes in Computer Science*, pages 75–92. Springer, 2013.
20. C. Gentry, A. Sahai, and B. Waters. Homomorphic encryption from learning with errors: Conceptually-simpler, asymptotically-faster, attribute-based. In *CRYPTO 2013, Proceedings, Part I*, volume 8042 of *LNCS*, pages 75–92. Springer, 2013.
21. S. Halevi, Y. Polyakov, and V. Shoup. An improved RNS variant of the BFV homomorphic encryption scheme. *IACR Cryptol. ePrint Arch.*, page 117, 2018.
22. M. Kim, D. Lee, J. Seo, and Y. Song. Accelerating HE operations from key decomposition technique. *preprint available at <https://eprint.iacr.org/2023/413.pdf>*, 2023.
23. I. Kundu, E. Cottle, F. Michel, J. Wilson, and N. New. The dawn of energy efficient computing: Optically accelerating the fast fourier transform core. In *Photonics in Switching and Computing 2021*. Optica Publishing Group, 2021.
24. Y. Lee, S. Heo, S. Cheon, S. Jeong, C. Kim, E. Kim, D. Lee, and H. Kim. Hecate: Performance-aware scale optimization for homomorphic encryption compiler. In *2022 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*, pages 193–204, 2022.
25. A. C. Mert, E. Öztürk, and E. Savas. FPGA implementation of a run-time configurable ntt-based polynomial multiplication hardware. *Microprocess. Microsystems*, 78:103219, 2020.
26. Ö. Özerk, C. Elgezen, A. C. Mert, E. Öztürk, and E. Savas. Efficient number theoretic transform implementation on GPU for homomorphic encryption. *J. Supercomput.*, 78(2):2840–2872, 2022.
27. O. Regev. On lattices, learning with errors, random linear codes, and cryptography. *J. ACM*, 56(6):34:1–34:40, 2009.
28. S. S. Roy, F. Turan, K. Järvinen, F. Vercauteren, and I. Verbauwhede. Fpga-based high-performance parallel architecture for homomorphic computing on encrypted data. In *25th IEEE International Symposium on High Performance Computer Architecture, HPCA 2019, Washington, DC, USA, February 16-20, 2019*, pages 387–398. IEEE, 2019.
29. A. Schönhage and V. Strassen. Schnelle Multiplikation grosser Zahlen. *Computing (Arch. Elektron. Rechnen)*, 7:281–292, 1971.
30. F. Turan, S. S. Roy, and I. Verbauwhede. HEAWS: an accelerator for homomorphic encryption on the amazon AWS FPGA. *IEEE Trans. Computers*, 69(8):1185–1196, 2020.
31. E. R. Türkoglu, A. S. Özcan, C. Ayduman, A. C. Mert, E. Öztürk, and E. Savas. An accelerated GPU library for homomorphic encryption operations of BFV scheme. In *IEEE International Symposium on Circuits and Systems, ISCAS 2022, Austin, TX, USA, May 27 - June 1, 2022*, pages 1155–1159. IEEE, 2022.

## A Normalization and reduction lemma proof

*Proof.* Let  $n = \lceil (L + M)/K \rceil$  be the number of algorithm iterations. Let  $\text{acc}^{(i)}$  be accumulator value after step 3 of iteration  $i$  and let  $\text{acc}^{(n+1)} = 0$ . At iteration  $i$ ,  $n \geq i \geq 1$ , we have:

$$\text{acc}^{(i)} = A_i + \epsilon_{i+1} \text{ and } R_i = \text{acc}^{(i)} - \epsilon_i 2^K,$$

here  $\epsilon_i = \lfloor \text{acc}^{(i)} 2^{-K} \rfloor$  and has integer values.

The evaluation of result polynomial at  $2^{-K}$  gives:

$$\phi_K(R) = \sum_{i=1}^{\ell} R_i 2^{-iK} = \sum_{i=1}^{\ell} (A_i + \epsilon_{i+1} - \epsilon_i 2^K) 2^{-iK}$$

Expanding the sum and simplifying common expressions we obtain:

$$\phi_K(R) = \sum_{i=1}^{\ell} A_i + \epsilon_{\ell+1} 2^{-\ell K} - \epsilon_1$$

Now, we will prove that  $\|(\phi_K(A) - \phi_K(R)) \bmod \mathcal{R}\|_{\infty}$  is smaller than  $2^{-L}$ . We have:

$$\phi_K(A) - \phi_K(R) = A_0 + \sum_{\ell < i} A_i 2^{-iK} - \epsilon_{(\ell+1)} 2^{-\ell K} + \epsilon_1.$$

Observe that terms  $A_0$  and  $\epsilon_1$  are integers and are reduced by  $\bmod \mathcal{R}$  operation, we obtain:

$$(\phi_K(A) - \phi_K(R)) \bmod \mathcal{R} = \sum_{\ell < i} A_i 2^{-iK} - \epsilon_{\ell+1} 2^{-\ell K}.$$

From relations (A) it is easy to see that  $A_i - \epsilon_i 2^K \equiv R_i - \epsilon_{i+1}$ .

Using previous relations we have:

$$\begin{aligned} (\phi_K(A) - \phi_K(R)) \bmod \mathcal{R} &= \\ &= 2^{-(\ell+1)K} (A_{\ell+1} + \epsilon_{\ell+1} 2^K) + \sum_{\ell+1 < i} A_i 2^{-iK} \\ &= 2^{-(\ell+1)K} R_{\ell+1} + \sum_{\ell+1 < i} A_i 2^{-iK} - \epsilon_{\ell+2} 2^{-(\ell+1)K} \\ &= \dots \\ &= \sum_{i=\ell+1}^n R_i 2^{-iK} + \sum_{n < i} A_i 2^{-iK} \end{aligned}$$

Looking at the infinity norm we have:

$$\begin{aligned} \|(\phi_K(A) - \phi_K(R)) \bmod \mathcal{R}\|_\infty &= \left\| \sum_{i=\ell+1}^n R_i 2^{-iK} + \sum_{n < i} A_i 2^{-iK} \right\|_\infty \\ &< \left\| \sum_{i=\ell+1}^n 2^{-iK+K} + \sum_{n < i} 2^{-iK+M} \right\|_\infty \\ &\leq 2^{-\ell K+1} \leq 2^{-L} \end{aligned}$$

*Element-wise operations are on integers in interval  $(-2^{M+1}, 2^{M+1})$*

The accumulator variable `acc` (at step 3) has the largest values during algorithm execution. We will prove that its magnitude (i.e. infinity norm) never exceeds  $2^{M+1}$ .

Algorithm step 3 increases accumulator value by at most  $2^M$  and step 5 divides the new value by  $2^K$ . After the first iteration, we have  $\|\mathbf{acc}^{(n)}\|_\infty < 2^M$ , after second iteration  $\|\mathbf{acc}^{(n-1)}\|_\infty < 2^{M-K} + 2^M$  and so on until the last iteration where we have  $\|\mathbf{acc}^{(1)}\|_\infty < \sum_{0 \leq i < n} 2^{M-iK}$ , which is the maximum accumulator magnitude attained during algorithm execution. We can rewrite the last expression as:

$$\|\mathbf{acc}^{(1)}\|_\infty < \sum_{0 \leq i < n} 2^{M-iK} = 2^M \cdot \sum_{0 \leq i < n} 2^{-iK} < 2^M \cdot 2 = 2^{M+1},$$

which proves the accumulator bound.

*Complexity* Algorithm steps 3-5 are executed  $\lceil (L+M)/K \rceil$  times. In each iteration 5 operations are performed: an addition (step 3), 3 shifts (2 in step 4 and 1 in step 5) and a subtraction (step 4). The overall complexity of the algorithm is  $5 \cdot \lceil (L+M)/K \rceil$  element-wise operations. □

## B Normalization and Conversion of Limb Sizes

In this section we present a general conversion and normalization algorithm with two limb size  $K$  and  $\tilde{K}$ . The algorithm is slightly more complex algorithm than the one with  $K = \tilde{K}$  described in Algorithm 1, since it must handle additional binary shifts to synchronize the limb sizes, and thus, the for loop on a single index  $k$  is replaced by two while loops and two indexes  $k, \tilde{k}$  that decrease at their respective speed. Besides that, it follows the same principle as the single-limb normalization.

---

**Algorithm 2** Conversion, Normalization and Reduction (from  $K$  to  $\tilde{K}$ ).

---

**Input:** A target precision of  $L$  bits, an input limb size  $K$  and an output limb size  $\tilde{K}$

**Input:** An input polynomial  $A(X, Y) = \sum_{k \in \mathbb{Z}_{\geq 0}} A_k(X) Y^k$  satisfying  $\|A\|_{\infty} \leq 2^B$

**Output:** A  $\tilde{K}$ -normalized and reduced polynomial  $R(X, Y)$  of degree  $\leq \lceil L/\tilde{K} \rceil$  in  $Y$  such that  $\|\phi_{\tilde{K}}(R) - \phi_K(A)\|_{\infty} \leq 2^{-L}$ .

- 1:  $k = \lceil (L + B)/K \rceil$ ,  $\tilde{k} = \lceil (L + B)/\tilde{K} \rceil$
- 2:  $\text{acc}(X) = \lfloor A_k(X) \cdot 2^{\tilde{k}\tilde{K} - kK} \rfloor$
- 3: **while**  $\tilde{k} \geq 1$  **do**
- 4:   **while**  $(\tilde{k} - 1)\tilde{K} < (k - 1)K$  **do**
- 5:      $\text{acc}(X) \leftarrow \text{acc}(X) + A_{k-1}(X) \cdot 2^{(k-1)K - (\tilde{k}-1)\tilde{K}}$
- 6:      $k \leftarrow k - 1$
- 7:   **end while**
- 8:    $R_{\tilde{k}}(X) \leftarrow \text{centermod}_{2^{\tilde{k}\tilde{K}}}(\text{acc}(X))$
- 9:    $\text{acc}(X) \leftarrow (\text{acc}(X) - R_{\tilde{k}}(X))/2^{\tilde{K}}$
- 10:    $\tilde{k} \leftarrow \tilde{k} - 1$
- 11: **end while**
- 12: Return  $R(X, Y) = \sum_{k=1}^{\lceil L/K \rceil} R_k(X) Y^k$

---

## C Complexity of half-external products and BFV/CKKS internal products

$N = 65536, L = 1761$						$N = 32768, L = 880$				
$K$	$\ell$	# DFTs	# SIMD prods	# IP DFTs	# IP SIMD prods	$\ell$	# DFTs	# SIMD prods	# IP DFTs	# IP SIMD prods
13	136	408	36992	952	50345	68	204	9248	476	15155
16	111	333	24642	777	31128	55	165	6050	385	8900
<b>19</b>	<b>93</b>	<b>279</b>	<b>17298</b>	<b>651</b>	<b>23541</b>	<b>47</b>	<b>141</b>	<b>4418</b>	<b>329</b>	<b>7160</b>
22	81	243	13122	567	19203	40	120	3200	280	5540
25	71	213	10082	497	16028	36	108	2592	252	4482
28	63	189	7938	441	10896	32	96	2048	224	3341
31	57	171	6498	399	9375	29	87	1682	203	2900
34	52	156	5408	364	8219	26	78	1352	182	2327
37	48	144	4608	336	7365	24	72	1152	168	1980
40	45	135	4050	315	6765	22	66	968	154	1661
43	41	123	3362	287	5822	21	63	882	147	1512
46	39	117	3042	273	5265	20	60	800	140	1370
49	36	108	2592	252	4482	18	54	648	126	1107
<b>52</b>	<b>34</b>	<b>102</b>	<b>2312</b>	<b>238</b>	<b>3995</b>	<b>17</b>	<b>51</b>	<b>578</b>	<b>119</b>	<b>986</b>

Number of operations in one half-external product (i.e. one relinearization, one keyswitch, or one automorphism) and in one BFV/CKKS internal product (IP). Complete version of Table 1.

## D Forward noise propagation theorems

All the theorems about the external products in the main Section have been provided in a backward propagation direction, which is the most useful for compilation purposes: we consider the full plaintext circuit and assign nodes and levels to each elementary operation one by one starting from the end. Each time, we start from the targeted output noise, and deduce the (larger) minimal required noise levels for the products. Whenever the level become too high (i.e. not efficient, or noise too small to be secure) , we can place a bootstrapping and continue backwards at a lower level.

In a more traditional interpreted language scenario, we already have existing ciphertexts and pre-existing relinearization keys, that we must just feed-in as input to a homomorphic product, at their current level of noise. In this case, we can use forward propagation formulae to deduce the output noise. In a real circuit, these forward-propagation formulae end-up forming a large system of constraints that is usually non-linear, and are quite hard to optimize by the compilers, compared to their backward-propagation counterpart. However, they are in general also easier to understand for a human being, so we provide them in this section.

Also, the main section theorems use the worst-case propagation, based on the infinity norm of the noise. These formulae are inconditionnally true, but predict a noise growth that is in bits twice as fast as the reality. As explained in Section 4, the practical noise growth can be easily obtained if we use the additional assumption that all `bivRLWE` noises are samples from independent Gaussian distributions of same variance, and they remain independent. Such independance can be enforced randomizing ciphertexts between each external or internal product. In practice, a very mild randomization are in general sufficient to turn exceptionally large correllations in dot products into the average-case behaviour.

This is the analogue of Theorem 3.1 with forward average-case and worst-case noise propagation:

**Theorem D.1 (half external product - forward propatation (secret  $\times$  public)).**

Let  $K, \tilde{K}$  be limb sizes and let  $s \in \mathcal{R}$  a small key. Let  $C_u$  be a `bivHalfRGSW`<sup>DFT</sup> <sub>$s, K, \tilde{K}$</sub>  of an integer polynomial  $u \in \mathcal{R}$ , composed of  $\tilde{\ell}$  ciphertexts of  $\ell$  limbs, satisfying one of these two bounds:

- noise variance  $V_{GSW}$ , (average case)
- noise infinity norm  $\varepsilon_{GSW}$  (worst-case)

and let  $a = \sum (A_i)Y^i \in \mathcal{R}[Y]^2$  be  $\tilde{K}$ -normalized.

$$C_f \square a = \text{normalizeReduce} \left( \text{iDFT} \left( \sum_{i=1}^{\tilde{\ell}} \text{DFT}(A_i) \cdot c_i \right) \right)$$

is a `bivRLWE` <sub>$s, K$</sub>  of  $u \cdot \phi_K(a)$  with, depending on which input noise condition above is satisfied:

- output noise variance  $V_{out} \leq \tilde{\ell} N 2^{2\tilde{K}-2} V_{GSW} + \|u\|_2^2 \varepsilon_{Dec}^2$ , (average case)
- noise infinity norm  $\varepsilon_{out} \leq \tilde{\ell} N 2^{\tilde{K}-1} \varepsilon_{GSW} + \|u\|_1 \varepsilon_{Dec}$  (worst-case)

where  $\varepsilon_{Dec} = 2^{-\tilde{K}\tilde{\ell}-1}$  when the degree of  $a$  is  $\leq \tilde{\ell}$  (approx decomposition), otherwise zero (exact decomposition).

Such computation requires

- $\tilde{\ell}$  bounded DFT's of the  $A_i$ 's of norm  $\leq 2^{\tilde{K}-1}$ ,
- $2\tilde{\ell}$  element-wise **admmul**'s in DFT domain for polynomials of norm  $\leq N\tilde{\ell}2^{\tilde{K}+K-2}$ ,
- $2\tilde{\ell}$  bounded DFT's for the results  $C_j(X)$  with norm bound  $N\tilde{\ell}2^{\tilde{K}+K-2}$ ,
- $2\tilde{\ell}$  element-wise additions/shifts/masks to normalize and truncate the result.

And this is the analogue of Corollary 3.2, much closer to the traditional presentation in [10], again with worst-case and average-case forward-propagation of noise.

**Corollary D.1 (full external product- forward propagation (secret  $\times$  secret)).**

Let  $K, \tilde{K}_A, \tilde{K}_B$  be limb sizes, let  $u \in \mathcal{R}$  be a small polynomial, let  $v \in \mathbb{T}_N[X]$  be a message, let  $s \in \mathcal{R}$  a small key, and

- Let  $C_u$  be a  $\mathbf{bivHalfRGSW}_{s,K,\tilde{K}_B}^{\text{DFT}}$  of  $u \in \mathcal{R}$ , composed of  $\tilde{\ell}_B$  ciphertexts of  $\ell$  limbs, and of noise variance (resp. amplitude) bounded by  $V_A$  (resp.  $\varepsilon_A$ )
- Let  $C_{su}$  be a  $\mathbf{bivHalfRGSW}_{s,K,\tilde{K}_A}^{\text{DFT}}$  of  $su \in \mathcal{R}$ , composed of  $\tilde{\ell}_A$  ciphertexts of  $\ell$  limbs, and of noise variance (resp. amplitude) bounded by  $V_A$  (resp.  $\varepsilon_A$ )
- $(A, B) \in \mathcal{R}[Y]^2$  is a  $\mathbf{bivRLWE}_{s,K}$  ciphertext, of noise variance (resp. amplitude) bounded by  $V_{in}$  (resp.  $\varepsilon_{in}$ ).

Then

$$(C_u, C_{su}) \square (a, b) = \mathbf{normalizeRed} \left( \mathbf{iDFT} \left( \sum_{i=1}^{\tilde{\ell}_B} \mathbf{DFT}(B_i) \cdot c_i - \sum_{i=1}^{\tilde{\ell}_A} \mathbf{DFT}(A_i) \cdot d_i \right) \right),$$

is a  $\mathbf{bivRLWE}_{S,K,2-L}$  encryption of  $u \cdot v$ . Here,  $A$  and  $B$  have been  $\tilde{K}_A$  and  $\tilde{K}_B$ -normalized, and truncated to degree respectively  $\tilde{\ell}_A$  and  $\tilde{\ell}_B$ .

Let  $\varepsilon_{out}$  and  $V_{out}$  denote the amplitude and variance of the output ciphertext, these bounds are satisfied:

$$\begin{aligned} \varepsilon_{out} &\leq \varepsilon_{in} + \tilde{\ell}_A N 2^{\tilde{K}_A-1} \varepsilon_A + \|us\|_1 2^{-\tilde{\ell}_A \tilde{K}_A-1} + \tilde{\ell}_B N 2^{\tilde{K}_B-1} \varepsilon_B + \|u\|_1 2^{-\tilde{\ell}_B \tilde{K}_B-1} \\ V_{out} &\leq V_{in} + \tilde{\ell}_A N 4^{\tilde{K}_A-1} V_A + \|us\|_2^2 4^{-\tilde{\ell}_A \tilde{K}_A-1} + \tilde{\ell}_B N 4^{\tilde{K}_B-1} V_B + \|u\|_2^2 4^{-\tilde{\ell}_B \tilde{K}_B-1} \end{aligned}$$

computing this encryption with  $\ell = \lceil (L'_2 + \log_2 N + 2)/K \rceil$  requires at most

- $\tilde{\ell}_A + \tilde{\ell}_B$  bounded DFT's of the  $A_i, B_i$ 's of norm  $\leq 2^{\tilde{K}-1}$ ,
- $2\ell(\tilde{\ell}_A + \tilde{\ell}_B)$  element-wise **admmul**'s in DFT domain for polynomials of norm  $\leq 2N\tilde{\ell}2^{\tilde{K}+K-2}$ ,
- $2\ell$  bounded DFT's for the results  $C_j(X)$  with norm bound  $2N\tilde{\ell}2^{\tilde{K}+K-2}$ ,
- $\tilde{\ell}_A + \tilde{\ell}_B + 2\ell$  element-wise additions/shifts or masks to normalize and truncate the input ciphertext and the final result.