

# A Full RNS Variant of FV like Somewhat Homomorphic Encryption Schemes

Jean-Claude Bajard<sup>1</sup>, Julien Eynard<sup>2</sup>, M. Anwar Hasan<sup>2</sup>, and Vincent Zucca<sup>1</sup>

<sup>1</sup>Sorbonne Universités, UPMC, CNRS, LIP6, Paris, France.

{jean-claude.bajard, vincent.zucca}@lip6.fr

<sup>2</sup>Department of Electrical and Computer Engineering, University of Waterloo.

{jeynard, ahasan}@uwaterloo.ca

**Abstract.** Since Gentry’s breakthrough work in 2009, homomorphic cryptography has received a widespread attention. Implementation of a fully homomorphic cryptographic scheme is however still highly expensive. Somewhat Homomorphic Encryption (SHE) schemes, on the other hand, allow only a limited number of arithmetical operations in the encrypted domain, but are more practical. Many SHE schemes have been proposed, among which the most competitive ones rely on Ring Learning With Errors (RLWE) and operations occur on high-degree polynomials with large coefficients. This work focuses in particular on the Chinese Remainder Theorem representation (a.k.a. Residue Number Systems) applied to the large coefficients. In SHE schemes like that of Fan and Vercauteren (FV), such a representation remains hardly compatible with procedures involving coefficient-wise division and rounding required in decryption and homomorphic multiplication. This paper suggests a way to entirely eliminate the need for multi-precision arithmetic, and presents techniques to enable a full RNS implementation of FV-like schemes. For dimensions between  $2^{11}$  and  $2^{15}$ , we report speed-ups from  $5\times$  to  $20\times$  for decryption, and from  $2\times$  to  $4\times$  for multiplication.

**Keywords:** Lattice-based Cryptography, Homomorphic Encryption, FV, Residue Number Systems, Software Implementation

## 1 Introduction

Cryptographers’ deep interests in lattices are for multiple reasons. Besides possessing highly desirable post-quantum security features, lattice-based cryptography relies on simple structures, allowing efficient asymptotic complexities, and is quite flexible in practice. In addition to encryption/signature schemes ([15, 25, 18, 7, 21, 22]), identity-based encryption [8], multilinear maps [10, 16], lattices are also involved in homomorphic encryption (HE). The discovery of this property by Gentry in 2009 [12], through the use of ideal rings, is a major breakthrough which has opened the door to many opportunities in terms of applications, especially when coupled with cloud computing.

HE is generally composed of a basic layer, which is a Somewhat Homomorphic Encryption scheme (SHE). Such a scheme allows us to compute a limited

number of additions and multiplications on ciphertexts. This can be explained by the fact that any ciphertext contains an inherent noise which increases after each homomorphic operation. Beyond a certain limit, the noise becomes too large to allow a correct decryption. This drawback may be tackled by using bootstrapping, which however constitutes a bottleneck in terms of efficiency. Further improvements of noise management [6, 5] have been suggested so that, in practice, and given an applicative context, it may be wiser to select an efficient SHE with parameters enabling a sufficient number of operations. For instance, schemes like **FV** [9] and **YASHE** [4] have been implemented and tested for evaluating the SIMON Feistel Cipher [17]. Among the today’s more practical SHE schemes, **FV** is arguably one of the most competitive. This scheme is being currently considered by major stakeholders, such as the European H2020 HEAT consortium [26], as a viable candidate for practical homomorphic encryption.

**Our Contribution** This work focuses on practical improvement of SHE schemes, in particular **FV**. Despite the fact that the security of **YASHE** has been called into question recently [1], this scheme can also benefit from the present work. These schemes handle elements of a polynomial ring  $\mathbb{Z}_q[X]/(X^n + 1)$ . The modulus  $q$  can be chosen as the product of some small moduli fitting with practical hardware requirements (machine word, etc.). This enables us to avoid the need of multi-precision arithmetic in almost the whole scheme. However, this CRT representation (a.k.a. Residue Number Systems, or RNS) is hardly compatible with a couple of core operations: coefficient-wise division and rounding, occurring in multiplication and decryption, and a noise management technique within homomorphic multiplication, relying on the access to a positional number system.

We show how to efficiently avoid any switch between RNS and the positional system for performing these operations. We present a full RNS variant of **FV** and analyse the new bounds on noise growth. A software implementation highlights the practical benefits of the new RNS variant.

It is important to note that this work is related to the arithmetic at the coefficient level. Thus, the security features of the original scheme are not modified.

**Outline** Section 2 provides some preliminaries about **FV** and RNS. Section 3 provides a full RNS variant of decryption. Section 4 gives a full RNS variant of homomorphic multiplication. Results of a software implementation are presented in Section 5. Finally, some conclusions are drawn.

## 2 Preliminaries

For the sake of clarity, proofs of lemmas, propositions and theorems within Sect. 3 and 4 are provided in an extended version [2].

**Context** High-level operations occur in a polynomial ring  $\mathcal{R} = \mathbb{Z}[X]/(X^n + 1)$  with  $n$  a power of 2.  $\mathcal{R}$  is one-to-one mapped to integer polynomials of degree

strictly smaller than  $n$ . Most of the time, elements of  $\mathcal{R}$  are denoted by lower-case boldface letters and identified by their coefficients. Polynomial arithmetic is done modulo  $(X^n + 1)$ . The infinity norm of  $\mathbf{a} = (a_0, \dots, a_{n-1}) \in \mathcal{R}$  is defined by  $\|\mathbf{a}\|_\infty = \max_{0 \leq i \leq n-1} (|a_i|)$ . Ciphertexts will be managed as polynomials (of degree 1) in  $\mathcal{R}[Y]$ . For  $\mathbf{ct} \in \mathcal{R}[Y]$ , we note  $\|\mathbf{ct}\|_\infty = \max_i \|\mathbf{ct}[i]\|_\infty$ ,  $\mathbf{ct}[i]$  being the coefficient of degree  $i$  in  $Y$ . The multiplicative law of  $\mathcal{R}[Y]$  is denoted by  $\star$ .

Behind lattice-based cryptosystems in general, and FV in particular, lies the principle of noisy encryption. Additionally to the plaintext, a ciphertext contains a noise (revealed by using the secret key) which grows after each homomorphic operation. Since the homomorphic multiplication involves multiplications in  $\mathcal{R}$ , it is crucial that the size of a product in  $\mathcal{R}$  does not increase too much. This increase is related to the ring constant  $\delta = \sup\{\|\mathbf{f} \cdot \mathbf{g}\|_\infty / \|\mathbf{f}\|_\infty \cdot \|\mathbf{g}\|_\infty : (\mathbf{f}, \mathbf{g}) \in (\mathcal{R} \setminus \{\mathbf{0}\})^2\}$ . It means that  $\|\mathbf{f} \cdot \mathbf{g}\|_\infty \leq \delta \|\mathbf{f}\|_\infty \cdot \|\mathbf{g}\|_\infty$ . For the specific ring  $\mathcal{R}$  used here,  $\delta$  is equal to  $n$ .

For our subsequent discussions on decryption and homomorphic multiplication, we denote the ‘‘Division and Rounding’’ in  $\mathcal{R}[Y]$  (depending on parameters  $t, q$  which are defined thereafter) by:

$$\text{DR}_i : \mathbf{ct} = \sum_{j=0}^i \mathbf{ct}[j]Y^j \in \mathcal{R}[Y] \mapsto \sum_{j=0}^i \left\lfloor \frac{t}{q} \mathbf{ct}[j] \right\rfloor Y^j \in \mathcal{R}[Y] . \quad (1)$$

The notation  $\lfloor \frac{t}{q} \mathbf{c} \rfloor$ , for any  $\mathbf{c} \in \mathcal{R}$  (e.g.  $\mathbf{ct}[j]$  in (1)), means a coefficient-wise division-and-rounding.

**Plaintext and Ciphertext Spaces** The plaintext space is determined by an integer parameter  $t$  ( $t \geq 2$ ). A message is an element of  $\mathcal{R}_t = \mathcal{R}/(t\mathcal{R})$ , i.e. a polynomial of degree at most  $n - 1$  with coefficients in  $\mathbb{Z}_t$ . The notation  $[\mathbf{m}]_t$  (resp.  $[\mathbf{m}]_t$ ) means that coefficients lie in  $[-t/2, t/2)$  (resp.  $[0, t)$ ). Ciphertexts will lie in  $\mathcal{R}_q[Y]$  with  $q$  a parameter of the scheme. On one side, some considerations about security imply a relationship between  $q$  and  $n$  which, for a given degree  $n$ , establish an upper bound for  $\log_2(q)$  (cf. (6) in [9]). On the other side, the ratio  $\Delta = \lfloor \frac{q}{t} \rfloor$  will basically determine the maximal number of homomorphic operations which can be done in a row to ensure a correct decryption.

**RNS Representation** Beyond the upper bound on  $\log_2(q)$  due to security requirements, the composition of  $q$  has no restriction. So,  $q$  can be chosen as a product of small pairwise coprime moduli  $q_1 \dots q_k$ . The reason for such a choice is the Chinese Remainder Theorem (CRT) which offers a ring isomorphism  $\mathbb{Z}_q \xrightarrow{\sim} \prod_{i=1}^k \mathbb{Z}_{q_i}$ . Thus, the CRT implies the existence of a non-positional number system (RNS) in which large integers (mod  $q$ ) are mapped to sets of small residues. Furthermore, the arithmetic modulo  $q$  over large integers can be substituted by  $k$  independent arithmetics in the small rings  $\mathbb{Z}_{q_i}$ . The isomorphism can be naturally extended to polynomials:  $\mathcal{R}_q \simeq \mathcal{R}_{q_1} \times \dots \times \mathcal{R}_{q_k}$ . It means that RNS can be used at the coefficient level to accelerate the arithmetic in  $\mathcal{R}_q$ .

In the rest of the paper, the letter  $q$  may refer either to the product  $q_1 \dots q_k$  or to the “RNS base”  $\{q_1, \dots, q_k\}$ . Symbol  $\nu$  denotes the “width” of the moduli. In other words, from now on, any modulus  $m$  (should it belong to  $q$  or to any other RNS base) is assumed to verify  $m < 2^\nu$ .

**Asymmetric Keys** The *secret key*  $\mathbf{s}$  is picked up in  $\mathcal{R}$  according to a discrete distribution  $\chi_{key}$  on  $\mathcal{R}$  (in practice, bounded by  $B_{key} = 1$ , i.e.  $\|\mathbf{s}\|_\infty \leq 1$ ).

For creating the public key, an “error” distribution  $\chi_{err}$  over  $\mathcal{R}$  is used. In practice, this is a discrete distribution statistically close to a Gaussian (with mean 0 and standard deviation  $\sigma_{err}$ ) truncated at  $B_{err}$  (e.g.  $B_{err} = 6\sigma_{err}$ ).  $\chi_{err}$  is related to the hardness of the underlying (search version of) RLWE problem (for which the purpose is, given samples  $([-(\mathbf{a}_i \mathbf{s} + \mathbf{e}_i)]_q, \mathbf{a}_i)$  with  $\mathbf{e}_i \leftarrow \chi_{err}$  and  $\mathbf{a} \leftarrow \mathcal{U}(\mathcal{R}_q)$ , to find  $\mathbf{s}$ ;  $\mathcal{U}(\mathcal{R}_q)$  is the uniform distribution on  $\mathcal{R}_q$ ). The *public key*  $\mathbf{pk}$  is created as follows: sample  $\mathbf{a} \leftarrow \mathcal{U}(\mathcal{R}_q)$  and  $\mathbf{e} \leftarrow \chi_{key}$ , then output  $\mathbf{pk} = (\mathbf{p}_0, \mathbf{p}_1) = ([-(\mathbf{a} \mathbf{s} + \mathbf{e})]_q, \mathbf{a})$ .

**Encryption, Addition, Inherent Noise of a Ciphertext** Encryption and homomorphic addition are already fully compliant with RNS arithmetic. They are recalled hereafter:

- $\text{Enc}_{\text{FV}}([\mathbf{m}]_t)$ : from samples  $\mathbf{e}_1, \mathbf{e}_2 \leftarrow \chi_{err}$ ,  $\mathbf{u} \leftarrow \chi_{key}$ , output  $\mathbf{ct} = (\mathbf{ct}[0], \mathbf{ct}[1]) = ([\Delta[\mathbf{m}]_t + \mathbf{p}_0 \mathbf{u} + \mathbf{e}_1]_q, [\mathbf{p}_1 \mathbf{u} + \mathbf{e}_2]_q)$ .
- $\text{Add}_{\text{FV}}(\mathbf{ct}_1, \mathbf{ct}_2)$ : output  $([\mathbf{ct}_1[0] + \mathbf{ct}_2[0]]_q, [\mathbf{ct}_1[1] + \mathbf{ct}_2[1]]_q)$ .

By definition, the *inherent noise* of  $\mathbf{ct}$  (encrypting  $[\mathbf{m}]_t$ ) is the polynomial  $\mathbf{v}$  such that  $[\mathbf{ct}(\mathbf{s})]_q = [\mathbf{ct}[0] + \mathbf{ct}[1]\mathbf{s}]_q = [\Delta[\mathbf{m}]_t + \mathbf{v}]_q$ . Thus, it is revealed by evaluating  $\mathbf{ct} \in \mathcal{R}_q[Y]$  on the secret key  $\mathbf{s}$ .

**Elementary Operations** A basic word will fit in  $\nu$  bits. In RNS, an “inner modular multiplication” (IMM) in a small ring like  $\mathbb{Z}_m$  is a core operation. If EM stands for an elementary multiplication of two words, in practice an IMM is costlier than an EM. But it can be well controlled. For instance, the moduli provided in `NFLib` library [19] (cf. Sect. 5) enable a modular reduction which reduces to one EM followed by a multiplication modulo  $2^\nu$ . Furthermore, the cost of an inner reduction can be limited by using lazy reduction, e.g. during RNS base conversions used throughout this paper. NTT and `invNTT` denote the Number Theoretic Transform and its inverse in a ring  $\mathcal{R}_m$  for a modulus  $m$ . They enable an efficient polynomial multiplication ( $\text{NTT}, \text{invNTT} \in \mathcal{O}(n \log_2(n))$ ).

### 3 Towards a Full RNS Decryption

This section deals with the creation of a variant of the original decryption function  $\text{Dec}_{\text{FV}}$ , which will only involve RNS representation. The definition of  $\text{Dec}_{\text{FV}}$  is recalled hereafter.

–  $\text{Dec}_{\text{FV}}(\text{ct} = (c_0, c_1) \in \mathcal{R}_q[Y])$ : compute  $[\text{DR}_0([\text{ct}(\mathbf{s})]_q)]_t = \left[ \left[ \frac{t}{q} [c_0 + c_1 \mathbf{s}]_q \right] \right]_t$ .

The idea is that computing  $[c_0 + c_1 \mathbf{s}]_q = [\Delta[\mathbf{m}]_t + \mathbf{v}]_q$  reveals the noise. If this noise is small enough, and given that  $[\mathbf{m}]_t$  has been scaled by  $\Delta$ , the function  $\text{DR}_0$  allows to cancel the noise while scaling down  $\Delta[\mathbf{m}]_t$  to recover  $[\mathbf{m}]_t$ . Concretely, decryption is correct as long as  $\|\mathbf{v}\|_\infty < (\Delta - |q|_t)/2$ , i.e. the size of the noise should not go further this bound after homomorphic operations.

The division-and-rounding operation makes  $\text{Dec}_{\text{FV}}$  hardly compatible with RNS at a first sight. Because RNS is of non positional nature, only exact integer division can be naturally performed (by multiplying by a modular inverse). But it is not the case here. And the rounding operation involves comparisons which require to switch from RNS to another positional system anyway, should it be a classical binary system or a mixed-radix one [11]. To get an efficient RNS variant of  $\text{Dec}_{\text{FV}}$ , we use an idea of [3]. To this end, we introduce relevant RNS tools.

### 3.1 Fast RNS Base Conversion

At some point, the decryption requires, among others, a polynomial to be converted from  $\mathcal{R}_q$  to  $\mathcal{R}_t$ . To achieve such kind of operations as efficiently as possible, we suggest to use a “fast base conversion”. In order to convert residues of  $x \in [0, q)$  from base  $q$  to a base  $\mathcal{B}$  (e.g.  $\{t\}$ ) coprime to  $q$ , we compute:

$$\text{FastBconv}(x, q, \mathcal{B}) = \left( \sum_{i=1}^k \left| x_i \frac{q_i}{q} \right|_{q_i} \times \frac{q}{q_i} \bmod m \right)_{m \in \mathcal{B}}. \quad (2)$$

This conversion is relatively fast. This is because the sum should ideally be reduced modulo  $q$  in order to provide the exact value  $x$ ; instead, (2) provides  $x + \alpha_x q$  for some integer  $\alpha_x \in [0, k-1]$ . Computing  $\alpha_x$  requires costly operations in RNS. So this step is by-passed, at the cost of an approximate result.

$\text{FastBconv}$  naturally extends to polynomials of  $\mathcal{R}$  by applying it coefficient-wise.

### 3.2 Approximate RNS Rounding

The above mentioned fast conversion allows us to efficiently compute an approximation of  $\lfloor \frac{t}{q} [c_0 + c_1 \mathbf{s}]_q \rfloor$  modulo  $t$ . The next step consists in correcting this approximation.

First, we remark that  $|\text{ct}(\mathbf{s})|_q$  can be used instead of  $[\text{ct}(\mathbf{s})]_q$ . Indeed, the difference between these two polynomials is a multiple of  $q$ . So, the division-and-rounding turns it into a polynomial multiple of  $t$ , which is cancelled by the last reduction modulo  $t$ . Second, a rounding would involve, at some point, a comparison. This is hardly compatible with RNS, so it is avoided. Therefore, we propose to simplify the computation, albeit at the price of possible errors, by replacing rounding by flooring. To this end, we use the following formula:

$$\left\lfloor \frac{t}{q} |\text{ct}(\mathbf{s})|_q \right\rfloor = \frac{t |\text{ct}(\mathbf{s})|_q - |t \cdot \text{ct}(\mathbf{s})|_q}{q}.$$

The division is now exact, so it can be done in RNS. Since this computation has to be done modulo  $t$ , the term  $t|\mathbf{ct}(\mathbf{s})|_q$  cancels. Furthermore, the term  $(|t \cdot \mathbf{ct}(\mathbf{s})|_q \bmod t)$  is obtained through a fast conversion.

Lemma 1 sums up the strategy by replacing  $|\mathbf{ct}(\mathbf{s})|_q$  by  $\gamma|\mathbf{ct}(\mathbf{s})|_q$ , where  $\gamma$  is an integer which will help in correcting the approximation error.

**Lemma 1.** *Let  $\mathbf{ct}$  be such that  $[\mathbf{ct}(\mathbf{s})]_q = \Delta[\mathbf{m}]_t + \mathbf{v} + q\mathbf{r}$ , and let  $\mathbf{v}_c := t\mathbf{v} - [\mathbf{m}]_t|_q|_t$ . Let  $\gamma$  be an integer coprime to  $q$ . Then, for  $m \in \{t, \gamma\}$ , the following equalities hold modulo  $m$ :*

$$\begin{aligned} \text{FastConv}(|\gamma t \cdot \mathbf{ct}(\mathbf{s})|_q, q, \{t, \gamma\}) \times | -q^{-1}|_m &= \left\lfloor \gamma \frac{t}{q} [\mathbf{ct}(\mathbf{s})]_q \right\rfloor - \mathbf{e} \\ &= \gamma([\mathbf{m}]_t + t\mathbf{r}) + \left\lfloor \gamma \frac{\mathbf{v}_c}{q} \right\rfloor - \mathbf{e} \end{aligned} \quad (3)$$

where each integer coefficient of the error polynomial  $\mathbf{e} \in \mathcal{R}$  lies in  $[0, k]$ .

The error  $\mathbf{e}$  is due to the fast conversion and the replacement of rounding by flooring. It is the same error for residues modulo  $t$  and  $\gamma$ . The residues modulo  $\gamma$  will enable a fast correction of it and of the term  $\lfloor \gamma \frac{\mathbf{v}_c}{q} \rfloor$  at a same time. Also, note that  $\mathbf{r}$  vanishes since it is multiplied by both  $t$  and  $\gamma$ .

### 3.3 Correcting the Approximate RNS Rounding

The next step is to show how  $\gamma$  in (3) can be used to correct the term  $(\lfloor \gamma \frac{\mathbf{v}_c}{q} \rfloor - \mathbf{e})$ . It can be done efficiently when the polynomial  $\mathbf{v}_c$  is such that  $\|\mathbf{v}_c\|_\infty \leq q(\frac{1}{2} - \varepsilon)$ , for some real number  $\varepsilon \in (0, 1/2]$ .

**Lemma 2.** *Let  $\|\mathbf{v}_c\|_\infty \leq q(\frac{1}{2} - \varepsilon)$ ,  $\mathbf{e} \in \mathcal{R}$  with coefficients in  $[0, k]$ , and  $\gamma$  an integer. Then,*

$$\gamma\varepsilon \geq k \Rightarrow \left[ \left\lfloor \gamma \frac{\mathbf{v}_c}{q} \right\rfloor - \mathbf{e} \right]_\gamma = \left\lfloor \gamma \frac{\mathbf{v}_c}{q} \right\rfloor - \mathbf{e} . \quad (4)$$

Lemma 2 enables an efficient and correct RNS rounding as long as  $k(\frac{1}{2} - \frac{\|\mathbf{v}_c\|_\infty}{q})^{-1} \sim \gamma$  has the size of a modulus. Concretely, one computes (3) and uses the centered remainder modulo  $\gamma$  to obtain  $\gamma([\mathbf{m}]_t + t\mathbf{r})$  modulo  $t$ , which reduces to  $\gamma[\mathbf{m}]_t \bmod t$ . And it remains to multiply by  $|\gamma^{-1}|_t$  to recover  $[\mathbf{m}]_t$ .

### 3.4 A Full RNS Variant of Dec<sub>FV</sub>

The new variant of the decryption is detailed in Alg. 1. The main modification for the proposed RNS decryption is due to Lem. 2. As stated by Thm. 1, given a  $\gamma$ , the correctness of rounding requires a new bound on the noise to make the  $\gamma$ -correction technique successful.

**Theorem 1.** *Let  $\text{ct}(\mathbf{s}) = \Delta[\mathbf{m}]_t + \mathbf{v} \pmod{q}$ . Let  $\gamma$  be a positive integer coprime to  $t$  and  $q$  such that  $\gamma > 2k/(1 - \frac{t|q|_t}{q})$ . For Alg. 1 returning  $[\mathbf{m}]_t$ , it suffices that  $\mathbf{v}$  satisfies the following bound:*

$$\|\mathbf{v}\|_\infty \leq \frac{q}{t} \left( \frac{1}{2} - \frac{k}{\gamma} \right) - \frac{|q|_t}{2}. \quad (5)$$

There is a trade-off between the size of  $\gamma$  and the bound in (5). Ideally,  $\gamma \sim 2k$  at the price of a (*a priori*) quite small bound on the noise. But by choosing  $\gamma \sim 2^{p+1}k$  for  $p < \nu - 1 - \lceil \log_2(k) \rceil$  (i.e.  $\gamma < 2^\nu$  is a standard modulus), the bound  $(\Delta(1 - 2^{-p}) - |q|_t)/2$  for a correct decryption should be close to the original bound  $(\Delta - |q|_t)/2$  for practical values of  $\nu$ . A concrete estimation of  $\gamma$  in Sect. 5.1 will show that  $\gamma$  can be chosen very close to  $2k$  in practice, and thus fitting on a basic word by far.

---

**Algorithm 1**  $\text{Dec}_{\text{RNS}}(\text{ct}, \mathbf{s}, \gamma)$

---

**Require:**  $\text{ct}$  an encryption of  $[\mathbf{m}]_t$ , and  $\mathbf{s}$  the secret key, both in base  $q$ ; an integer  $\gamma$  coprime to  $t$  and  $q$

**Ensure:**  $[\mathbf{m}]_t$

- 1: **for**  $m \in \{t, \gamma\}$  **do**
  - 2:    $\mathbf{s}^{(m)} \leftarrow \text{FastBconv}(|\gamma t \cdot \text{ct}(\mathbf{s})|_q, q, \{m\}) \times | -q^{-1}|_m \pmod{m}$
  - 3: **end for**
  - 4:  $\tilde{\mathbf{s}}^{(\gamma)} \leftarrow [\mathbf{s}^{(\gamma)}]_\gamma$
  - 5:  $\mathbf{m}^{(t)} \leftarrow [(\mathbf{s}^{(t)} - \tilde{\mathbf{s}}^{(\gamma)}) \times |\gamma^{-1}|_t]_t$
  - 6: **return**  $\mathbf{m}^{(t)}$
- 

### 3.5 Staying in RNS Is Asymptotically Better

In any decryption technique,  $(\text{ct}(\mathbf{s}) \pmod{q})$  has to be computed first. To optimize this polynomial product, one basically performs  $k\text{NTT} \rightarrow kn\text{IMM} \rightarrow k\text{invNTT}$ . For next steps, a simple strategy is to compute  $(\lfloor \frac{t}{q} [\text{ct}(\mathbf{s}) \rfloor_q \pmod{t})$  by doing an RNS-to-binary conversion in order to perform the division and rounding. By denoting  $\mathbf{x}_i = |\text{ct}(\mathbf{s}) \frac{q_i}{q}|_{q_i}$ , one computes  $\sum_{i=1}^k \mathbf{x}_i \frac{q}{q_i} \pmod{q}$ , compares it to  $q/2$  so as to center the result, and performs division and rounding. Hence, the division-and-rounding requires  $\mathcal{O}(k^2 n)\text{EM}$ . In practice, security analysis (cf. e.g. [9, 4, 17]) requires that  $k\nu = \lceil \log_2(q) \rceil \in \mathcal{O}(n)$ . So, the cost of leaving RNS to access a positional system is dominant in the asymptotic computational complexity.

Staying in RNS enables to get a better asymptotic complexity. Indeed, it is easy to see that Alg. 1 requires  $\mathcal{O}(kn)$  operations (excluding the polynomial product). Thus, the cost of NTT is dominant in this case. By considering  $k \in \mathcal{O}(n)$ , we deduce  $\mathcal{C}(\text{Dec}_{\text{FV}}) \in \mathcal{O}(n^3)$ , while  $\mathcal{C}(\text{Dec}_{\text{RNS}}) \in \mathcal{O}(n^2 \log_2(n))$ . But the hidden constant in “ $k \in \mathcal{O}(n)$ ” is small, and the NTT, common to both variants, should avoid any noticeable divergence (cf. 5.4) for practical ranges for parameters.

In order to provide optimized RNS variants of decryption, we make two remarks. First, the reduction modulo  $q$  is unnecessary. Indeed, any extra multiple of  $q$  in the sum  $\sum_{i=1}^k \mathbf{x}_i \frac{q}{q_i}$  is multiplied by  $\frac{t}{q}$ , making the resulting term a multiple of  $t$ , which is not affected by the rounding and is finally cancelled modulo  $t$ . Second, it is possible to precompute  $\frac{t}{q}$  as a multiprecision floating point number in order to avoid a costly integer division. But given the first remark, it suffices to precompute the floating point numbers  $Q_i \sim \frac{t}{q_i}$  with  $2\nu + \log_2(k) - \log_2(t)$  bits ( $\sim 2$  words) of precision. In this case, using standard double or quadruple (depending on  $\nu$ ) precision is sufficient. Finally, it is sufficient to compute  $\lfloor \sum_{i=1}^k \mathbf{x}_i Q_i \rfloor \bmod t$ . This represents about  $2kn\text{EM}$ . Reducing modulo  $t$  is nearly free of cost when  $t$  is a power of 2.

A second optimized RNS variant, with only integer arithmetic, is based on Alg. 1, in which  $\gamma$  is assumed to be coprime to  $t$ . It is possible to be slightly more efficient by noticing that the coprimality assumption can be avoided. This is because the division by  $\gamma$  is exact. To do it, the `for` loop can be done modulo  $\gamma \times t$ . For instance, even if  $t$  is a power of 2, one can choose  $\gamma$  as being a power of 2 too, and use the following lemma to finish the decryption efficiently.

**Lemma 3.** *Let  $\gamma$  be a power of 2. Let  $\mathbf{z} := \lfloor \gamma \lfloor \mathbf{m} \rfloor_t + \lfloor \gamma \frac{v_c}{q} \rfloor - e \rfloor_{\gamma t}$  coming from (3) when computed modulo  $\gamma t$ . If  $\gamma$  satisfies (4), then ( $\gg$  denotes the right bit-shifting, and  $\mathbf{v}_1$  is the polynomial with all its coefficients being equal to 1)*

$$\left[ (\mathbf{z} + \frac{\gamma}{2} \mathbf{v}_1) \gg \log_2(\gamma) \right]_t = \lfloor \mathbf{m} \rfloor_t . \quad (6)$$

Lemma 3 can be adapted to other values for  $\gamma$ . The important remark is that  $\lfloor \mathbf{m} \rfloor_t$  is contained in the  $\lfloor \log_2(t) \rfloor + 1$  most significant bits of  $(\mathbf{z} + \frac{\gamma}{2} \mathbf{v}_1) \bmod \gamma t$ . So, by choosing  $\gamma$  as a power of 2, a simple bit shifting enables to recover these bits. Finally, as soon as  $\gamma t$  fits in 1 word, the cost of such variant (besides the polynomial product) reduces to  $kn\text{IMM}$ , or simply to  $kn\text{EM}$  modulo  $2^{\log_2(\gamma t)}$  whenever  $t$  is a power of 2.

*Remark 1.* In previous discussion, the product  $\gamma t$  is assumed fitting in one machine word to simplify complexity analysis. However, for some applications, the plaintext modulus  $t$  can be bigger than a machine word (e.g. homomorphic neural networks [13], where  $t > 2^{80}$ ). In such cases, either the plaintexts directly lie in  $\mathcal{R}_t$ , or  $t$  can be decomposed in a product of smaller moduli  $t_1, \dots, t_\ell$ , enabling the use of RNS for encoding plaintexts (and then allowing better homomorphic multiplicative depth for a given dimension  $n$ ). In the first case, the optimized RNS decryption (given by Lem. 3) remains available, but the residues modulo  $t$  should be handled with several words. In the second case, a plaintext is recovered by decrypting its residue modulo each of the  $t_i$ . These  $\ell$  decryptions can be done as in Lem. 3, by using  $\gamma$  as a power of 2 (whatever the  $t_i$ 's are). Finally, the plaintext is reconstructed from residues modulo the  $t_i$ 's by using a classical RNS to binary conversion. However, this conversion is only related to the way the plaintexts are encoded. This is not handled by RNS decryption described in this paper, which only deals with representation of ciphertexts (i.e. modulo  $q$ ).



## 4 Towards a Full RNS Homomorphic Multiplication

### 4.1 Preliminaries About $\text{Mult}_{\text{FV}}$

Below, we recall the main mechanisms of the homomorphic multiplication  $\text{Mult}_{\text{FV}}$  from [9]. More precisely, we focus on the variant with version 1 for relinearisation step. First, two functions, of which the purpose is to limit a too rapid noise growth during a multiplication, are recalled (these functions will be denoted as in [4]). They are applicable to any  $\mathbf{a} \in \mathcal{R}$ , for any radix  $\omega$ , and with the subsequent parameter  $\ell_{\omega,q} = \lceil \log_{\omega}(q) \rceil + 1$ .  $\mathcal{D}_{\omega,q}$  is a decomposition in radix base  $\omega$ , while  $\mathcal{P}_{\omega,q}$  gets back powers of  $\omega$  which are lost within the decomposition process.

$$\forall \mathbf{a} \in \mathcal{R}, \begin{cases} \mathcal{D}_{\omega,q}(\mathbf{a}) = ([\mathbf{a}]_{\omega}, [[\mathbf{a}\omega^{-1}]_{\omega}], \dots, [[\mathbf{a}\omega^{-(\ell_{\omega,q}-1)}]_{\omega}]) \in \mathcal{R}_{\omega}^{\ell_{\omega,q}} \\ \mathcal{P}_{\omega,q}(\mathbf{a}) = ([\mathbf{a}]_q, [\mathbf{a}\omega]_q, \dots, [\mathbf{a}\omega^{\ell_{\omega,q}-1}]_q) \in \mathcal{R}_q^{\ell_{\omega,q}} \end{cases} . \quad (7)$$

In particular, for any  $(\mathbf{a}, \mathbf{b}) \in \mathcal{R}^2$ ,  $\langle \mathcal{D}_{\omega,q}(\mathbf{a}), \mathcal{P}_{\omega,q}(\mathbf{b}) \rangle \equiv \mathbf{a}\mathbf{b} \pmod{q}$ .

Next,  $\text{Mult}_{\text{FV}}$  is built as follows ( $\text{rlk}_{\text{FV}}$  is a public relinearisation key):

- $\text{rlk}_{\text{FV}} = ([\mathcal{P}_{\omega,q}(s^2) - (\vec{e} + s\vec{a})]_q, \vec{a})$  where  $\vec{e} \leftarrow \chi_{err}^{\ell_{\omega,q}}$ ,  $\vec{a} \leftarrow \mathcal{U}(\mathcal{R}_q)^{\ell_{\omega,q}}$ ,
- $\text{Relin}_{\text{FV}}(\mathbf{c}_0, \mathbf{c}_1, \mathbf{c}_2, \text{rlk}_{\text{FV}})$ :  
compute  $([\mathbf{c}_0 + \langle \mathcal{D}_{\omega,q}(\mathbf{c}_2), \text{rlk}_{\text{FV}}[0] \rangle]_q, [\mathbf{c}_1 + \langle \mathcal{D}_{\omega,q}(\mathbf{c}_2), \text{rlk}_{\text{FV}}[1] \rangle]_q)$ ,
- $\text{Mult}_{\text{FV}}(\mathbf{ct}_1, \mathbf{ct}_2)$ : denote  $\mathbf{ct}_{\star} = \mathbf{ct}_1 \star \mathbf{ct}_2$  (degree-2 element of  $\mathcal{R}[Y]$ ),
  - Step 1:  $\widetilde{\mathbf{ct}}_{mult} = [\text{DR}_2(\mathbf{ct}_{\star})]_q = ([\text{DR}_0(\mathbf{ct}_{\star}[i])]_q)_{i \in \{0,1,2\}}$ ,
  - Step 2:  $\mathbf{ct}_{mult} = \text{Relin}_{\text{FV}}(\widetilde{\mathbf{ct}}_{mult})$ .

There are two main obstacles to a full RNS variant. First, the three calls to  $\text{DR}_0$  in Step 1, for which the context is different than for the decryption. While in the decryption we are working with a noise whose size can be controlled, and while we are reducing a value from  $q$  to  $\{t\}$ , here the polynomial coefficients of the product  $\mathbf{ct}_1 \star \mathbf{ct}_2$  have kind of random size modulo  $q$  (for each integer coefficient) and have to be reduced towards  $q$ . Second, the function  $\mathcal{D}_{\omega,q}$  (in  $\text{Relin}_{\text{FV}}$ ) requires, by definition, an access to a positional system (in radix base  $\omega$ ), which is hardly compatible with RNS.

### 4.2 Auxiliary RNS Bases

Step 1 requires to use enough moduli to contain any product, in  $\mathcal{R}[Y]$  (i.e. on  $\mathbb{Z}$ ), of degree-1 elements from  $\mathcal{R}_q[Y]$ . So, we need an auxiliary base  $\mathcal{B}$ , additionally to the base  $q$ . We assume that  $\mathcal{B}$  contains  $\ell$  moduli (while  $q$  owns  $k$  elements). A sufficient size for  $\ell$  will be given later. An extra modulus  $m_{\text{sk}}$  is added to  $\mathcal{B}$  to create an extended base  $\mathcal{B}_{\text{sk}}$ . It will be used for a transition between the new steps 1 and 2. Computing the residues of ciphertexts in  $\mathcal{B}_{\text{sk}}$  is done through a fast conversion from  $q$ . In order to reduce the extra multiples of  $q$  (called “ $q$ -overflows” in further discussions) this conversion can produce, a single-modulus base  $\tilde{m}$  is introduced. All these bases are assumed to be pairwise coprime.

**Reducing (mod  $q$ ) a Ciphertext in  $\mathcal{B}_{\text{sk}}$**  A `FastBconv` from  $q$  can create  $q$ -overflows (i.e. unnecessary multiples of  $q$ ) in the output. To limit the impact on noise growth (because of division by  $q$  in step 1), we give an efficient way to reduce a polynomial  $\mathbf{c} + q\mathbf{u}$  in  $\mathcal{B}_{\text{sk}}$ . It should be done prior to each multiplication. For that purpose, we use the residues modulo  $\tilde{m}$  as it is described in Alg. 2.

---

**Algorithm 2**  $\text{SmMR}_{\mathbf{q}_{\tilde{m}}}((\mathbf{c}''_m)_{m \in \mathcal{B}_{\text{sk}} \cup \{\tilde{m}\}})$ : Small Montgomery Reduction mod  $q$

---

**Require:**  $\mathbf{c}''$  in  $\mathcal{B}_{\text{sk}} \cup \{\tilde{m}\}$

**Ensure:**  $\mathbf{c}'$  in  $\mathcal{B}_{\text{sk}}$ , with  $\mathbf{c}' \equiv \mathbf{c}''\tilde{m}^{-1} \pmod{q}$ ,  $\|\mathbf{c}'\|_\infty \leq \frac{\|\mathbf{c}''\|_\infty}{\tilde{m}} + \frac{q}{2}$

1:  $\mathbf{r}_{\tilde{m}} \leftarrow \lfloor -\mathbf{c}''_{\tilde{m}}/q \rfloor_{\tilde{m}}$

2: **for**  $m \in \mathcal{B}_{\text{sk}}$  **do**

3:      $\mathbf{c}'_m \leftarrow \lfloor (\mathbf{c}''_m + q\mathbf{r}_{\tilde{m}})\tilde{m}^{-1} \rfloor_m$

4: **end for**

5: **return**  $\mathbf{c}'$  in  $\mathcal{B}_{\text{sk}}$

---

**Lemma 4.** *On input  $\mathbf{c}''_m = \lfloor \tilde{m}\mathbf{c} \rfloor_q + q\mathbf{u}|_m$  for all  $m \in \mathcal{B}_{\text{sk}} \cup \{\tilde{m}\}$ , with  $\|\mathbf{u}\|_\infty \leq \tau$ , and given a parameter  $\rho > 0$ , then Alg. 2 returns  $\mathbf{c}'$  in  $\mathcal{B}_{\text{sk}}$  with  $\mathbf{c}' \equiv \mathbf{c} \pmod{q}$  and  $\|\mathbf{c}'\|_\infty \leq \frac{q}{2}(1 + \rho)$  if  $\tilde{m}$  satisfies:*

$$\tilde{m}\rho \geq 2\tau + 1 . \quad (8)$$

To use this fast reduction, the ciphertexts have to be handled in base  $q$  through the Montgomery [20] representation with respect to  $\tilde{m}$  (i.e.  $\lfloor \tilde{m}\mathbf{c} \rfloor_q$  instead of  $\lfloor \mathbf{c} \rfloor_q$ ). This can be done for free of cost during the base conversions (in (2), multiply residues of  $\mathbf{c}$  by precomputed  $\lfloor \frac{\tilde{m}q_i}{q} \rfloor_{q_i}$  instead of  $\lfloor \frac{q_i}{q} \rfloor_{q_i}$ ). Since  $\{\tilde{m}\}$  is a single-modulus base, conversion of  $\mathbf{r}_{\tilde{m}}$  from  $\{\tilde{m}\}$  to  $\mathcal{B}_{\text{sk}}$  (l. 3 of Alg. 2) is a simple copy-paste when  $\tilde{m} < m_i$ . Finally, if  $\text{SmMR}_{\mathbf{q}_{\tilde{m}}}$  is performed right after a `FastBconv` from  $q$  (for converting  $\lfloor \tilde{m}\mathbf{c} \rfloor_q$ ),  $\tau$  is nothing but  $k$ .

### 4.3 Adapting the First Step

We recall that originally this step is the computation of  $[\text{DR}_2(\mathbf{ct}_*)]_q$ . Unlike the decryption, a  $\gamma$ -correction technique does not guarantee an exact rounding. Indeed, for the decryption we wanted to get  $\text{DR}_0([\mathbf{ct}(\mathbf{s})]_q)$ , and through  $\mathbf{s}$  we had access to the noise of  $\mathbf{ct}$ , on which we have some control. In the present context, we cannot ensure a condition like  $\| \lfloor t \cdot \mathbf{ct}_* \rfloor_q \|_\infty \leq q(\frac{1}{2} - \varepsilon)$ , for some  $\varepsilon^{-1} \sim 2^\nu$ , which would enable the use of an efficient  $\gamma$ -correction. Thus, we suggest to perform a simple uncorrected RNS flooring. For that purpose, we define:

$$\forall \mathbf{a} \in \mathcal{R}, \text{fastRNSFloor}_q(\mathbf{a}, m) := (\mathbf{a} - \text{FastBconv}(\lfloor \mathbf{a} \rfloor_q, q, m)) \times \lfloor q^{-1} \rfloor_m \pmod{m}.$$

Algorithm 2 should be executed first. Consequently, by Lem. 4, if  $\tilde{m}$  satisfies the bound in (8) for a given parameter  $\rho > 0$ , we assume having, in  $\mathcal{B}_{\text{sk}}$ , the residues of  $\mathbf{ct}'_i$  ( $\equiv \mathbf{ct}_i \pmod{q}$ ) such that:

$$\|\mathbf{ct}'_* := \mathbf{ct}'_1 \star \mathbf{ct}'_2\|_\infty \leq \delta \frac{q^2}{2} (1 + \rho)^2 . \quad (9)$$

The parameter  $\rho$  is determined in practice. Notice that, in base  $q$ ,  $ct'_i$  and  $ct_i$  are equal.

**Lemma 5.** *Let the residues of  $ct'_i \equiv ct_i \pmod q$  be given in base  $q \cup \mathcal{B}_{sk}$ , with  $\|ct'_i\|_\infty \leq \frac{q}{2}(1 + \rho)$  for  $i \in \{1, 2\}$ . Let  $ct'_* = ct'_1 \star ct'_2$ . Then, for  $j \in \{0, 1, 2\}$ ,*

$$fastRNSFloor_q(t \cdot ct'_*[j], \mathcal{B}_{sk}) = \left\lfloor \frac{t}{q} ct'_*[j] \right\rfloor + \mathbf{b}_j \text{ in } \mathcal{B}_{sk}, \text{ with } \|\mathbf{b}_j\|_\infty \leq k. \quad (10)$$

A first part of the noise growth is detailed in the following proposition.

**Proposition 1.** *Let  $\widetilde{ct}_{mult} = DR_2(ct'_*)$  with (9) satisfied, and  $r_\infty := \frac{1+\rho}{2}(1 + \delta B_{key}) + 1$ . Let  $\mathbf{v}_i$  be the inherent noise of  $ct'_i$ . Then  $\widetilde{ct}_{mult}(\mathbf{s}) = \Delta[\mathbf{m}_1 \mathbf{m}_2]_t + \widetilde{\mathbf{v}}_{mult} \pmod q$  with:*

$$\begin{aligned} \|\widetilde{\mathbf{v}}_{mult}\|_\infty &< \delta t(r_\infty + \frac{1}{2})(\|\mathbf{v}_1\|_\infty + \|\mathbf{v}_2\|_\infty) + \frac{\delta}{2} \min \|\mathbf{v}_i\|_\infty + \delta t|q|_t(r_\infty + 1) \\ &+ \frac{1}{2}(3 + |q|_t + \delta B_{key}(1 + \delta B_{key})). \end{aligned} \quad (11)$$

#### 4.4 Transitional Step

Lemma 5 states that we have got back  $DR_2(ct'_*) + \mathbf{b}$  in  $\mathcal{B}_{sk}$  so far, where we have denoted  $(\mathbf{b}_0, \mathbf{b}_1, \mathbf{b}_2)$  by  $\mathbf{b}$ . To perform the second step of multiplication, we need to convert it in base  $q$ . However, the conversion has to be exact because extra multiples of  $M = m_1 \dots m_\ell$  cannot be tolerated.  $m_{sk}$  allows us to perform a complete Shenoy and Kumaresan like conversion [24]. The next lemma describes such kind of conversion for a more general context where the input can be either positive or negative, and can be larger, in absolute value, than  $M$ .

**Lemma 6.** *Let  $\mathcal{B}$  be an RNS base and  $m_{sk}$  be a modulus coprime to  $M = \prod_{m \in \mathcal{B}} m$ . Let  $x$  be an integer such that  $|x| < \lambda M$  (for some real number  $\lambda \geq 1$ ) and whose residues are given in  $\mathcal{B}_{sk}$ . Let's assume that  $m_{sk}$  satisfies  $m_{sk} \geq 2(|\mathcal{B}| + \lceil \lambda \rceil)$ . Let  $\alpha_{sk,x}$  be the following integer:*

$$\alpha_{sk,x} := \left[ (FastBconv(x, \mathcal{B}, \{m_{sk}\}) - x_{sk})M^{-1} \right]_{m_{sk}}. \quad (12)$$

Then, for  $x$  being either positive or negative, the following equality holds:

$$FastBconvSK(x, \mathcal{B}_{sk}, q) := (FastBconv(x, \mathcal{B}, q) - \alpha_{sk,x}M) \pmod q = x \pmod q. \quad (13)$$

Consequently, since  $\|DR_2(ct'_*) + \mathbf{b}\|_\infty \leq \delta t \frac{q}{2}(1 + \rho)^2 + \frac{1}{2} + k$ , we can establish the following proposition.

**Proposition 2.** *Given a positive real number  $\lambda$ , let  $m_{sk}$  and  $\mathcal{B}$  be such that:*

$$\lambda M > \delta t \frac{q}{2}(1 + \rho)^2 + \frac{1}{2} + k, \quad m_{sk} \geq 2(|\mathcal{B}| + \lceil \lambda \rceil). \quad (14)$$

Let's assume that  $DR_2(ct'_*) + \mathbf{b}$  is given in  $\mathcal{B}_{sk}$ , with  $\|\mathbf{b}\|_\infty \leq k$ . Then,

$$FastBconvSK(DR_2(ct'_*) + \mathbf{b}, \mathcal{B}_{sk}, q) = (DR_2(ct'_*) + \mathbf{b}) \pmod q.$$

#### 4.5 Adapting the Second Step

At this point,  $\widetilde{\mathbf{ct}}_{mult} + \mathbf{b} = (\bar{\mathbf{c}}_0, \bar{\mathbf{c}}_1, \bar{\mathbf{c}}_2)$  is known in base  $q$  ( $\widetilde{\mathbf{ct}}_{mult} := \text{DR}_2(\mathbf{ct}'_*)$ ). We recall that the original second step of homomorphic multiplication would be done as follows:

$$\mathbf{ct}_{mult} = ([\bar{\mathbf{c}}_0 + \langle \mathcal{D}_{\omega,q}(\bar{\mathbf{c}}_2), \mathcal{P}_{\omega,q}(\mathbf{s}^2) - (\vec{\mathbf{e}} + \mathbf{s}\vec{\mathbf{a}}) \rangle]_q, [\bar{\mathbf{c}}_1 + \langle \mathcal{D}_{\omega,q}(\bar{\mathbf{c}}_2), \vec{\mathbf{a}} \rangle]_q) \quad (15)$$

where  $\vec{\mathbf{e}} \leftarrow \chi_{err}^{\ell_{\omega,q}}$ ,  $\vec{\mathbf{a}} \leftarrow \mathcal{U}(\mathcal{R}_q)^{\ell_{\omega,q}}$ . The decomposition of  $\bar{\mathbf{c}}_2$  in radix  $\omega$  enables a crucial reduction of the noise growth due to the multiplications by the terms  $\mathbf{e}_i + \mathbf{s}\mathbf{a}_i$ . It cannot be done directly in RNS as is. Indeed, it would require a costly switch between RNS and positional representation in radix  $\omega$ . However, we can do something very similar. We recall that we can write  $\bar{\mathbf{c}}_2 = \sum_{i=1}^k |\bar{\mathbf{c}}_2^{\frac{q_i}{q}}|_{q_i} \times \frac{q}{q_i} \pmod{q}$ . If  $\omega$  has the same order of magnitude than  $2^\nu$  (size of moduli in  $q$ ), we obtain a similar limitation of the noise growth by using the vectors  $\xi_q(\bar{\mathbf{c}}_2) = (|\bar{\mathbf{c}}_2^{\frac{q_1}{q}}|_{q_1}, \dots, |\bar{\mathbf{c}}_2^{\frac{q_k}{q}}|_{q_k})$  and  $\mathcal{P}_{\text{RNS},q}(\mathbf{s}^2) = (|\mathbf{s}^2_{\frac{q_1}{q}}|_q, \dots, |\mathbf{s}^2_{\frac{q_k}{q}}|_q)$ , both in  $\mathcal{R}^k$ . This is justified by the following lemma.

**Lemma 7.** *For any  $\mathbf{c} \in \mathcal{R}$ ,  $\langle \xi_q(\mathbf{c}), \mathcal{P}_{\text{RNS},q}(\mathbf{s}^2) \rangle \equiv \mathbf{c}\mathbf{s}^2 \pmod{q}$ .*

Thus, we replace the public relinearisation key  $\mathbf{rlk}_{\text{FV}}$  by the following one:  $\mathbf{rlk}_{\text{RNS}} = ([\mathcal{P}_{\text{RNS},q}(\mathbf{s}^2) - (\vec{\mathbf{e}} + \mathbf{s}\vec{\mathbf{a}})]_q, \vec{\mathbf{a}})$ . The next lemma helps for providing a bound on the extra noise introduced by this step.

**Lemma 8.** *Let  $\vec{\mathbf{e}} \leftarrow \chi_{err}^k$ ,  $\vec{\mathbf{a}} \leftarrow \mathcal{U}(\mathcal{R}_q)^k$ , and  $\mathbf{c} \in R$ . Then,*

$$\|(\langle \xi_q(\mathbf{c}), -(\vec{\mathbf{e}} + \vec{\mathbf{a}}\mathbf{s}) \rangle + \langle \xi_q(\mathbf{c}), \vec{\mathbf{a}} \rangle \mathbf{s}) \pmod{q}\|_\infty < \delta B_{err} k 2^\nu . \quad (16)$$

*Remark 2.* It is still possible to add a second level of decomposition (like in original approach, but applied on the residues) to limit a bit more the noise growth. Furthermore, Sect. 4.6 details how the size of  $\mathbf{rlk}_{\text{RNS}}$  can be reduced in a similar way that  $\mathbf{rlk}_{\text{FV}}$  could be through the method described in ([4], 5.4).

Finally, the output of the new variant of multiplication is the following one:

$$\mathbf{ct}_{mult} = \left( [\bar{\mathbf{c}}_0 + \langle \xi_q(\bar{\mathbf{c}}_2), \mathcal{P}_{\text{RNS},q}(\mathbf{s}^2) - (\vec{\mathbf{e}} + \vec{\mathbf{a}}\mathbf{s}) \rangle]_q, [\bar{\mathbf{c}}_1 + \langle \xi_q(\bar{\mathbf{c}}_2), \vec{\mathbf{a}} \rangle]_q \right) . \quad (17)$$

**Proposition 3.** *Let  $\mathbf{ct}_{mult}$  be as in (17), and  $\mathbf{v}_{mult}$  (resp.  $\widetilde{\mathbf{v}}_{mult}$ ) the inherent noise of  $\mathbf{ct}_{mult}$  (resp.  $\widetilde{\mathbf{ct}}_{mult}$ ). Then  $\mathbf{ct}_{mult}(\mathbf{s}) = \Delta[\mathbf{m}_1\mathbf{m}_2]_t + \mathbf{v}_{mult} \pmod{q}$  with:*

$$\|\mathbf{v}_{mult}\|_\infty < \|\widetilde{\mathbf{v}}_{mult}\|_\infty + k(1 + \delta B_{key}(1 + \delta B_{key})) + \delta k B_{err} 2^{\nu+1} . \quad (18)$$

Algorithm 3 depicts the scheme of the new full RNS variant  $\mathbf{Mult}_{\text{RNS}}$ .

---

**Algorithm 3** RNS homomorphic multiplication  $\text{Mult}_{\text{RNS}}$ 


---

**Require:**  $\text{ct}_1, \text{ct}_2$  in  $q$ 
**Ensure:**  $\text{ct}_{\text{mult}}$  in  $q$ 

S0: Convert fast  $\text{ct}_1$  and  $\text{ct}_2$  from  $q$  to  $\mathcal{B}_{\text{sk}} \cup \{\tilde{m}\}$ :  $\rightsquigarrow \text{ct}_i'' = \text{ct}_i + q\text{-overflows}$ 

S1: Reduce  $q$ -overflows in  $\mathcal{B}_{\text{sk}}$ :  $(\text{ct}'_i \text{ in } \mathcal{B}_{\text{sk}}) \leftarrow \text{SmMRq}_{\tilde{m}}(((\text{ct}''_i)_m)_{m \in \mathcal{B}_{\text{sk}} \cup \{\tilde{m}\}})$ 

S2: Compute the product  $\text{ct}'_{\star} = \text{ct}'_1 \star \text{ct}'_2$  in  $q \cup \mathcal{B}_{\text{sk}}$ 

S3: Convert fast from  $q$  to  $\mathcal{B}_{\text{sk}}$  to achieve first step (approximate rounding) in  $\mathcal{B}_{\text{sk}}$ :

$$(\widetilde{\text{ct}}_{\text{mult}} + \mathbf{b} = \text{DR}_2(\text{ct}'_{\star}) + \mathbf{b} \text{ in } \mathcal{B}_{\text{sk}}) \leftarrow \dots \leftarrow \text{FastBconv}(t \cdot \text{ct}'_{\star}, q, \mathcal{B}_{\text{sk}})$$

S4: Convert exactly from  $\mathcal{B}_{\text{sk}}$  to  $q$  to achieve transitional step:

$$(\widetilde{\text{ct}}_{\text{mult}} + \mathbf{b} \text{ in } q) \leftarrow \text{FastBconvSK}(\widetilde{\text{ct}}_{\text{mult}} + \mathbf{b}, \mathcal{B}_{\text{sk}}, q)$$

S5: Perform second step (relinearisation) in  $q$ :

$$\text{ct}_{\text{mult}} \leftarrow \text{Relin}_{\text{RNS}}(\widetilde{\text{ct}}_{\text{mult}} + \mathbf{b}) \bmod (q_1, \dots, q_k)$$


---

#### 4.6 Reducing the Size of the Relinearization Key $\text{rlk}_{\text{RNS}}$

In [4], Sect. 5.4, a method to significantly reduce the size of the public evaluation key  $\text{evk}$  is described (by truncating the ciphertext) and it is applicable to the original FV scheme. We provide an efficient adaptation of such kind of optimization to the RNS variant of the relinearisation step.

We recall that the relinearisation is applied to a degree-2 ciphertext denoted here by  $(\mathbf{c}_0, \mathbf{c}_1, \mathbf{c}_2)$ . The initial suggestion was to set to zero, say, the  $i$  lowest significant components of the vector  $\mathcal{D}_{\omega, q}(\mathbf{c}_2)$ . Doing so is equivalent to replacing  $\mathbf{c}_2$  by  $\mathbf{c}'_2 = \omega^i \lfloor \mathbf{c}_2 \omega^{-i} \rfloor = \mathbf{c}_2 - |\mathbf{c}_2|_{\omega^i}$ . Thus, only the  $\ell_{\omega, q} - i$  most significant components of  $\text{rlk}_{\text{FV}}[0]$  (and then of  $\text{rlk}_{\text{FV}}[1]$ ) are required (in other words, when  $\text{rlk}_{\text{FV}}[0]$  is viewed as an  $(\ell_{q, \omega}, k)$  RNS matrix by decomposing each component in base  $q$ ,  $ik$  entries are set to zero like this). This optimization causes a greater noise than the one in Lem. 4 of [4]. Given  $(\mathbf{c}_0, \mathbf{c}_1, \mathbf{c}_2)$  decryptable under  $\mathbf{s}$ , the relinearisation step provides the following ciphertext:

$$(\tilde{\mathbf{c}}_0, \tilde{\mathbf{c}}_1) := (\mathbf{c}_0 + \langle \mathcal{D}_{\omega, q}(\mathbf{c}'_2), \mathcal{P}_{\omega, q}(\mathbf{s}^2) - (\vec{\mathbf{e}} + \vec{\mathbf{a}}\mathbf{s}) \rangle, \mathbf{c}_1 + \langle \mathcal{D}_{\omega, q}(\mathbf{c}'_2), \vec{\mathbf{a}} \rangle) .$$

Thus,  $(\tilde{\mathbf{c}}_0, \tilde{\mathbf{c}}_1)(\mathbf{s}) = \mathbf{c}_0 + \mathbf{c}_1\mathbf{s} + \mathbf{c}'_2\mathbf{s}^2 - \langle \mathcal{D}_{\omega, q}(\mathbf{c}'_2), \vec{\mathbf{e}} \rangle \bmod q$ . Consequently, the extra noise comes from the following term:

$$\| -|\mathbf{c}_2|_{\omega^i} \mathbf{s}^2 - \langle \mathcal{D}_{\omega, q}(\mathbf{c}'_2), \vec{\mathbf{e}} \rangle \|_{\infty} = \| -|\mathbf{c}_2|_{\omega^i} \mathbf{s}^2 - \sum_{j=i}^{\ell_{\omega, q}-1} \mathcal{D}_{\omega, q}(\mathbf{c}_2)_j \mathbf{e}_j \|_{\infty} \quad (19)$$

$$< \delta^2 \omega^i B_{\text{key}}^2 + (\ell_{\omega, q} - i) \delta \omega B_{\text{err}} .$$

In the present RNS variant, the computation of  $\lfloor \mathbf{c}_2 \omega^{-i} \rfloor$  is not straightforward. This could be replaced by  $\lfloor \mathbf{c}_2 (q_1 \dots q_i)^{-1} \rfloor$  through a Newton like interpolation (also known as mixed-radix conversion [11]). Though the result would be quite similar to the original optimization in terms of noise growth, its efficiency is not satisfying. Indeed, despite  $ik$  entries of the RNS matrix  $\text{rlk}_{\text{RNS}}[0]$  can be set to zero like this, such a Newton interpolation is intrinsically sequential, while the division by  $\omega^i$  followed by a flooring is simply achieved by an immediate zeroing of the lowest significant coefficients in radix  $\omega$  representation.

For our approach, we rely on the fact that  $\text{rlk}_{\text{RNS}}$  contains the RLWE-encryptions of the polynomials  $\lfloor \mathbf{s}^2 \frac{q}{q_j} \rfloor_q$ . Then, we notice that only the  $j^{\text{th}}$ -residue

of  $|\mathbf{s}^2 \frac{q}{q_j}|_q$  can be non zero. So, let's assume that we want to cancel  $ik$  entries in  $\mathbf{rlk}_{\text{RNS}}[0]$  (as it has been done in  $\mathbf{rlk}_{\text{FV}}$  with the previous optimization). Then we choose, for each index  $j$ , a subset of index-numbers  $\mathcal{I}_j \subseteq [1, k] \setminus \{j\}$  with cardinality  $i$  (i.e. at line  $j$  of  $\mathbf{rlk}_{\text{RNS}}$ , choose  $i$  columns, except the diagonal one; these terms will be set to zero). Next, for each  $j$ , we introduce an RLWE-encryption of  $|\mathbf{s}^2 \frac{q}{q_j q_{\mathcal{I}_j}}|_q$ , where  $q_{\mathcal{I}_j} = \prod_{s \in \mathcal{I}_j} q_s$ , which is  $(|\mathbf{s}^2 \frac{q}{q_j q_{\mathcal{I}_j}} - (e_j + \mathbf{s}\mathbf{a}_j)|_q, \mathbf{a}_j)$ . So far, the underlying security features are still relevant. Now, it remains to multiply this encryption by  $q_{\mathcal{I}_j}$ , which gives in particular  $|\mathbf{s}^2 \frac{q}{q_j} - q_{\mathcal{I}_j}(e_j + \mathbf{s}\mathbf{a}_j)|_q$ . This is the  $j^{\text{th}}$ -line of the new matrix  $\mathbf{rlk}'_{\text{RNS}}[0]$ . It is clear that this line contains zeros at columns index-numbered by  $\mathcal{I}_j$ .  $\mathbf{rlk}_{\text{RNS}}[1] = (\mathbf{a}_1, \dots, \mathbf{a}_k)$  is replaced by  $\mathbf{rlk}'_{\text{RNS}}[1] = (|q_{\mathcal{I}_1} \mathbf{a}_1|_q, \dots, |q_{\mathcal{I}_k} \mathbf{a}_k|_q)$ .

Let's analyse the new noise growth. By evaluating in  $\mathbf{s}$  the output of relinearisation with this new  $\mathbf{rlk}'_{\text{RNS}}$ , we obtain:

$$\begin{aligned} & \mathbf{c}_0 + \langle \xi_q(\mathbf{c}_2), \mathbf{rlk}'_{\text{RNS}}[0] \rangle + (\mathbf{c}_1 + \langle \xi_q(\mathbf{c}_2), \mathbf{rlk}'_{\text{RNS}}[1] \rangle) \mathbf{s} \\ &= \mathbf{c}_0 + \sum_{j=1}^k |\mathbf{c}_2 \frac{q_j}{q}|_{q_j} \left( \mathbf{s}^2 \frac{q}{q_j} - q_{\mathcal{I}_j}(e_j + \mathbf{s}\mathbf{a}_j) \right) + \left( \mathbf{c}_1 + \sum_{j=1}^k |\mathbf{c}_2 \frac{q_j}{q}|_{q_j} q_{\mathcal{I}_j} \mathbf{a}_j \right) \mathbf{s} \\ &= \mathbf{c}_0 + \mathbf{c}_1 \mathbf{s} + \mathbf{c}_2 \mathbf{s}^2 - \sum_{j=1}^k |\mathbf{c}_2 \frac{q_j}{q}|_{q_j} q_{\mathcal{I}_j} e_j . \end{aligned}$$

Consequently, the cancellation of  $ik$  terms in the public matrix  $\mathbf{rlk}_{\text{RNS}}[0]$  by using this method causes an extra noise growth bounded by (this can be fairly compared to (19) in the case where  $\omega = 2^\nu$ , i.e.  $k = \ell_{\omega, q}$ ):

$$\| \sum_{j=1}^k |\mathbf{c}_2 \frac{q_j}{q}|_{q_j} q_{\mathcal{I}_j} e_j \|_\infty < \sum_{j=1}^k \delta q_j q_{\mathcal{I}_j} B_{err} < \delta k 2^{\nu(i+1)} B_{err} .$$

Therefore, the truncation of ciphertexts can be efficiently adapted to RNS representation without causing more significant noise growth.

#### 4.7 About Computational Complexity

In a classical multi-precision (MP) variant, for the purpose of efficiency the multiplication can involve NTT-based polynomial multiplication to implement the ciphertext product (e.g. [23] for such kind of implementation on FPGA). This approach requires the use of a base  $\mathcal{B}'$  (besides  $q$ ) with  $|\mathcal{B}'| = k + 1$  for storing the product. Notice that, in RNS variant, we also have  $|\mathcal{B}_{sk}| = k + 1$  (the same amount of information is kept, but used differently). Thus, it is easy to show that RNS and MP variants involve the same number of NTT and invNTT operations in the case  $\ell_{\omega, q} = k$  (e.g. the keys  $\mathbf{rlk}_{\text{FV}}$  and  $\mathbf{rlk}_{\text{RNS}}$  have the same size in this situation). In other words, the same number of polynomial products is performed in both cases when  $\omega = 2^\nu$ .

Above all, the RNS variant enables us to decrease the cost of other parts of the computation. Despite the fact that the asymptotic computational complexity of these parts remains identical for both variants, i.e.  $\mathcal{O}(k^2 n)$  elementary multiplications, the RNS variant only involves single-precision integer arithmetic.

To sum up, because of a complexity of  $\mathcal{O}(k^2 n \log_2(n))$  due to the NTT's, we keep the same asymptotic computational complexity  $\mathcal{C}(\text{Mult}_{\text{FV}}) \underset{n \rightarrow +\infty}{\sim} \mathcal{C}(\text{Mult}_{\text{RNS}})$ .

However, the most important fact is that multi-precision multiplications within MP variant are replaced in RNS by fast base conversions, which are simple matrix-vector products. Thus,  $\text{Mult}_{\text{RNS}}$  retains all the benefits of RNS properties and is highly parallelizable.

## 5 Software Implementation

The C++ `NFLlib` library [19] was used for efficiently implementing the arithmetic in  $\mathcal{R}_q$ . It provides an efficient NTT-based product in  $\mathcal{R}_q$  for  $q$  a product of 30 or 62-bit prime integers, and with degree  $n$  a power of 2, up to  $2^{15}$ .

### 5.1 Concrete Examples of Parameter Settings

In this part, we analyse which depth can be reached in a multiplicative tree, and for which parameters.

The initial noise is at most  $V = B_{err}(1 + 2\delta B_{key})$  [17]. The output of a tree of depth  $L$  has a noise bounded by  $C_{\text{RNS},1}^L V + LC_{\text{RNS},1}^{L-1} C_{\text{RNS},2}$  (cf. [4], Lem. 9) with, for the present RNS variant:

$$\begin{cases} C_{\text{RNS},1} = 2\delta^2 t^{\frac{(1+\rho)}{2}} B_{key} + \delta t(4 + \rho) + \frac{\delta}{2}, \\ C_{\text{RNS},2} = (1 + \delta B_{key})(\delta t|q|_t^{\frac{1+\rho}{2}} + \delta B_{key}(k + \frac{1}{2})) + 2\delta t|q|_t + k(\delta B_{err}2^{\nu+1} + 1) \\ \quad + \frac{1}{2}(3 + |q|_t). \end{cases} \tag{20}$$

We denote by  $L_{\text{RNS}} = \max\{L \in \mathbb{N} \mid C_{\text{RNS},1}^L V + LC_{\text{RNS},1}^{L-1} C_{\text{RNS},2} \leq \frac{q}{t}(\frac{1}{2} - \frac{k}{\gamma}) - \frac{|q|_t}{2}\}$  the depth allowed by  $\text{Mult}_{\text{RNS}}$ , when  $\text{Dec}_{\text{RNS}}$  is used for decryption.

For an 80-bit security level and parameters  $B_{key} = 1$ ,  $\sigma_{err} = 8$ ,  $B_{err} = 6\sigma_{err}$ , we consider the security analysis in [17], which provides ranges for  $(\log_2(q), n)$  (cf. [17], Tab. 2). We analyse our parameters by using the moduli given in `NFLlib`, because those were used for concrete testing. For a 32-bit (resp. 64) implementation, a set of 291 30-bit (resp. 1000 62-bit) moduli is available. These moduli are chosen to enable an efficient modular reduction (cf. [19], Alg. 2).

$n$	$k$	$t$	$L_{\text{RNS}} (L_{\text{std}})$	$\rho$	$\tilde{m}$	$\lceil \log_2(m_{\text{sk}}) \rceil$	$\gamma$
$2^{11}$	3	2	2 (2)	5	(no need)	18	7
		$2^{10}$	1 (1)	5	(no need)	27	7
$2^{12}$	6	2	5 (6)	11	(no need)	21	13
		$2^{10}$	4 (4)	10	2	29	54
$2^{13}$	13	2	13 (13)	$\frac{1}{3}$	81	15	36
		$2^{10}$	9 (9)	13	3	31	58
$2^{14}$	26	2	25 (25)	$\frac{1}{2}$	106	17	53
		$2^{10}$	19 (19)	1	53	27	53
$2^{15}$	53	2	50 (50)	$\frac{1}{16}$	1712	20	753
		$2^{10}$	38 (38)	$\frac{1}{2}$	214	30	107

Table 1: Examples of parameter settings, by using 30-bit moduli of `NFLlib`.

Table 1 lists parameters when  $q$  and  $\mathcal{B}$  are built with the 30-bit moduli of `NFLlib`. These parameters were determined by choosing the largest  $\rho$  (up to  $2k - 1$ ) allowing to reach the depth  $L_{\text{RNS}}$ .  $L_{\text{std}}$  corresponds to the bounds given in [17] for a classical approach. Sufficient sizes for  $\gamma$  and  $m_{\text{sk}}$  (allowing, for  $m_{\text{sk}}$ , to have  $|\mathcal{B}| = k$  through (14) after having chosen for  $q$  the  $k$  greatest available moduli) are provided. For these specific parameters, the new bounds on the noise for RNS variant cause a smaller depth in only one case.

## 5.2 Influence of $\tilde{m}$ Over Noise Growth

After a fast conversion from  $q$ , ciphertexts in  $\mathcal{B}_{\text{sk}}$  can contain  $q$ -overflows and verify  $\|\text{ct}'_i\|_\infty < \frac{q}{2}(1+\tau)$ . In a multiplicative tree without any addition,  $\tau \leq 2k-1$ . By applying Alg. 2, this bound decreases to  $\frac{q}{2}(1+\rho)$ , for some  $0 < \rho \leq 2k-1$ . Having  $\rho = 2k-1$  in Tab. 1 means it is unnecessary to use  $\text{SmMRq}_{\tilde{m}}$  to reach the best depth. It happens only three times. Most of the time, doing such reduction is necessary before a multiplication so as to reach the highest depth. Moreover, choosing a lower  $\rho$  (i.e. higher  $\tilde{m}$ ) than necessary can decrease the size of  $\gamma$  and  $m_{\text{sk}}$  (as shown in Tab. 1, one set of parameters leads to  $\lceil \log_2(m_{\text{sk}}) \rceil = 31$ , avoiding the use of a 30-bit modulus; this can be solved by taking a larger  $\tilde{m}$ ).

To illustrate the impact of  $\text{SmMRq}_{\tilde{m}}$ , Fig. 1 depicts the noise growth for  $\tilde{m} \in \{0, 2^8, 2^{16}\}$ . Given Tab. 1,  $\tilde{m} = 2^8$  is sufficient in such scenario to reach  $L_{\text{RNS}} = 13$ . Against a computation with no reduction at all ( $\tilde{m} = 0$ , implying  $L_{\text{RNS}} = 11$  in this case), choosing  $\tilde{m} = 2^8$  implies an average reduction of 25%. By using  $\tilde{m} = 2^{16}$ , we gain around 32%. Therefore,  $\text{SmMRq}_{\tilde{m}}$  has been systematically integrated within the implementation of  $\text{Mult}_{\text{RNS}}$  for timing measurements in next part.

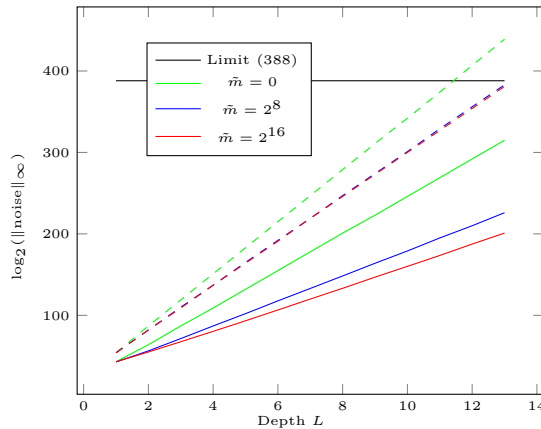


Fig. 1: Example of noise growth;  $n = 2^{13}$ ,  $\log_2(q) = 390$  ( $\nu = 30$ ,  $k = 13$ ),  $t = 2$ ,  $\sigma_{err} = 8$ ,  $B_{key} = 1$  (dashed line: bound using (20); plain line: measurements).

## 5.3 Some Remarks

*Convenient  $\tilde{m}$  and  $\gamma$*  Given values of  $\rho$  in Tab. 1,  $\tilde{m} = 2^8$  (resp.  $\tilde{m} = 2^{16}$ ) satisfies, by far, any set of analysed parameters. This enables an efficient and straightforward modular arithmetic through standard types like `uint8_t` (resp. `uint16_t`) and a type conversion towards the signed `int8_t` (resp. `int16_t`) immediately gives the centered remainder. Parameter analysis with such  $\tilde{m}$  shows that  $\gamma = 2^8$  is sufficient to ensure a correct decryption for all configurations. A reduction modulo  $\gamma$  can then be achieved by a simple type conversion to `uint8_t`.

*Tested Algorithms* The code<sup>1</sup> we compared with was implemented in the context of HEAT [26]. It is based on `NFLlib` too. Multi-precision arithmetic is handled with `GMP 6.1.0` [14], and multiplications by  $\frac{t}{q}$  are performed with integer divisions. `MultMP` and `DecMP` denote functions from this code.

<sup>1</sup> <https://github.com/CryptoExperts/FV-NFLlib>



$\text{Mult}_{\text{RNS}}$  was implemented as described by Alg. 3. Could the use of  $\text{SmMRq}_{\tilde{m}}$  be avoided to reach the maximal theoretical depth, it was however systematically used. Its cost is negligible and it enables a noticeable decrease of noise growth.

Two variants of  $\text{Dec}_{\text{RNS}}$  (cf. Sect. 3.5) have been implemented. Depending on  $\nu$ , the one with floating point arithmetic (named  $\text{Dec}_{\text{RNS-f1p}}$  thereafter) uses `double` (resp. `long double`) for double (resp. quadruple) precision, and then does not rely on any other external library at all.

### 5.4 Results

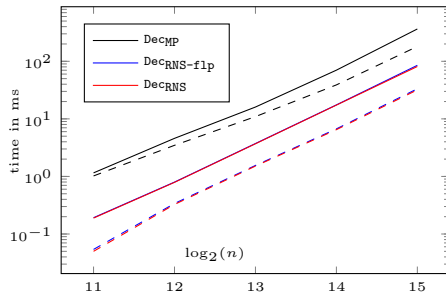


Fig. 2: Decryption time ( $t = 2^{10}$ ), with  $\nu = 30$  (plain lines) and  $\nu = 62$  (dashed lines).

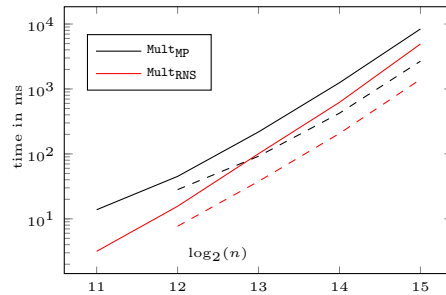


Fig. 3: Multiplication time ( $t = 2^{10}$ ), with  $\nu = 30$  (plain lines) and  $\nu = 62$  (dashed lines).

The tests have been run on a laptop, running under Fedora 22 with Intel<sup>®</sup> Core<sup>™</sup> i7-4810MQ CPU @ 2.80GHz and using GNU compiler g++ version 5.3.1. Hyper-Threading and Turbo Boost were deactivated. The timings have been measured on  $2^{12}$  decryptions/multiplications for each configuration.

Figure 2 presents timings for  $\text{Dec}_{\text{MP}}$ ,  $\text{Dec}_{\text{RNS}}$  and  $\text{Dec}_{\text{RNS-f1p}}$ , and Fig. 3 depicts timings for  $\text{Mult}_{\text{MP}}$  and  $\text{Mult}_{\text{RNS}}$ . Both figures gather data for two modulus sizes:  $\nu = 30$  and  $\nu = 62$ . Step 2 of  $\text{Mult}_{\text{MP}}$  uses a decomposition in radix-base  $\omega = 2^{32}$  when  $\nu = 30$ , and  $\omega = 2^{62}$  when  $\nu = 62$ . The auxiliary bases  $\mathcal{B}_{sk}$  and  $\mathcal{B}'$  involved in  $\text{Mult}_{\text{RNS}}$  and  $\text{Mult}_{\text{MP}}$  contain  $k + 1$  moduli each. Table 2 shows which values of  $k$  have been tested (depending on  $n$ ). Multiplication timing for  $(n, \nu, k) = (2^{11}, 62, 1)$  is not given since  $L = 1$  already causes decryption failures.

$\log_2(n)$	11	12	13	14	15
$k (\nu = 30)$	3	6	13	26	53
$k (\nu = 62)$	1	3	6	12	25

Table 2: Parameter  $k$  used in the tests (i.e.  $\lceil \log_2(q) \rceil = k\nu$ ).

It has to be noticed that the performance are mainly due to `NFLlib`. The contributions appear in the speed-ups between RNS and MP variants.

In Fig. 3, the convergence of complexities of  $\text{Mult}_{\text{RNS}}$  and  $\text{Mult}_{\text{MP}}$  (as explained in Sect. 4.7) is noticeable. The new algorithm presented in this paper allows speed-ups from 4.3 to 1.7 for degree  $n$  from  $2^{11}$  to  $2^{15}$  when  $\nu = 30$ , and from 3.6 to 1.9 for  $n$  from  $2^{12}$  to  $2^{15}$  when  $\nu = 62$ .

In Fig. 2, the two variants of decryption described in 3.5 are almost equally fast. Indeed, they perform the same number of elementary (floating point or integer) operations. Between degree  $2^{11}$  and  $2^{15}$ , the RNS variants allow speed-ups varying from 6.1 to 4.4 when  $\nu = 30$ , and from 20.4 to 5.6 when  $\nu = 62$ . All the implemented decryption functions take as input a ciphertext in NTT representation. Thus, only one `invNTT` is performed (after the product of residues) within each decryption. As explained (cf. 3.5), despite a better asymptotic computational complexity for RNS decryption, the efficiency remains in practice highly related to this `invNTT` procedure, even maybe justifying the slight convergence between MP and RNS decryption times observed in Fig. 2.

## 6 Conclusion

In this paper, the somewhat homomorphic encryption scheme `FV` has been fully adapted to Residue Number Systems. Prior to this work, RNS was used to accelerate polynomial additions and multiplications. However, the decryption and the homomorphic multiplication involve operations at the coefficient level which are hardly compatible with RNS, such as division and rounding.

Our proposed solutions overcome these incompatibilities, without modifying the security features of the original scheme. As a consequence, we have provided a SHE scheme which only involves RNS arithmetic. It means that only single-precision integer arithmetic is required, and the new variant fully benefits from the properties of RNS, such as parallelization.

The proposed scheme has been implemented in software using C++. Because arithmetic on polynomials (in particular polynomial product) is not concerned by the new optimizations provided here, the implementation has been based on the `NFLlib` library, which embeds a very efficient NTT-based polynomial product. Our implementation has been compared to a classical version of `FV` (based on `NFLlib`, and `GMP`). For degrees from  $2^{11}$  to  $2^{15}$ , the new decryption (resp. homomorphic multiplication) offers speed-ups from 20 to 5 (resp. 4 to 2) folds for cryptographic parameters.

Further work should demonstrate the high potential of the new variant by exploiting all the concurrency properties of RNS, in particular through dedicated hardware implementations.

## Acknowledgement

We thank anonymous reviewers for their helpful comments and remarks.

This work was partially supported by the European Union’s H2020 Programme under grant agreement # ICT-644209, and by the Natural Sciences and Engineering Research Council (NSERC).

## References

1. M. Albrecht, S. Bai, and L. Ducas. A Subfield Lattice Attack on Overstretched NTRU Assumptions: Cryptanalysis of some FHE and Graded Encoding Schemes.

- Cryptology ePrint Archive, Report 2016/127, 2016. <http://eprint.iacr.org/2016/127>.
2. J.-C. Bajard, J. Eynard, M. A. Hasan, and V. Zucca. A Full RNS Variant of FV like Somewhat Homomorphic Encryption Schemes. Cryptology ePrint Archive, Report 2016/510, 2016. <http://eprint.iacr.org/2016/510>.
  3. J.-C. Bajard, J. Eynard, N. Merkiche, and T. Plantard. RNS Arithmetic Approach in Lattice-Based Cryptography: Accelerating the "Rounding-off" Core Procedure. In *22nd IEEE Symposium on Computer Arithmetic, ARITH 2015, Lyon, France, June 22-24, 2015*, pages 113–120. IEEE, 2015.
  4. J. W. Bos, K. E. Lauter, J. Loftus, and M. Naehrig. Improved Security for a Ring-Based Fully Homomorphic Encryption Scheme. In Martijn Stam, editor, *Cryptography and Coding - 14th IMA International Conference, IMACC 2013, Oxford, UK, December 17-19, 2013. Proceedings*, volume 8308 of *Lecture Notes in Computer Science*, pages 45–64. Springer, 2013.
  5. Z. Brakerski. Fully Homomorphic Encryption without Modulus Switching from Classical GapSVP. Cryptology ePrint Archive, Report 2012/078, 2012. <http://eprint.iacr.org/2012/078>.
  6. Z. Brakerski, C. Gentry, and V. Vaikuntanathan. (Leveled) Fully Homomorphic Encryption without Bootstrapping. In Shafi Goldwasser, editor, *Innovations in Theoretical Computer Science 2012, Cambridge, MA, USA, January 8-10, 2012*, pages 309–325. ACM, 2012.
  7. L. Ducas, A. Durmus, T. Lepoint, and V. Lyubashevsky. Lattice Signatures and Bimodal Gaussians. In Ran Canetti and Juan 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 40–56. Springer, 2013.
  8. L. Ducas, V. Lyubashevsky, and T. Prest. Efficient Identity-Based Encryption over NTRU Lattices. In Palash Sarkar and Tetsu Iwata, editors, *Advances in Cryptology - ASIACRYPT 2014 - 20th International Conference on the Theory and Application of Cryptology and Information Security, Kaoshiung, Taiwan, R.O.C., December 7-11, 2014, Proceedings, Part II*, volume 8874 of *Lecture Notes in Computer Science*, pages 22–41. Springer, 2014.
  9. J. Fan and F. Vercauteren. Somewhat Practical Fully Homomorphic Encryption. Cryptology ePrint Archive, Report 2012/144, 2012. <http://eprint.iacr.org/2012/144>.
  10. S. Garg, C. Gentry, and S. Halevi. Candidate Multilinear Maps from Ideal Lattices. In Thomas Johansson and Phong Q. Nguyen, editors, *Advances in Cryptology - EUROCRYPT 2013, 32nd Annual Intern. Conf. on the Theory and Applications of Cryptographic Techniques, Athens, Greece, May 26-30, 2013. Proceedings*, volume 7881 of *Lecture Notes in Computer Science*, pages 1–17. Springer, 2013.
  11. H. L. Garner. The Residue Number System. In *Papers Presented at the March 3-5, 1959, Western Joint Computer Conference, IRE-AIEE-ACM '59 (Western)*, pages 146–153, New York, NY, USA, 1959. ACM.
  12. C. Gentry. Fully Homomorphic Encryption Using Ideal Lattices. In *Proceedings of the Forty-first Annual ACM Symposium on Theory of Computing, STOC '09*, pages 169–178, New York, NY, USA, 2009. ACM.
  13. R. Gilad-Bachrach, N. Dowlin, K. Laine, K. E. Lauter, M. Naehrig, and J. Wernsing. CryptoNets: Applying Neural Networks to Encrypted Data with High Throughput and Accuracy. In Maria-Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of the 33rd International Conference on Machine Learning*,

- ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 201–210. JMLR.org, 2016.
14. T. Granlund and the GMP development team. *GNU MP: The GNU Multiple Precision Arithmetic Library*, 6.1.0 edition, 2015. <http://gmplib.org/>.
  15. J. Hoffstein, J. Pipher, and J. H. Silverman. NTRU: A Ring-Based Public Key Cryptosystem. In Joe Buhler, editor, *Algorithmic Number Theory, Third Intern. Symp., ANTS-III, Portland, Oregon, USA, June 21-25, 1998, Proceedings*, volume 1423 of *Lecture Notes in Computer Science*, pages 267–288. Springer, 1998.
  16. A. Langlois, D. Stehlé, and R. Steinfeld. GGHLite: More Efficient Multilinear Maps from Ideal Lattices. In Phong Q. Nguyen and Elisabeth Oswald, editors, *Advances in Cryptology - EUROCRYPT 2014 - 33rd Annual International Conference on the Theory and Applications of Cryptographic Techniques, Copenhagen, Denmark, May 11-15, 2014. Proceedings*, volume 8441 of *Lecture Notes in Computer Science*, pages 239–256. Springer, 2014.
  17. T. Lepoint and M. Naehrig. A Comparison of the Homomorphic Encryption Schemes FV and YASHE. In David Pointcheval and Damien Vergnaud, editors, *Progress in Cryptology - AFRICACRYPT 2014 - 7th International Conference on Cryptology in Africa, Marrakesh, Morocco, May 28-30, 2014. Proceedings*, volume 8469 of *Lecture Notes in Computer Science*, pages 318–335. Springer, 2014.
  18. V. Lyubashevsky. Lattice Signatures without Trapdoors. In *EUROCRYPT*, volume 7237 of *Lecture Notes in Computer Science*, pages 738–755. Springer, 2012.
  19. C. Aguilar Melchor, J. Barrier, S. Guelton, A. Guinet, M.-O. Killijian, and T. Lepoint. NTLlib: NTT-Based Fast Lattice Library. In Kazue Sako, editor, *Topics in Cryptology - CT-RSA 2016 - The Cryptographers' Track at the RSA Conference 2016, San Francisco, CA, USA, February 29 - March 4, 2016, Proceedings*, volume 9610 of *Lecture Notes in Computer Science*, pages 341–356. Springer, 2016.
  20. P. L. Montgomery. Modular Multiplication without Trial Division. *Mathematics of Computation*, 44(170):519–521, 1985.
  21. T. Oder, T. Pöppelmann, and T. Güneysu. Beyond ECDSA and RSA: Lattice-based Digital Signatures on Constrained Devices. In *DAC*, pages 110:1–110:6. ACM, 2014.
  22. C. Peikert. Lattice Cryptography for the Internet. In Michele Mosca, editor, *Post-Quantum Cryptography - 6th International Workshop, PQCrypto 2014, Waterloo, ON, Canada, October 1-3, 2014. Proceedings*, volume 8772 of *Lecture Notes in Computer Science*, pages 197–219. Springer, 2014.
  23. S. S. Roy, K. Järvinen, F. Vercauteren, V. S. Dimitrov, and I. Verbauwhede. Modular Hardware Architecture for Somewhat Homomorphic Function Evaluation. In Tim Güneysu and Helena Handschuh, editors, *Cryptographic Hardware and Embedded Systems - CHES 2015 - 17th International Workshop, Saint-Malo, France, September 13-16, 2015, Proceedings*, volume 9293 of *Lecture Notes in Computer Science*, pages 164–184. Springer, 2015.
  24. A. P. Shenoy and R. Kumaresan. Fast Base Extension Using a Redundant Modulus in RNS. *IEEE Trans. Computers*, 38(2):292–297, 1989.
  25. D. Stehlé, R. Steinfeld, K. Tanaka, and K. Xagawa. Efficient Public Key Encryption Based on Ideal Lattices. In Mitsuru Matsui, editor, *Advances in Cryptology - ASIACRYPT 2009, 15th Intern. Conf. on the Theory and Application of Cryptology and Information Security, Tokyo, Japan, December 6-10, 2009. Proceedings*, volume 5912 of *Lecture Notes in Computer Science*, pages 617–635. Springer, 2009.
  26. European Union. Homomorphic Encryption, Applications and Technology (HEAT). <https://heat-project.eu>. H2020-ICT-2014-1, Project reference: 644209.

## A Proofs

### A.1 Lemma 1

According to (2),  $\text{FastBconv}(|t\gamma \cdot \text{ct}(\mathbf{s})|_q, q, \{t, \gamma\})$  provides  $|t\gamma \cdot \text{ct}(\mathbf{s})|_q + q\mathbf{a}$ , where each coefficient  $\mathbf{a}_i$  is an integer lying in  $[0, k-1]$ . Let  $m$  be  $t$  or  $\gamma$ . Then,

$$\begin{aligned} & \text{FastBconv}(|t\gamma \cdot \text{ct}(\mathbf{s})|_q, q, \{t, \gamma\}) \times |-q^{-1}|_m \bmod m \\ &= (|t\gamma \cdot \text{ct}(\mathbf{s})|_q + q\mathbf{a}) \times |-q^{-1}|_m \bmod m \\ &= \frac{t\gamma[\text{ct}(\mathbf{s})]_q - |t\gamma \cdot \text{ct}(\mathbf{s})|_q - q\mathbf{a}}{q} \bmod m \text{ (exact division)} \\ &= \left( \left\lfloor \frac{t\gamma[\text{ct}(\mathbf{s})]_q}{q} \right\rfloor - \mathbf{a} \right) \bmod m \\ &= \left( \left\lfloor \frac{t\gamma[\text{ct}(\mathbf{s})]_q}{q} \right\rfloor - \mathbf{e} \right) \bmod m \end{aligned}$$

where  $\mathbf{e}_i \in \{\mathbf{a}_i, \mathbf{a}_i + 1\}$ , i.e.  $\mathbf{e}_i \in [0, k]$ . To conclude the proof, it suffices to use the equality  $\Delta t = q - |q|_t$ . That way, one can write  $t[\text{ct}(\mathbf{s})]_q = q([\mathbf{m}]_t + t\mathbf{r}) + \mathbf{v}_c$ , and the second equality of (3) follows.

### A.2 Lemma 2

By hypothesis, we have  $-\gamma(\frac{1}{2} - \varepsilon) - k \leq (\gamma\frac{\mathbf{v}_c}{q} - \mathbf{e})_i \leq \gamma(\frac{1}{2} - \varepsilon)$  for any  $i \in [0, n-1]$ . To have  $\lfloor \gamma\frac{\mathbf{v}_c}{q} \rfloor - \mathbf{e} = \lfloor \gamma\frac{\mathbf{v}_c}{q} \rfloor - \mathbf{e}$ , we then require that  $-\lfloor \frac{\gamma}{2} \rfloor - \frac{1}{2} \leq (\gamma\frac{\mathbf{v}_c}{q} - \mathbf{e})_i < \lfloor \frac{\gamma-1}{2} \rfloor + \frac{1}{2}$ . Hence, a sufficient condition is given by:

$$\begin{cases} \gamma(\frac{1}{2} - \varepsilon) < \lfloor \frac{\gamma-1}{2} \rfloor + \frac{1}{2} \\ -\lfloor \frac{\gamma}{2} \rfloor - \frac{1}{2} \leq -\gamma(\frac{1}{2} - \varepsilon) - k \end{cases} \Leftrightarrow (\gamma \text{ odd}) \begin{cases} \gamma\varepsilon > 0 \\ \gamma\varepsilon \geq k \end{cases}, (\gamma \text{ even}) \begin{cases} \gamma\varepsilon > \frac{1}{2} \\ \gamma\varepsilon \geq k - \frac{1}{2} \end{cases} .$$

### A.3 Theorem 1

According to Lem. 2, the  $\gamma$ -correction technique works as long as  $\gamma(\frac{1}{2} - \frac{\|\mathbf{v}_c\|_\infty}{q}) \geq k \Leftrightarrow \|\mathbf{v}_c\|_\infty \leq q(\frac{1}{2} - \frac{k}{\gamma})$ . Moreover,  $\|\mathbf{v}_c\|_\infty = \|t\mathbf{v} - |q|_t[\mathbf{m}]_t\|_\infty \leq t\|\mathbf{v}\|_\infty + |q|_t\frac{t}{2}$ . Then, the bound (5) follows. The lower bound on  $\gamma$  guarantees that the bound (5) for the noise is positive.

### A.4 Lemma 3

Let's denote  $\tilde{\mathbf{v}}_c := \lfloor \gamma\frac{\mathbf{v}_c}{q} \rfloor - \mathbf{e}$ . By computing (3) modulo  $\gamma t$ , we obtain  $\mathbf{z} = |\gamma[\mathbf{m}]_t + \tilde{\mathbf{v}}_c|_{\gamma t}$ . We notice that  $\mathbf{z} = |\gamma|\mathbf{m}|_t + \tilde{\mathbf{v}}_c|_{\gamma t}$ . Indeed,  $\gamma|\mathbf{m}|_t = \gamma([\mathbf{m}]_t + t\mathbf{a})$ , where  $\mathbf{a}_i \in \{0, 1\}$ . Thus,  $\gamma t\mathbf{a}$  vanishes modulo  $\gamma t$ . More generally,  $|\gamma|\mathbf{m}|_t + \tilde{\mathbf{v}}_c|_{\gamma t} = \gamma t\mathbf{b} + \gamma|\mathbf{m}|_t + \tilde{\mathbf{v}}_c$  with  $\mathbf{b}_i \in \{0, 1\}$  ( $\mathbf{b}_i = 1 \Leftrightarrow ((\gamma|\mathbf{m}|_t)_i = 0 \text{ and } (\tilde{\mathbf{v}}_c)_i < 0)$ ).

The integer  $\gamma$  has been chosen such that each coefficient of  $\tilde{\mathbf{v}}_c$  lies in  $[-\gamma/2, \gamma/2)$ . Thus, it is clear that any coefficient of  $\tilde{\mathbf{v}}_c + \frac{\gamma}{2}\mathbf{v}_1$  lies in  $[0, \gamma)$ .

Consequently, we have  $\mathbf{z} + \frac{\gamma}{2}\mathbf{v}_1 = \gamma t\mathbf{b} + \gamma|\mathbf{m}|_t + \mathbf{w}$  where  $\mathbf{w}_i \in [0, \gamma)$  and  $\mathbf{b}_i \geq 0$  for any  $i$ . Furthermore, it can be noticed that  $|\mathbf{z} + \frac{\gamma}{2}\mathbf{v}_1|_{\gamma t} = \gamma|\mathbf{m}|_t + \mathbf{w}$ . Finally, since  $\gamma$  is a power of 2, we have  $(\mathbf{z} + \frac{\gamma}{2}\mathbf{v}_1) \gg \log_2(\gamma) = t\mathbf{b} + |\mathbf{m}|_t$  (and also that  $|\mathbf{z} + \frac{\gamma}{2}\mathbf{v}_1|_{\gamma t} \gg \log_2(\gamma) = |\mathbf{m}|_t$ ). By applying a centered reduction modulo  $t$ , (6) follows.

#### A.5 Lemma 4

Algorithm 2 performs in  $\mathcal{B}_{\text{sk}}$  the computation of  $\frac{[\mathbf{c}\bar{m}]_q + q\mathbf{u} + q[-([\mathbf{c}\bar{m}]_q + q\mathbf{u})/q]_{\bar{m}}}{\bar{m}}$ . This quantity is clearly congruent to  $\mathbf{c}$  modulo  $q$ . In accordance with hypothesis (8), its norm is bounded by  $\frac{q(1/2 + \tau + \bar{m}/2)}{\bar{m}} \leq \frac{q}{2}(1 + \rho)$ .

#### A.6 Lemma 5

We recall that  $\text{FastBconv}(|t \cdot \mathbf{ct}'_\star[j]|_q, q, \mathcal{B}_{\text{sk}})$  outputs  $|t \cdot \mathbf{ct}'_\star[j]|_q + q\mathbf{u}$  with  $\|\mathbf{u}\|_\infty \leq k - 1$ . Then, the proof is complete by using the general equalities  $\frac{x - |x|_q}{q} = \lfloor \frac{x}{q} \rfloor = \lfloor \frac{x}{q} \rfloor + \tau$ ,  $\tau \in \{-1, 0\}$ .

#### A.7 Proposition 1

The following noise analysis is inspired from the one provided in [4]. Some of the tools and bounds from there are re-used here. In the following, we write  $\mathbf{ct}'_i = (\mathbf{c}'_{i,0}, \mathbf{c}'_{i,1})$ . In particular, we have  $\mathbf{ct}'_\star = (\mathbf{c}'_{1,0}\mathbf{c}'_{2,0}, \mathbf{c}'_{1,1}\mathbf{c}'_{2,0} + \mathbf{c}'_{1,0}\mathbf{c}'_{2,1}, \mathbf{c}'_{1,1}\mathbf{c}'_{2,1})$ . By hypothesis, each  $\mathbf{c}'_{i,j}$  satisfies  $\|\mathbf{c}'_{i,j}\|_\infty \leq \frac{q}{2}(1 + \rho)$ . In particular, the bound in (9) comes from the fact that  $\|\mathbf{c}'_{1,1}\mathbf{c}'_{2,0} + \mathbf{c}'_{1,0}\mathbf{c}'_{2,1}\|_\infty \leq \delta \frac{q^2}{2}(1 + \rho)^2$ .

Since  $\mathbf{ct}'_i = \mathbf{ct}_i \bmod q$  and  $\|\mathbf{ct}'_i\|_\infty \leq \frac{q}{2}(1 + \rho)$ , and by using  $\|\mathbf{s}\|_\infty \leq B_{\text{key}}$ ,  $\|\mathbf{v}_i\|_\infty < \frac{\Delta}{2} < \frac{q}{2t}$ ,  $\|\Delta[\mathbf{m}]_t\|_\infty < \frac{t\Delta}{2} < \frac{q}{2}$  and  $t \geq 2$ , we can write:

$$\mathbf{ct}'_i(\mathbf{s}) = \mathbf{c}'_{i,0} + \mathbf{c}'_{i,1}\mathbf{s} = \Delta[\mathbf{m}_i]_t + \mathbf{v}_i + q\mathbf{r}_i$$

with  $\|\mathbf{r}_i\|_\infty < r_\infty := \frac{1+\rho}{2}(1 + \delta B_{\text{key}}) + 1$ . For our purpose, we use some convenient notations and bounds from [4], but applicable to the present context:

$$\begin{cases} \mathbf{v}_1\mathbf{v}_2 = [\mathbf{v}_1\mathbf{v}_2]_\Delta + \Delta\mathbf{r}_v, & \|\mathbf{r}_v\|_\infty < \frac{\delta}{2} \min \|\mathbf{v}_i\|_\infty + \frac{1}{2} , \\ [\mathbf{m}_1]_t[\mathbf{m}_2]_t = [\mathbf{m}_1\mathbf{m}_2]_t + t\mathbf{r}_m, & \|\mathbf{r}_m\|_\infty < \frac{1}{2}\delta t . \end{cases} \quad (21)$$

By noticing that  $\mathbf{ct}'_\star(\mathbf{s}) = (\mathbf{ct}'_1 \star \mathbf{ct}'_2)(\mathbf{s}) = \mathbf{ct}'_1(\mathbf{s}) \times \mathbf{ct}'_2(\mathbf{s})$ , we obtain:

$$\begin{aligned} \mathbf{ct}'_\star(\mathbf{s}) &= \Delta^2[\mathbf{m}_1]_t[\mathbf{m}_2]_t + \Delta([\mathbf{m}_1]_t\mathbf{v}_2 + [\mathbf{m}_2]_t\mathbf{v}_1) \\ &\quad + q\Delta([\mathbf{m}_1]_t\mathbf{r}_2 + [\mathbf{m}_2]_t\mathbf{r}_1) + \mathbf{v}_1\mathbf{v}_2 + q^2\mathbf{r}_1\mathbf{r}_2 + q(\mathbf{v}_1\mathbf{r}_2 + \mathbf{v}_2\mathbf{r}_1) . \end{aligned}$$

Then, by using (21) and  $\Delta t = q - |q|_t$ , we deduce that:

$$\begin{aligned} \frac{t}{q}\mathbf{ct}'_\star(\mathbf{s}) &= \Delta[\mathbf{m}_1\mathbf{m}_2]_t + q(\mathbf{r}_m + [\mathbf{m}_1]_t\mathbf{r}_2 + [\mathbf{m}_2]_t\mathbf{r}_1 + t\mathbf{r}_1\mathbf{r}_2) \\ &\quad - \frac{|q|_t\Delta}{q}[\mathbf{m}_1\mathbf{m}_2]_t + (\frac{|q|_t}{q} - 2)|q|_t\mathbf{r}_m \\ &\quad + (1 - \frac{|q|_t}{q})([\mathbf{m}_1]_t\mathbf{v}_2 + [\mathbf{m}_2]_t\mathbf{v}_1) - |q|_t([\mathbf{m}_1]_t\mathbf{r}_2 + [\mathbf{m}_2]_t\mathbf{r}_1) \\ &\quad + \frac{t}{q}[\mathbf{v}_1\mathbf{v}_2]_\Delta + (1 - \frac{|q|_t}{q})\mathbf{r}_v + t(\mathbf{v}_1\mathbf{r}_2 + \mathbf{v}_2\mathbf{r}_1) . \end{aligned}$$

Thus,

$$\widetilde{\text{ct}}_{mult}(\mathbf{s}) = \text{DR}_2(\text{ct}'_*)(\mathbf{s}) = \frac{t}{q}\text{ct}'_*(\mathbf{s}) + \mathbf{r}_a = \Delta[\mathbf{m}_1\mathbf{m}_2]_t + \tilde{\mathbf{v}}_{mult} \pmod{q}$$

with  $\mathbf{r}_a := (\text{DR}_2(\text{ct}'_*) - \frac{t}{q}\text{ct}'_*)(\mathbf{s}) = \sum_{i=0}^2 \left( \lfloor \frac{t}{q}\text{ct}'_*[i] \rfloor - \frac{t}{q}\text{ct}'_*[i] \right) \mathbf{s}^i$  and

$$\begin{aligned} \tilde{\mathbf{v}}_{mult} := & \mathbf{r}_a - \frac{|q|_t \Delta}{q} [\mathbf{m}_1\mathbf{m}_2]_t + \left( \frac{|q|_t}{q} - 2 \right) |q|_t \mathbf{r}_m \\ & + \left( 1 - \frac{|q|_t}{q} \right) ([\mathbf{m}_1]_t \mathbf{v}_2 + [\mathbf{m}_2]_t \mathbf{v}_1) - |q|_t ([\mathbf{m}_1]_t \mathbf{r}_2 + [\mathbf{m}_2]_t \mathbf{r}_1) \\ & + \frac{t}{q} [\mathbf{v}_1\mathbf{v}_2]_\Delta + \left( 1 - \frac{|q|_t}{q} \right) \mathbf{r}_v + t(\mathbf{v}_1 \mathbf{r}_2 + \mathbf{v}_2 \mathbf{r}_1) . \end{aligned} \quad (22)$$

Below, some useful bounds are given.

$$\begin{cases} \| [\mathbf{m}_1]_t \mathbf{v}_2 + [\mathbf{m}_2]_t \mathbf{v}_1 \|_\infty \leq \frac{\delta t}{2} (\|\mathbf{v}_1\|_\infty + \|\mathbf{v}_2\|_\infty) , \\ \| \mathbf{v}_1 \mathbf{r}_2 + \mathbf{v}_2 \mathbf{r}_1 \|_\infty < \delta r_\infty (\|\mathbf{v}_1\|_\infty + \|\mathbf{v}_2\|_\infty) . \end{cases}$$

Next, we set a bound for each term of (22), then it suffices to put them all together to obtain (11).

$$\begin{aligned} \| \mathbf{r}_a \|_\infty & \leq \frac{1}{2} (1 + \delta B_{key} + \delta^2 B_{key}^2) , \\ \| -\frac{|q|_t \Delta}{q} [\mathbf{m}_1\mathbf{m}_2]_t + \frac{t}{q} [\mathbf{v}_1\mathbf{v}_2]_\Delta \|_\infty & \leq \frac{|q|_t \Delta t}{2q} + \frac{t \Delta}{2q} < \frac{1}{2} (|q|_t + 1) , \\ \| \left( \frac{|q|_t}{q} - 2 \right) |q|_t \mathbf{r}_m \|_\infty & \leq 2|q|_t \|\mathbf{r}_m\|_\infty < \delta t |q|_t , \\ \| \left( 1 - \frac{|q|_t}{q} \right) \mathbf{r}_v \|_\infty & \leq \|\mathbf{r}_v\|_\infty < \frac{\delta}{2} \min \|\mathbf{v}_i\|_\infty + \frac{1}{2} , \\ \| t(\mathbf{v}_1 \mathbf{r}_2 + \mathbf{v}_2 \mathbf{r}_1) \|_\infty & < \delta t r_\infty (\|\mathbf{v}_1\|_\infty + \|\mathbf{v}_2\|_\infty) , \\ \| \left( 1 - \frac{|q|_t}{q} \right) ([\mathbf{m}_1]_t \mathbf{v}_2 + [\mathbf{m}_2]_t \mathbf{v}_1) \|_\infty & < \frac{\delta t}{2} (\|\mathbf{v}_1\|_\infty + \|\mathbf{v}_2\|_\infty) , \\ \| -|q|_t ([\mathbf{m}_1]_t \mathbf{r}_2 + [\mathbf{m}_2]_t \mathbf{r}_1) \|_\infty & < |q|_t \delta t r_\infty . \end{aligned}$$

### A.8 Lemma 6

We set  $\mathcal{B} = \{m_1, \dots, m_\ell\}$  ( $|\mathcal{B}| = \ell$ ). The case  $x \geq 0$  and  $\lambda = 1$  is the classical case of Shenoy and Kumaresan's conversion. We recall that, by definition,  $\text{FastBconv}(x, \mathcal{B}, \cdot) = \sum_{i=1}^{\ell} |x \frac{m_i}{M}|_{m_i} \frac{M}{m_i}$ . There exists an integer  $0 \leq \alpha \leq \ell - 1$  such that  $\sum_{i=1}^{\ell} |x \frac{m_i}{M}|_{m_i} \frac{M}{m_i} = x + \alpha M$ . By inverting this equality modulo  $m_{\text{sk}}$ , it would suffice, in this case, that  $m_{\text{sk}} > \ell$  to enable us to recover  $\alpha$  by noticing that  $\alpha = |\alpha|_{m_{\text{sk}}} = |\alpha_{\text{sk}, x}|_{m_{\text{sk}}}$ .

Let's consider the general case  $|x| < \lambda M$ . Here, the possible negativity of  $x$  is the reason why we have to compute a centered remainder modulo  $m_{\text{sk}}$  in (12).

First, we notice that the residues of  $x$  in  $\mathcal{B}$  are actually those of  $|x|_M = x + \mu M \in [0, M)$ , where  $\mu$  is an integer lying in  $[-\lfloor \lambda \rfloor, \lceil \lambda \rceil]$ . Therefore, we deduce that

$$\sum_{i=1}^{\ell} |x \frac{m_i}{M}|_{m_i} \frac{M}{m_i} = |x|_M + \alpha_{\text{sk}, |x|_M} M = (x + \mu M) + \alpha_{\text{sk}, (x + \mu M)} M$$

with  $0 \leq \alpha_{\text{sk}, (x + \mu M)} \leq \ell - 1$ . Denoting  $|x|_{m_{\text{sk}}}$  by  $x_{\text{sk}}$ , it follows that the quantity computed in (12) is the following one:

$$\begin{aligned} [(\sum_{i=1}^{\ell} |x \frac{m_i}{M}|_{m_i} \frac{M}{m_i} - x_{\text{sk}}) M^{-1}]_{m_{\text{sk}}} & = [(\sum_{i=1}^{\ell} |x \frac{m_i}{M}|_{m_i} \frac{M}{m_i} - (x_{\text{sk}} + \mu M)) M^{-1} + \mu]_{m_{\text{sk}}} \\ & = [\alpha_{\text{sk}, (x + \mu M)} + \mu]_{m_{\text{sk}}} . \end{aligned}$$

It remains to show that  $[\alpha_{\mathbf{sk},(x+\mu M)} + \mu]_{m_{\mathbf{sk}}} = \alpha_{\mathbf{sk},(x+\mu M)} + \mu$  or, in other words, that  $-\lfloor \frac{m_{\mathbf{sk}}}{2} \rfloor \leq \alpha_{\mathbf{sk},(x+\mu M)} + \mu \leq \lfloor \frac{m_{\mathbf{sk}}-1}{2} \rfloor$ . But by hypothesis on  $m_{\mathbf{sk}}$ , and because  $\ell \geq 1$ , we can write  $m_{\mathbf{sk}} \geq 2(\ell + \lceil \lambda \rceil) \geq 2\lceil \lambda \rceil + 1$ . Then,

$$\begin{cases} \alpha_{\mathbf{sk},(x+\mu M)} + \mu \leq \ell - 1 + \lceil \lambda \rceil \leq \frac{m_{\mathbf{sk}}}{2} - 1 \leq \lfloor \frac{m_{\mathbf{sk}}-1}{2} \rfloor , \\ \alpha_{\mathbf{sk},(x+\mu M)} + \mu \geq -\lceil \lambda \rceil \geq -\frac{m_{\mathbf{sk}}-1}{2} \geq -\lfloor \frac{m_{\mathbf{sk}}}{2} \rfloor . \end{cases}$$

Thus,  $[\alpha_{\mathbf{sk},(x+\mu M)} + \mu]_{m_{\mathbf{sk}}} = \alpha_{\mathbf{sk},(x+\mu M)} + \mu$ , and it follows that, by computing the right member of (13), we obtain the following equalities:

$$\begin{aligned} \sum_{i=1}^{\ell} |x \frac{m_i}{M}|_{m_i} \frac{M}{m_i} - [\alpha_{\mathbf{sk},(x+\mu M)} + \mu]_{m_{\mathbf{sk}}} M &= (x + \mu M) + \alpha_{\mathbf{sk},(x+\mu M)} M \\ &\quad - (\alpha_{\mathbf{sk},(x+\mu M)} + \mu) M \\ &= x . \end{aligned}$$

### A.9 Proposition 2

By using (9), we have  $\|\text{DR}_2(\mathbf{ct}'_*)\|_{\infty} \leq \|\frac{t}{q} \mathbf{ct}'_*\|_{\infty} + \frac{1}{2} \leq \delta t \frac{q}{2} (1 + \rho)^2 + \frac{1}{2}$ . Lemma 6 concludes the proof.

### A.10 Lemma 8

First, we have

$$\langle \xi_q(\mathbf{c}), -(\vec{\mathbf{e}} + \vec{\mathbf{d}} \mathbf{s}) \rangle + \mathbf{s} \langle \xi_q(\mathbf{c}), \vec{\mathbf{d}} \rangle \pmod{q} = -\langle \xi_q(\mathbf{c}), \vec{\mathbf{e}} \rangle \pmod{q} .$$

Second, by using  $q_i < 2^{\nu}$ , we obtain

$$\begin{aligned} \|\langle \xi_q(\mathbf{c}), \vec{\mathbf{e}} \rangle\|_{\infty} &= \|\sum_{i=1}^k |\mathbf{c}^{\frac{q_i}{q}}|_{q_i} \mathbf{e}_i\|_{\infty} < \sum_{i=1}^k \delta q_i \|\mathbf{e}_i\|_{\infty} \leq \delta B_{err} \sum_{i=1}^k q_i \\ &< \delta B_{err} k 2^{\nu} . \end{aligned}$$

### A.11 Proposition 3

At this point, we are evaluating  $\text{Relin}_{\text{RNS}}(\widetilde{\mathbf{ct}}_{mult} + \mathbf{b})$ . We denote  $\widetilde{\mathbf{ct}}_{mult} = (\mathbf{c}_0, \mathbf{c}_1, \mathbf{c}_2)$ . A first remark is that we can write

$$\xi_q(\mathbf{c}_2 + \mathbf{b}_2) = \xi_q(\mathbf{c}_2) + \xi_q(\mathbf{b}_2) - \vec{\mathbf{u}}$$

where  $\vec{\mathbf{u}} = (\mathbf{u}_1, \dots, \mathbf{u}_k)$  is such that, for any  $(i, j) \in [1, k] \times [0, n-1]$ ,  $(\mathbf{u}_i)_j \in \{0, q_i\}$ . In particular, it can be noticed that  $\|\mathbf{b}_2 \frac{q_i}{q}|_{q_i} - \mathbf{u}_i\|_{\infty} < q_i$ , and that  $\langle \vec{\mathbf{u}}, (\mathbf{s}^2 \frac{q}{q_1}, \dots, \mathbf{s}^2 \frac{q}{q_k}) \rangle \equiv 0 \pmod{q}$ . This latter fact leads to the following equality:

$$\langle \vec{\mathbf{u}}, \text{rlk}_{\text{RNS}}[0] \rangle + \langle \vec{\mathbf{u}}, \text{rlk}_{\text{RNS}}[1] \rangle \mathbf{s} = \langle \vec{\mathbf{u}}, \vec{\mathbf{e}} \rangle = \sum_{i=1}^k \mathbf{u}_i \mathbf{e}_i \pmod{q} .$$

Consequently, we have that, in  $\mathcal{R}_q \times \mathcal{R}_q$ ,

$$\begin{aligned} \text{Relin}_{\text{RNS}}(\widetilde{\mathbf{ct}}_{mult} + \mathbf{b}) &= \text{Relin}_{\text{RNS}}(\widetilde{\mathbf{ct}}_{mult}) + \text{Relin}_{\text{RNS}}(\mathbf{b}) \\ &\quad - (\langle \vec{\mathbf{u}}, \text{rlk}_{\text{RNS}}[0] \rangle, \langle \vec{\mathbf{u}}, \text{rlk}_{\text{RNS}}[1] \rangle) . \end{aligned}$$



Thus, a part of the noise comes from the following extra term:

$$\begin{aligned}
& \left\| \mathbf{Relin}_{\text{RNS}}(\mathbf{b})(\mathbf{s}) - \langle \vec{\mathbf{u}}, \mathbf{rlk}_{\text{RNS}}[0] \rangle - \langle \vec{\mathbf{u}}, \mathbf{rlk}_{\text{RNS}}[1] \rangle \mathbf{s} \pmod{q} \right\|_{\infty} \\
&= \left\| \mathbf{b}_0 + \langle \xi_q(\mathbf{b}_2) - \vec{\mathbf{u}}, \mathbf{rlk}_{\text{RNS}}[0] \rangle + \mathbf{s}(\mathbf{b}_1 + \langle \xi_q(\mathbf{b}_2) - \vec{\mathbf{u}}, \mathbf{rlk}_{\text{RNS}}[1] \rangle) \pmod{q} \right\|_{\infty} \\
&= \left\| \mathbf{b}_0 + \mathbf{b}_1 \mathbf{s} + \mathbf{b}_2 \mathbf{s}^2 - \sum_{i=1}^k (|\mathbf{b}_2 \frac{q_i}{q}|_{q_i} - \mathbf{u}_i) \mathbf{e}_i \pmod{q} \right\|_{\infty} \\
&< k(1 + \delta B_{key} + \delta^2 B_{key}^2) + \delta B_{err} \sum_{i=1}^k q_i \\
&< k(1 + \delta B_{key} + \delta^2 B_{key}^2) + \delta B_{err} k 2^{\nu} .
\end{aligned}$$

Lemma 8 brings the rest of the noise.

## B Additional elements about RNS homomorphic multiplication

### B.1 Combining two levels of decomposition within step 2

To reduce the noise growth due to the relinearisation step a bit more, we can integrate another level of decomposition in radix  $\omega$  where  $\omega = 2^{\theta} \ll 2^{\nu}$  as efficiently as in the original scheme by doing it on the residues, because they are handled through the classical binary positional system. By denoting  $\ell_{\omega, 2^{\nu}} = \lceil \frac{\nu}{\theta} \rceil$ , each polynomial  $|\mathbf{c} \frac{q_i}{q}|_{q_i}$  is decomposed into the vector of polynomials  $(\left[ \left[ |\mathbf{c} \frac{q_i}{q}|_{q_i} \omega^{-z} \right]_{\omega} \right]_{z \in [0, \dots, \ell_{\omega, 2^{\nu}} - 1]})$ , and the new decomposition function is defined by:

$$\mathcal{D}_{\text{RNS}, \omega, q}(\mathbf{c}) = \left( \mathbf{d}_i^z = \left[ \left[ |\mathbf{c} \frac{q_i}{q}|_{q_i} \omega^{-z} \right]_{\omega} \right]_{z \in [0, \dots, \ell_{\omega, 2^{\nu}} - 1]} \right)_{i \in [1, k]} .$$

Therefore, each term  $|\mathbf{s}^2 \frac{q}{q_i} - (\mathbf{e}_i + \mathbf{s} \mathbf{a}_i)|_{q_j}$  in  $\mathbf{rlk}_{\text{RNS}}[0]$  has to be replaced by

$$\left( |\mathbf{s}^2 \frac{q}{q_i} \omega^z - (\mathbf{e}_i^z + \mathbf{s} \mathbf{a}_i^z)|_{q_j} \right)_z, \quad z = 0, \dots, \ell_{\omega, 2^{\nu}} - 1, \quad \mathbf{e}_i^z \leftarrow \chi_{err}, \mathbf{a}_i^z \leftarrow \mathcal{U}(\mathcal{R}_q) .$$

It follows that the extra noise is now bounded by:

$$\left\| \langle \mathcal{D}_{\text{RNS}, \omega, q}(\mathbf{c}), \vec{\mathbf{e}} \rangle \right\|_{\infty} = \left\| \sum_{i=1}^k \sum_{z=0}^{\ell_{\omega, 2^{\nu}} - 1} \mathbf{d}_i^z \mathbf{e}_i^z \right\|_{\infty} < \delta B_{err} \omega k \ell_{\omega, 2^{\nu}} .$$

In other words, the term  $2^{\nu}$  in (16) is replaced by  $\omega \ell_{\omega, 2^{\nu}}$ . Even if such optimization may be not really relevant in practice, this is another argument for saying that any trick and technique of the original scheme can be efficiently adapted to a pure RNS approach.

### B.2 Some details about complexity

We analyse the cost of a multi-precision variant, in order to estimate the benefits of the new RNS variant of multiplication in terms of computational cost.

The product  $\mathbf{ct}_\star = \mathbf{ct}_1 \star \mathbf{ct}_2 = (\mathbf{c}_{1,0}, \mathbf{c}_{1,1}) \star (\mathbf{c}_{2,0}, \mathbf{c}_{2,1})$  in MP variant is advantageously performed in RNS, in order to benefit from NTT. So, the MP variant considered here is assumed to use a second base  $\mathcal{B}'$  such that  $qM' > \|\mathbf{ct}_\star\|_\infty$ . By taking centered remainders modulo  $q$ , we consider  $\|\mathbf{ct}_i\|_\infty \leq \frac{q}{2}$ . Then  $\mathcal{B}'$  must verify in particular that  $\|\mathbf{ct}_\star[1] = \mathbf{c}_{1,0}\mathbf{c}_{2,1} + \mathbf{c}_{1,1}\mathbf{c}_{2,0}\|_\infty \leq 2\delta\frac{q^2}{4} < qM'$ . Thus,  $|\mathcal{B}'|$  has to be at least equal to  $k+1$  (notice that, in RNS variant, we also need to have  $|\mathcal{B}_{sk}| = k+1$  for dealing with the multiplication).

The conversion, from  $q$  to  $\mathcal{B}'$ , of each  $\mathbf{ct}_i$  has to be as exact as possible in order to reduce the noise growth. It can be done by computing  $[\sum_{i=1}^k |\mathbf{c}_{a,b}|_{q_i} \times (|\frac{q_i}{q}|_{q_i} \frac{q}{q_i})]_q$  in  $\mathcal{B}'$ . The terms  $(|\frac{q_i}{q}|_{q_i} \frac{q}{q_i})$  are precomputable and their size is  $k$  words ( $\log_2(q)$  bits). Thus, the sum involves  $k^2 nEM$ . The reduction modulo  $q$  can be performed by using an efficient reduction as described in [19], reducing to around 2 multiplications of  $k$ -word integers, that is  $\mathcal{O}(k^{1+\varepsilon})nEM$  (where  $\varepsilon$  stands for complexity of multi-precision multiplication in radix-base  $2^\nu$ ; e.g.  $\varepsilon = 1$  for the schoolbook multiplication). Next, the  $k$ -word value is reduced modulo each 1-word element of  $\mathcal{B}'$ , through around  $2knEM$  for the whole set of coefficients. Finally, this procedure has to be made four times. Its total cost is around  $(4k^2 + \mathcal{O}(k^{1+\varepsilon}))nEM$ .

Next, the product  $\mathbf{ct}_1 \star \mathbf{ct}_2$  is done in  $q \cup \mathcal{B}'$ . First,  $4(2k+1)NTT$  are applied. Second, by using a Karatsuba like trick, the product is achieved by using only  $3 \times (2k+1)nIMM$ . Third,  $3(2k+1)\text{invNTT}$  are applied to recover  $\mathbf{ct}_\star = (\mathbf{c}_{\star,0}, \mathbf{c}_{\star,1}, \mathbf{c}_{\star,2})$  in coefficient representation.

The next step is the division and rounding of the three polynomials  $\mathbf{c}_{\star,i}$ 's. A lift from  $q \cup \mathcal{B}'$  to  $\mathbb{Z}$  is required, for a cost of  $3(2k+1)^2 nEM$ .  $\frac{t}{q}$  can be precomputed with around  $3k+1$  words of precision to ensure a correct rounding. Thus, a product  $\frac{t}{q} \times \mathbf{c}_{\star,i}$  is achieved with  $\mathcal{O}(k^{1+\varepsilon})nEM$ . After, the rounding of  $\mathbf{c}_{\star,0}$  and  $\mathbf{c}_{\star,1}$  are reduced in RNS base  $q$  by  $2 \times 2knEM$ .

The relinearisation step (15) can be done in each RNS channel of  $q$ . By assuming that  $\omega = 2^\nu$ , we would have  $\ell_{\omega,q} = k$ . The computation of the vector  $\mathcal{D}_{\omega,q}(\lfloor \frac{t}{q} \mathbf{c}_2 \rfloor)$  reduces to shifting. The two scalar products in  $\text{Relin}_{\text{FV}}$ , with an output in coefficient representation, require  $k\ell_{\omega,q}NTT + 2k^2 nIMM + 2k\text{invNTT}$ . Thus, the total cost is at most the following one:

$$\begin{aligned} \text{Cost}(\text{Mult}_{\text{MP}}) &= (k\ell_{\omega,q} + 8k + 4)NTT + (8k + 3)\text{invNTT} \\ &\quad + [2k^2 + 6k + 3]nIMM + [40k^2 + \mathcal{O}(k^{1+\varepsilon})]nEM. \end{aligned}$$

Let's analyse the cost of Alg. 3. The fast conversions at **S0** from  $q$  to  $\mathcal{B}_{sk} \cup \{\tilde{m}\}$  require  $4 \times k(k+2)nIMM$ . Next, the reduction of  $q$ -overflows at step **S1** requires  $4 \times (k+1)nIMM$ . The product of ciphertexts  $\mathbf{ct}'_1 \star \mathbf{ct}'_2$  (**S2**) in  $q \cup \mathcal{B}_{sk}$  requires the same cost as for MP variant, that is  $4(2k+1)NTT + 3(2k+1)nIMM + 3(2k+1)\text{invNTT}$ .

Let's analyse the cost of steps **S3**, **S4** and **S5**. With adequate pre-computed data, the base conversion in **S3** can integrate the flooring computation in  $\mathcal{B}_{sk}$ . So, **S3** is achievable with  $3 \times k(k+1)nIMM$ . The exact **FastBconvSK** at step **S4** basically reduces to a fast conversion from  $\mathcal{B}$  to  $q_{sk}$ , followed by a second one

from  $m_{sk}$  to  $q$ . So, this is achieved with  $3 \times k(k + 2)n\text{IMM}$ . In **S5**, we already have the vector  $\xi_q(c_2)$  which is involved in the fast conversion in step **S3**. Indeed, the function  $\xi_q$  is an automorphism of  $\mathcal{R}_q$ . So, data in  $q$  can stay in this form throughout the computations. The two scalar products in  $\text{Relin}_{\text{RNS}}$  involve  $k^2\text{NTT} + 2k^2n\text{IMM} + 2k\text{invNTT}$ , exactly like the relinearisation step in MP variant. And finally  $2kn\text{IMM}$  are needed to manage the Montgomery representation, in  $q$ , with respect to  $\tilde{m}$ .

$$\begin{aligned} \text{Cost}(\text{Mult}_{\text{RNS}}) &= (k^2 + 8k + 4)\text{NTT} + (8k + 3)\text{invNTT} \\ &\quad + [10k^2 + 25k + 7]n\text{IMM}. \end{aligned}$$

To summarize, the RNS variant decreases the computational cost of the whole homomorphic multiplication algorithm except the parts concerning polynomial multiplications. Also, it involves as many NTT and invNTT as the MP variant. Even by considering an optimized multi-precision multiplication algorithm in MP variant (with sub-quadratic complexity), the asymptotic computational complexity remains dominated by the  $(k^2 + \mathcal{O}(k))n\text{NTT}$ . Finally, the MP and RNS variants are asymptotically equivalent when  $n \rightarrow +\infty$ , but the pure RNS variant is more flexible in terms of parallelization.

### C Timing results for decryption and multiplication

Tables 3 and 4 summarize the values chosen for  $\tilde{m}$  for all the configurations which have been implemented and tested, and they list sufficient sizes for  $m_{sk}$  and  $\gamma$  in these cases.

$n$	$k$	$t$	$\tilde{m}$	$L_{\text{RNS}}$	$\lceil \log_2(m_{sk}) \rceil$	$\gamma$
$2^{11}$	3	2	$2^8$	2	13	7
		$2^{10}$		1	22	7
$2^{12}$	6	2	$2^8$	5	14	13
		$2^{10}$		4	23	13
$2^{13}$	13	2	$2^8$	13	15	27
		$2^{10}$		9	24	27
$2^{14}$	26	2	$2^8$	25	17	53
		$2^{10}$		19	26	53
$2^{15}$	53	2	$2^{16}$	50	20	112
		$2^{10}$		38	29	107

Table 3: Parameters based on 30-bit moduli of `NFLlib` and depending on an  $a$  priori chosen  $\tilde{m}$ .

$n$	$k$	$t$	$\tilde{m}$	$L_{\text{RNS}} (L_{\text{std}})$	$\lceil \log_2(m_{sk}) \rceil$	$\gamma$
$2^{11}$	1	2	$2^8$	0 (0)	—	—
		$2^{10}$		0 (0)	—	—
$2^{12}$	3	2	$2^8$	4 (5)	14	7
		$2^{10}$		3 (3)	23	7
$2^{13}$	6	2	$2^8$	11 (11)	15	13
		$2^{10}$		8 (8)	24	13
$2^{14}$	12	2	$2^8$	23 (23)	16	25
		$2^{10}$		17 (17)	25	25
$2^{15}$	25	2	$2^{16}$	47 (47)	17	51
		$2^{10}$		37 (37)	26	52

Table 4: Parameters based on 62-bit moduli of `NFLlib` and depending on an  $a$  priori chosen  $\tilde{m}$ .

Table 5 lists timings and speed-ups of RNS vs MP variants. Also, timings of an RNS decryption and multiplication including the use of SIMD (Single Instruction Multiple Data) have been added.

In RNS variants (as well decryption as multiplication), the replacement of division and rounding by base conversions (i.e. matrix-vector multiplications) allows to benefit, easily and naturally, from concurrent computation. An RNS vector-matrix multiplication naturally owns two levels of parallelization: along the RNS channels, and along the dimension of the result. In `NFLlib`, an element of  $\mathcal{R}_q$  is stored in a (32-byte aligned) array `_data` in which the  $n$  first values are the coefficients of the polynomial in  $\mathcal{R}_{q_1}$ , and so forth and so on. Advanced Vector Extensions (`AVX2`) have been used to accelerate the computations.

An `AVX2` register is handled (for our purpose) through the type `_m256i`. This enables us to handle either 8x32-bit, or 4x64-bit, or again 16x16-bit integers concurrently. Given the configuration of the `_data` array, reading/writing communications between 256-bit `AVX2` registers and `_data` are the most efficient (through `_mm256_store_si256` and `_mm256_load_si256` intrinsics, which require the 32-byte alignment of `_data`) when the base conversion is parallelized along the dimension  $n$ . This is the way it has been implemented and tested. Moreover, this has been only done within the 32-bit implementations, for which more convenient intrinsic instructions are available.

Regarding the timings, the impact of `AVX2` remains quite moderate. This is because it is used to accelerate the parts of algorithms besides the NTT-based polynomial products which constitute the main cost. For decryption, the RNS variants, floating point and integer, are already very efficient, whereas the performance of `AVX2` variant depends on time-consuming loading operations from/to vectorial registers, explaining the small differences. About the multiplication, as expected the timings are converging when  $n$  grows, because of the cost of NTT's.

$n$	$\nu$	$k$	variant	Decryption (ms)	Speed-up	Multiplication (ms)	Speed-up
$2^{11}$	30	3	MP	1.153	—	13.809	—
			RNS-flp	0.192	6.005	—	—
			RNS	0.189	6.101	3.159	4.371
	62	1	RNS-AVX2	0.188	6.133	2.710	5.096
			MP	1.020	—	—	—
			RNS-flp	0.054	18.880	—	—
$2^{12}$	30	6	RNS	0.050	20.390	—	—
			MP	4.587	—	45.055	—
			RNS-flp	0.798	5.748	—	—
	62	3	RNS	0.789	5.814	15.614	2.886
			RNS-AVX2	0.775	5.919	13.737	3.280
			MP	3.473	—	28.168	—
$2^{13}$	30	13	RNS-flp	0.339	10.245	—	—
			RNS	0.326	10.653	7.688	3.664
			RNS-AVX2	3.637	4.413	88.589	2.462
	62	6	MP	16.051	—	218.103	—
			RNS-flp	3.732	4.301	—	—
			RNS	3.691	4.349	100.625	2.167
$2^{14}$	30	26	RNS-AVX2	3.637	4.413	88.589	2.462
			MP	10.945	—	92.093	—
			RNS-flp	1.552	7.052	—	—
	62	12	RNS	1.513	7.234	37.738	2.440
			MP	70.154	—	1,249.400	—
			RNS-flp	17.497	4.009	—	—
$2^{15}$	30	53	RNS	17.333	4.047	622.596	2.007
			RNS-AVX2	16.818	4.171	617.846	2.022
			MP	38.910	—	424.014	—
	62	25	RNS-flp	6.702	5.806	—	—
			RNS	6.494	5.992	206.511	2.053
			MP	364.379	—	8,396.080	—
$2^{15}$	30	53	RNS-flp	85.165	4.279	—	—
			RNS	81.225	4.486	4,923.220	1.705
			RNS-AVX2	72.665	5.015	5,063.920	1.658
	62	25	MP	180.848	—	2,680.535	—
			RNS-flp	33.310	5.429	—	—
			RNS	31.895	5.670	1,406.960	1.905

Table 5: Timing results.