



**HAL**  
open science

## An asymptotically faster version of FV supported on HPR

Jean-Claude Bajard, Julien Eynard, Paulo Martins, Leonel Sousa, Vincent Zucca

► **To cite this version:**

Jean-Claude Bajard, Julien Eynard, Paulo Martins, Leonel Sousa, Vincent Zucca. An asymptotically faster version of FV supported on HPR. ARITH-2020 - 27th IEEE International Symposium on Computer Arithmetic, Jun 2020, Portland, United States. hal-02883325

**HAL Id: hal-02883325**

**<https://hal.sorbonne-universite.fr/hal-02883325>**

Submitted on 29 Jun 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# An asymptotically faster version of FV supported on HPR

Jean-Claude Bajard<sup>\*</sup>, Julien Eynard<sup>†</sup>, Paulo Martins<sup>‡</sup>, Leonel Sousa<sup>‡</sup> and Vincent Zucca<sup>§</sup>

<sup>\*</sup>Sorbonne Université, CNRS, Inria, Institut de Mathématiques de Jussieu – Paris Rive Gauche, France.

jean-claude.bajard@sorbonne-universite.fr

<sup>†</sup>Univ. Grenoble Alpes, CEA, LETI, DSYS, CESTI, F-38000 Grenoble.

julien.eynard@cea.fr

<sup>‡</sup>INESC-ID, Instituto Superior Técnico, Universidade de Lisboa.

paulo.sergio@netcabo.pt, las@inesc-id.pt

<sup>§</sup>imec-COSIC KU Leuven.

vincent.zucca@kuleuven.be

**Abstract**—State-of-the-art implementations of homomorphic encryption exploit the Fan and Vercauteren (FV) scheme and the Residue Number System (RNS). While the RNS breaks down large integer arithmetic into smaller independent channels, its non-positional nature makes operations such as division and rounding hard to implement, and makes the representation of small values inefficient. In this work, we propose the application of the Hybrid Position-Residues Number System representation to the FV scheme. This is a positional representation of large radix where the digits are represented in RNS. It inherits the benefits from RNS and allows to accelerate the critical division and rounding operations while also making the representation of smaller values more compact. This directly benefits the decryption and the homomorphic multiplication procedures, reducing their asymptotic complexity, in dimension  $n$ , from  $\mathcal{O}(n^2 \log n)$  to  $\mathcal{O}(n \log n)$  and from  $\mathcal{O}(n^3 \log n)$  to  $\mathcal{O}(n^3)$ , respectively and has resulted in noticeable speedups when experimentally compared to related art RNS implementations.

**Index Terms**—Fan-Vercauteren scheme, Residue Number System, Homomorphic Encryption

## I. INTRODUCTION

Since Gentry’s uncovering of Fully Homomorphic Encryption (FHE) [1], a large amount of research has been dedicated to improving its efficiency. While the Fan and Vercauteren (FV) [2] proposal is one of the most efficient homomorphic schemes, it requires computationally expensive operations on polynomials with large coefficients. The Residue Number System (RNS) [3] allows to break down large integer arithmetic into many independent smaller channels. With it, additions and multiplications are computed independently for each channel, resulting in a reduction of complexity when compared to traditional positional representations. As a result, state-of-the-art implementations of FV are based on the RNS [4], [5].

However, the non-positional nature of RNS makes it more difficult to implement operations such as division and rounding, which are required by both FV homomorphic multiplication and decryption procedures. Furthermore, while for positional representations the size of the representation of a number is proportional to its magnitude, the RNS representation of a number is as large as the largest number representable in it. During the relinearisation procedure, which dominates

the asymptotic complexity of the homomorphic multiplication, polynomial coefficients need to be decomposed to values of smaller magnitude. Hence, after applying this decomposition, the RNS becomes a sub-optimal representation.

This work proposes to exploit instead the Hybrid Position-Residues Number System (HPR) [6] representation of numbers to accelerate the cryptographic operations of FV. This hybrid representation is of a positional nature, but each digit is represented in RNS. It not only inherits the parallel nature of RNS but it also makes the relinearisation, and division and rounding operations more efficient. As a result, in dimension  $n$ , the complexity of FV decryption is asymptotically reduced from  $\mathcal{O}(n^2 \log n)$  to  $\mathcal{O}(n \log n)$ . Similarly, the complexity of homomorphic multiplications is reduced from  $\mathcal{O}(n^3 \log n)$  to  $\mathcal{O}(n^3)$ . Complexity gains are translated in practice into speedups of up to 8.0 and 8.6 for the decryption operation and of up to 2.0 and 1.6 for the homomorphic multiplication, for the parameters considered in this article and two state-of-the-art schemes, with an upward trend for larger parameters.

The remainder of this article is organised as follows: the FV cryptosystem and the HPR representation are described in Section II. Based on these techniques, our construction is given in Section III. Afterwards, Section IV reviews related art, and Sections V and VI compare the proposed techniques with state-of-the-art methods from theoretical and experimental perspectives, respectively.

## II. BACKGROUND

The ring  $\mathcal{R} = \mathbb{Z}[X]/(X^n + 1)$  can be thought of as a set of polynomials with integer coefficients and degree smaller than  $n$ .  $\mathcal{R}_q = \mathcal{R}/q\mathcal{R}$  denotes the polynomials in  $\mathcal{R}$  whose coefficients are reduced modulo an integer  $q$ . When  $q \equiv 1 \pmod{2n}$ , elements of  $\mathcal{R}_q$  can be represented in the Number Theoretic Transform (NTT)-domain, which allows for coefficient-wise additions and multiplications [7]. Henceforth, boldface letters are used to denote polynomials and the infinity norm  $\|\cdot\|_\infty$  is the largest absolute value of the polynomial’s coefficients. The expansion factor associated to  $\mathcal{R}$  is defined as  $\delta_{\mathcal{R}} = \sup \left\{ \frac{\|\mathbf{a} \cdot \mathbf{b}\|_\infty}{\|\mathbf{a}\|_\infty \|\mathbf{b}\|_\infty}, \text{ for } \mathbf{a}, \mathbf{b} \in \mathcal{R} - \{0\} \right\}$ .  $\mathbf{a} \leftarrow D$  and

$\mathbf{a} \leftarrow \chi$  are used to denote that  $\mathbf{a}$  is drawn uniformly at random from the set  $D$  or according to a distribution  $\chi$ , respectively. The notation  $\cdot|_q$  (resp.  $[\cdot]_q$ ) is used to denote the residue modulo  $q$  in  $[0, q)$  (resp. in  $[-q/2, q/2)$ ). Moreover, for  $a \in \mathbb{Q}$ ,  $b = \lfloor a \rfloor$  is used to denote the closest integer  $b \in \mathbb{Z}$  to  $a$ . The computational complexity is herein evaluated in terms of the amount of Elementary Modular Multiplications (EMMs), whose operands fit a single machine word, floating point operations (FPOs) and forward and inverse NTTs required.

### A. FV cryptosystem

A private-key of the FV scheme [2] is generated as a polynomial  $\mathbf{s} \leftarrow \chi_{key}$ , and the corresponding public-key is computed as:

$$\mathbf{b} = ([-(\mathbf{a} \cdot \mathbf{s} + \mathbf{e})]_q, \mathbf{a}) \in \mathcal{R}_q^2$$

with  $\mathbf{a} \leftarrow \mathcal{R}_q$  drawn uniformly at random and  $\mathbf{e} \leftarrow \chi_{err}$ .  $\chi_{key}$  is such that for  $\mathbf{s} \leftarrow \chi_{key}$ ,  $\|\mathbf{s}\|_\infty \leq B_{key}$  for a small  $B_{key}$ , while the distribution  $\chi_{err}$  has standard deviation  $\sigma_{err}$  such that elements sampled from this distribution have their infinite norm smaller than  $B_{err} = 6\sigma_{err}$  with very high probability.

A relinearisation key is defined with respect to a pair of functions  $\mathfrak{D}$  and  $\mathfrak{P}$ , such that  $\langle [\mathfrak{D}(\mathbf{a}), \mathfrak{P}(\mathbf{b})] \rangle_q = [\mathbf{a} \cdot \mathbf{b}]_q$  and the infinite norm of  $\mathfrak{D}(\mathbf{a})$  is “small”. The relinearisation key, which can be interpreted as a vector of pseudo-encryptions of  $\mathfrak{P}(\mathbf{s}^2)$ , is defined as:

$$\overrightarrow{\text{rlk}} = (\langle [\mathfrak{P}(\mathbf{s}^2) - (\overrightarrow{\mathbf{e}} + \overrightarrow{\mathbf{d}} \cdot \mathbf{s})] \rangle_q, \overrightarrow{\mathbf{d}}) \in \mathcal{R}_q^{2l}$$

where  $\overrightarrow{\mathbf{e}} \leftarrow \chi_{err}^l$  and  $\overrightarrow{\mathbf{d}} \leftarrow \mathcal{R}_q^l$ . A ciphertext  $\text{ct} = (c_0, \mathbf{c}_1) \in \mathcal{R}_q^2$ , encrypting  $[\mathbf{m}]_t \in \mathcal{R}_t$ , satisfies  $c_0 + \mathbf{c}_1 \cdot \mathbf{s} = \Delta[\mathbf{m}]_t + \mathbf{v} \bmod q$ , where  $\Delta = \lfloor q/t \rfloor$ .  $\text{ct}$  can be correctly decrypted as long as the noise  $\mathbf{v}$  satisfies  $\|\mathbf{v}\|_\infty < \Delta/2 - |q|_t/2$ , through:

$$\mathbf{m} = \left\lfloor \left\lfloor \frac{t}{q} [c_0 + \mathbf{c}_1 \cdot \mathbf{s}]_q \right\rfloor \right\rfloor_t. \quad (1)$$

The homomorphic addition of two ciphertexts corresponds to the pairwise addition of their coefficients modulo  $q$ . The homomorphic multiplication of  $\text{ct}$  and  $\text{ct}'$  encrypting  $[\mathbf{m}]_t$  and  $[\mathbf{m}']_t$  is performed in two steps. First, we compute the three elements ciphertext:

$$\widehat{\text{ct}}_{\text{mult}} = \left( \left\lfloor \left\lfloor \frac{t}{q} c_0 \cdot c'_0 \right\rfloor \right\rfloor_q, \left\lfloor \left\lfloor \frac{t}{q} (c_0 \cdot c'_1 + c_1 \cdot c'_0) \right\rfloor \right\rfloor_q, \left\lfloor \left\lfloor \frac{t}{q} c_1 \cdot c'_1 \right\rfloor \right\rfloor_q \right) \in \mathcal{R}_q^3 \quad (2)$$

satisfying  $\widehat{c}_0 + \widehat{c}_1 \mathbf{s} + \widehat{c}_2 \mathbf{s}^2 = \Delta[\mathbf{m} \cdot \mathbf{m}']_t + \widehat{\mathbf{v}} \bmod q$ . In a second step,  $\widehat{\text{ct}}_{\text{mult}}$  is relinearised, by multiplying  $\mathfrak{D}(\widehat{c}_2)$  by  $\overrightarrow{\text{rlk}}$  and adding the result to  $(\widehat{c}_0, \widehat{c}_1)$ :

$$\text{ct}_{\text{mult}} = \left( \left[ \widehat{c}_0 + \langle \mathfrak{D}(\widehat{c}_2), \overrightarrow{\text{rlk}}_0 \rangle \right]_q, \left[ \widehat{c}_1 + \langle \mathfrak{D}(\widehat{c}_2), \overrightarrow{\text{rlk}}_1 \rangle \right]_q \right).$$

The result  $\text{ct}_{\text{mult}}$ , can be decrypted with (1) to produce  $[\mathbf{m} \cdot \mathbf{m}']_t$ . As homomorphic operations are applied, the noise associated with ciphertexts grows, limiting the number of operations that can be applied.

### B. Hybrid Position-Residue Number System

Using [6] and considering a modulus  $q$  such that  $q = p^d$ , one can represent polynomials  $\mathbf{a} \in \mathcal{R}_q$  in a positional system of large radix base  $p$  – i.e. such that  $\mathbf{a} = \sum_{i=0}^{d-1} \mathbf{a}_i p^i$  with  $(\mathbf{a}_0, \dots, \mathbf{a}_{d-1}) \in (\mathcal{R}_p)^d$ . We view elements represented in HPR as elements of  $\mathcal{R}_p[Y]$ , leading to:

$$\mathbf{a}(Y) = \sum_{i=0}^{d-1} \mathbf{a}_i Y^i$$

which produces the value of  $\mathbf{a}$  when evaluated at  $Y = p$ .

Moreover, if the modulus  $p$  is chosen as a product of small distinct primes,  $p = \prod_{p_i \in \mathcal{P}} p_i$  with  $\mathcal{P} = \{p_1, \dots, p_k\}$ , the Chinese Remainder Theorem (CRT) states that there is a ring isomorphism between  $\mathcal{R}_p$  and  $\mathcal{R}_{p_1} \times \dots \times \mathcal{R}_{p_k}$ . In this case, the digits  $\mathbf{a}_i$ s can be represented in RNS as  $f_{\mathcal{P}}(\mathbf{a}_i) = ([\mathbf{a}_i]_{p_1}, \dots, [\mathbf{a}_i]_{p_k})$  which allows to break down the arithmetic over the large coefficients into several smaller channels and thus enhance performance. The RNS representation can be reversed, and  $\mathbf{a}_i$  recovered, by computing:

$$f_{\mathcal{P}}^{-1}([\mathbf{a}_i]_{p_1}, \dots, [\mathbf{a}_i]_{p_k}) = \left[ \sum_{i=1}^k \left[ [\mathbf{a}_i]_{p_i} \frac{p_i}{p} \right]_{p_i} \frac{p}{p_i} \right]_p = \sum_{i=1}^k \left[ \mathbf{a}_i \frac{p_i}{p} \right]_{p_i} \frac{p}{p_i} - \alpha_i p \quad (3)$$

with  $\|\alpha_i\|_\infty \leq (k-1)/2$ . However, since the HPR is of a positional nature, one has to consider the propagation of carries over the digits after additions or multiplications. In order to compute the carries, the digits should be able to grow larger than  $p$ . Thus, the digits arithmetic cannot be performed directly modulo  $p$  and a second RNS basis  $\mathcal{B} = \{b_1, \dots, b_k\}$  is adjoined to  $\mathcal{P}$ .

1) *Carry Propagation*: In order to approximate the carry caused by the  $i^{\text{th}}$  digit,  $\nu_i = \lfloor \frac{\mathbf{a}_i}{p} \rfloor$ , we add a modulus  $b_{sk}$  to  $\mathcal{B}$  to form  $\mathcal{B}_{sk}$ , and use a Shenoy-Kumaresan approach [8]. First, we use a fast base extension (FBE),

$$\text{FBE}(\mathbf{a}_i, \mathcal{P}, \mathcal{B}_{sk}) = \left( \sum_{i=1}^k \left[ \mathbf{a}_i \frac{p_i}{p} \right]_{p_i} \cdot \frac{p}{p_i} \bmod b \right)_{b \in \mathcal{B}_{sk}}$$

to obtain  $[\mathbf{a}_i]_p + \alpha_i \cdot p$  in base  $\mathcal{B}_{sk}$ . Then, the approximated carry  $\tilde{\nu}_i = (\mathbf{a}_i - \text{FBE}(\mathbf{a}_i, \mathcal{P}, \mathcal{B}_{sk}))p^{-1} = \lfloor \frac{\mathbf{a}_i}{p} \rfloor - \alpha_i$  is produced in  $\mathcal{B}_{sk}$ . To finalise the procedure, the value of  $\tilde{\nu}_i$  is extended exactly to  $\mathcal{P}$ . An inexact extension first computes  $\text{FBE}(\tilde{\nu}_i, \mathcal{B}, \mathcal{P}) = \tilde{\nu}_i + \alpha'_i \cdot b$  in basis  $\mathcal{P}$  with  $b = b_1 \cdot \dots \cdot b_k$  and  $\|\alpha'_i\|_\infty \leq (k-1)/2$ . Then, the extra residue modulo  $b_{sk}$  is used to correct  $\alpha'_i \cdot b$ , by computing  $\alpha'_i$  with:

$$[\alpha'_i]_{b_{sk}} = \left[ \frac{1}{b} \left( \sum_{i=1}^{\ell} \left[ \mathbf{a}_i \frac{b_i}{b} \right]_{b_i} \frac{b}{b_i} - [\tilde{c}_i]_{b_{sk}} \right) \right]_{b_{sk}}.$$

With this procedure, values  $\|\tilde{\nu}_i\| < \lambda B$ , for some  $\lambda > 0$ , can be exactly represented and extended from  $\mathcal{B}$  so long as  $b_{sk} \geq 2(k+\lambda)$  [4, Lemma 6]. For the parameters considered in this article,  $b_{sk}$  can have the same bitwidth as the other

moduli in  $\mathcal{P}$  and  $\mathcal{B}$ , which is enough to get the residues of  $[\mathbf{a}_i]_p + \alpha_i \cdot p$  and of  $\tilde{\nu}_i = \nu_i - \alpha_i$  in  $\mathcal{P} \cup \mathcal{B}_{\text{sk}}$ .

Afterwards,  $\tilde{\nu}_i$  can be added to the digit of index  $i+1$  and one can re-run the procedure to compute the carry  $\tilde{\nu}_{i+1}$  caused by the digit of index  $i+1$ . If one assumes an element to belong to  $\mathcal{R}_q$ , the representation of the Most Significant Digit (MSD) in basis  $\mathcal{B}$  can be discarded. The MSD is implicitly multiplied by  $p^{d-1}$ , so any value larger than  $p/2$  in absolute value would lead to representations larger than  $q/2 = p \times p^{d-1}/2$ , which are redundant.

**Computational cost.** The computation of a carry for a single HPR digit requires a first fast base extension from  $\mathcal{P}$  to the  $k+1$  moduli of  $\mathcal{B}_{\text{sk}}$  ( $n(k(k+2))$  EMMs), followed by the product by  $1/p$  in base  $\mathcal{B}_{\text{sk}}$  ( $n(k+1)$  EMMs), and a Shenoy-Kumaresan extension ( $n(k^2+3k+1)$  EMMs). Finally, since, when operating modulo  $q$ , carries are consecutively computed for all the digits but the last, the carry propagation procedure applied to an element in  $\mathcal{R}_p^d$  requires the following amount of EMMs:

$$\begin{cases} \mathcal{C}_{1\_carry} &= n(2k^2 + 6k + 2) \text{ EMMs} \\ \mathcal{C}_{\text{all\_carries}} &= (d-1) \cdot \mathcal{C}_{1\_carry} \end{cases}.$$

### III. PROPOSED HPR-BASED CONSTRUCTION

In this section, we propose an efficient HPR variant of the FV scheme. Key generation, encryption and homomorphic addition take place as in the original scheme [2], but with polynomial coefficients represented in HPR.

#### A. Proposed HPR-based Decryption Algorithm

In this section, an HPR-based method for FV decryption is proposed. We assume that the HPR representation of a ciphertext  $(\mathbf{c}_0, \mathbf{c}_1) \in \mathcal{R}_q^2$ , with:

$$\mathbf{c}_i(Y) = \sum_{j=0}^{d-1} \mathbf{c}_{i,j} Y^j$$

having digits satisfying  $\|\mathbf{c}_{i,j}\| \leq \beta \cdot p$  for  $i = 0, 1$  and  $0 \leq j \leq d-2$ :  $\beta$  may take the value of  $k/2$  if carry propagation has been applied; or might be even slightly larger for cases where the ciphertext represents an homomorphic sum. Lemma 1 proves that, for the considered setting, the value of  $\mathbf{c}_0(p) + \mathbf{c}_1(p) \cdot \mathbf{s}$  can be efficiently approximated having only access to the mod- $\mathcal{P}$  representation of the MSD of  $\mathbf{c}_0$  and  $\mathbf{c}_1$  without causing a large error. This reduces by a factor of  $d$  the number of NTTs one needs to compute.

**Lemma 1.** For  $i = 0, 1$ , let  $\mathbf{c}_i$  be such that  $\|\mathbf{c}_{i,j}\|_\infty \leq \beta p$  for  $j \in \llbracket 0, d-2 \rrbracket$ , then:

$$[\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s}]_p = \frac{\Delta[\mathbf{m}]_t + \mathbf{v} + \mathbf{q}\mathbf{u}}{p^{d-1}} + \mathbf{e} + \epsilon \cdot p \quad (4)$$

for some  $\epsilon \in \mathcal{R}$ , with the error  $\mathbf{e}$  verifying  $\|\mathbf{e}\|_\infty \leq \frac{p}{p-1} \beta (1 + \delta_{\mathcal{R}} B_{\text{key}})$ .

*Proof.* In Appendix A.  $\square$

Decryption consists of efficiently computing (5) modulo  $t$  based on the values of  $\mathbf{c}_{0,d-1}$  and  $\mathbf{c}_{1,d-1}$ .

$$\left\lceil t \cdot \frac{[\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s}]_p}{p} \right\rceil = [\mathbf{m}]_t + \lceil \tilde{\mathbf{e}} \rceil + t(\mathbf{u} + \epsilon) \quad (5)$$

with  $\tilde{\mathbf{e}} = (t \cdot \mathbf{v} - |q|_t [\mathbf{m}]_t + t p^{d-1} \epsilon) / q$  which has a norm bounded up by:

$$\|\tilde{\mathbf{e}}\|_\infty \leq \frac{t}{q} \|\mathbf{v}\|_\infty + \frac{|q|_t}{2q} + \frac{t}{p-1} \beta (1 + \delta_{\mathcal{R}} B_{\text{key}}) \quad (6)$$

We use a strategy similar to [4] to compute (5) with a fast basis extension but without any error. By scaling the computation in (5) by an integer  $\gamma$ , any extension errors will be detected since they will be nonzero modulo  $\gamma$ . After efficiently getting  $[\gamma t(\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s})]_p + \alpha p$  through a fast conversion from  $\mathcal{P}$  to  $\{\gamma t\}$ , one can compute:

$$\begin{aligned} \mathbf{r} &= (\gamma t[\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s}]_p - ([\gamma t(\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s})]_p + \alpha p)) / p \\ &= (\gamma t[\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s}]_p - \gamma t([\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s}]_p) / p - \alpha \\ &\quad + (\gamma t([\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s}]_p) - \gamma t([\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s}]_p) / p) / p \\ &= \gamma \left\lceil t \frac{[\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s}]_p}{p} \right\rceil - \alpha + \left\lceil \gamma \frac{t([\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s}]_p)}{p} \right\rceil \\ &= \gamma([\mathbf{m}]_t + \lceil \tilde{\mathbf{e}} \rceil) + \left\lceil \gamma \frac{t([\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s}]_p)}{p} \right\rceil - \alpha \\ &= \gamma[\mathbf{m}]_t + \gamma \lceil \tilde{\mathbf{e}} \rceil + \left\lceil \gamma \frac{[p\tilde{\mathbf{e}}]_p}{p} \right\rceil - \alpha \text{ mod } \gamma t. \end{aligned} \quad (7)$$

Now, through the application of Lemma 2 to (7),  $[\mathbf{m}]_t$  can be computed as  $(\mathbf{r} - [\mathbf{r}]_\gamma) / \gamma \text{ mod } t$ .

**Lemma 2.** Let  $\gamma > k/2$  and  $\|\alpha\|_\infty \leq (k-1)/2$ . If we have

$$\|\mathbf{v}\|_\infty < \frac{q}{2t} \left( 1 - \frac{k}{\gamma} - \frac{2t\beta}{p-1} (1 + \delta_{\mathcal{R}} B_{\text{key}}) \right) - \frac{|q|_t}{2} \quad (8)$$

then  $\lceil \tilde{\mathbf{e}} \rceil = 0$ , and  $\left\lceil \left\lceil \gamma \frac{[p\tilde{\mathbf{e}}]_p}{p} \right\rceil - \alpha \right\rceil_\gamma = \left\lceil \gamma \frac{[p\tilde{\mathbf{e}}]_p}{p} \right\rceil - \alpha$ . Hence, the correction technique from [4] works.

*Proof.* In Appendix B.  $\square$

#### B. Summarising

Here we summarise the above-described decryption steps. First,  $\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s}$  is computed as an approximation to  $(\mathbf{c}_0(p) + \mathbf{c}_1(p) \cdot \mathbf{s}) / p^{d-1}$  with  $k_{\text{NTT}} + n_{\text{EMM}} + k_{\text{NTT}}^{-1}$ . Then (7) is computed modulo  $\gamma t$ . Notice that since  $\gamma t \cdot [\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s}]_p = 0 \pmod{\gamma t}$ , the expression for  $\mathbf{r}$  can be simplified as:

$$\begin{aligned} \mathbf{r} &= \frac{-([\gamma t(\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s})]_p + \alpha p)}{p} \text{ mod } \gamma t \\ &= -\text{FBE}([\gamma t(\mathbf{c}_{0,d-1} + \mathbf{c}_{1,d-1} \cdot \mathbf{s})], \mathcal{P}, \gamma t) \cdot p^{-1} \text{ mod } \gamma t \end{aligned}$$

which costs  $kn_{\text{EMM}} + kn_{\text{EMM}} \gamma t$  to compute if  $\gamma t$  fits one machine word. Then the message is extracted by using the centred remainder of  $\mathbf{r}$  modulo  $\gamma$  since  $[\mathbf{m}]_t = (\mathbf{r} - [\mathbf{r}]_\gamma) / \gamma \text{ mod } t$  under the conditions of Lemma 2.

Notice that as in [4], if  $\gamma$  is chosen as a power of two, division by  $\gamma$  can be replaced with a shift. Hence, the decryption algorithm has a computational cost of:

$$k \text{ NTT} + 2kn \text{ EMM} + k \text{ NTT}^{-1} + kn \text{ EMM}_{\gamma t}.$$

### C. Proposed HPR-based Homomorphic Multiplication

In this section, an HPR-based method for FV homomorphic multiplication is proposed. We consider two ciphertexts  $\text{ct} = (\mathbf{c}_0(Y), \mathbf{c}_1(Y))$  and  $\text{ct}' = (\mathbf{c}'_0(Y), \mathbf{c}'_1(Y))$ , with their elements represented in HPR, encrypting respectively  $\mathbf{m}$  and  $\mathbf{m}'$ , such that  $\mathbf{c}_0 + \mathbf{c}_1 \cdot \mathbf{s} = \Delta[\mathbf{m}]_t + \mathbf{v} + \mathbf{r}q$  and  $\mathbf{c}'_0 + \mathbf{c}'_1 \cdot \mathbf{s} = \Delta[\mathbf{m}']_t + \mathbf{v}' + \mathbf{r}'q$  where  $\mathbf{r}$  and  $\mathbf{r}'$  satisfy:

$$\|\mathbf{r}\|_\infty, \|\mathbf{r}'\|_\infty \leq \underbrace{\frac{\delta_{\mathcal{R}} B_{key} + 1}{2}}_{r_\infty} + 1. \quad (9)$$

The polynomial digits are assumed to be bounded in a similar way to how they were bounded during decryption in Section III-A, i.e. for  $(i, j) \in \{0, 1\} \times \{0, \dots, d-1\}$ ,  $\|\mathbf{c}_{i,j}\|_\infty, \|\mathbf{c}'_{i,j}\|_\infty \leq \beta \cdot p$  with  $\beta \simeq k/2$ .

Until now, we only needed to represent the digit of index  $d-1$  in basis  $\mathcal{P}$ , since we were dealing with elements of  $\mathcal{R}$  modulo  $q$ . However, during homomorphic multiplication, the first product in (2) is not reduced modulo  $q$ . Therefore we have to extend the digits of index  $d-1$  to the basis  $\mathcal{P} \times \mathcal{B}_{sk}$ . This extension is done exactly with the approach of [5], denoted HPS in the rest of the paper. First,  $\text{FBE}(\mathbf{c}_{i,d-1}, \mathcal{P}, \mathcal{B}_{sk}) = [\mathbf{c}_{i,d-1}]_p + \boldsymbol{\alpha}_i \cdot p$  is computed. Then,  $\boldsymbol{\alpha}_i$  is produced as:

$$\boldsymbol{\alpha}_i = \left\lfloor \sum_{j=1}^k \frac{[\mathbf{c}_{i,d-1} \frac{p_j}{p}]_{p_j}}{p_j} \right\rfloor \quad (10)$$

by using floating point operations to compute the quotients and by rounding to the nearest integer. Finally, the term  $\boldsymbol{\alpha}_i \cdot p$  is subtracted from the intermediate result. The lifting procedure of one digit has a practical complexity of  $n(k^2 + 3k + 1)$  EMMs and  $n(k + 1)$  Floating-Point Operations (FPOs).

Note that the approximation of (10) with floating point arithmetic may lead to an erroneous computation of  $\boldsymbol{\alpha}_i$ . There is a region  $\mathbb{Z} + \frac{1}{2} \pm \epsilon$ , where  $\epsilon$  is related to the used precision and the number of moduli in the basis  $\mathcal{P}$  ( $|\epsilon| < 2^{-50}$  for IEEE 754 double precision and  $k < 8$ ), to which the value computed in (10) may belong before rounding, causing a possible error of at most 1. While one can detect when this happens, and redo the computation with a higher precision to reduce  $\epsilon$ , we have chosen, as in [5], to skip this error detection. We will discuss the possible effects of this choice in Remark 1.

At this point we aim at computing (2). To do so we start by computing the NTTs of all the digits of the two input ciphertexts ( $4d(2k + 1)$  NTTs in total). Then the product  $(\tilde{\mathbf{c}}_0(Y), \tilde{\mathbf{c}}_1(Y), \tilde{\mathbf{c}}_2(Y))$  is computed with a Karatsuba pattern.

The products of polynomials in  $Y$  are performed as usual, i.e. in the power basis, and the products of digits, which are in NTT form, are made component-wise. Since the product is not reduced modulo  $q$ , no reduction modulo  $Y^d$  is done. Thus we obtain polynomials of degree  $(2d - 2)$  in  $Y$ . The division by  $q = p^d$  and the rounding can be approximated by computing only the  $d - 1$  MSDs of the products. Nonetheless, in order to reduce the noise added by this approximation, we also keep the carry coming from the digits of index  $d - 1$ .

**Lemma 3.** Let  $i \in \{0, 1, 2\}$  and  $[\tilde{\mathbf{c}}_{i,d-1}/p] - \boldsymbol{\alpha}_{i,d-1}$  be the carry coming from  $\tilde{\mathbf{c}}_{i,d-1}$  with  $\|\boldsymbol{\alpha}_{i,d-1}\|_\infty \leq (k - 1)/2$  (see Section II-B1). We have:

$$\left\lfloor \frac{\tilde{\mathbf{c}}_{i,d-1}}{p} \right\rfloor - \boldsymbol{\alpha}_{i,d-1} + \tilde{\mathbf{c}}_{i,d} + \dots + \tilde{\mathbf{c}}_{i,2d-2} \cdot p^{d-2} = \frac{\tilde{\mathbf{c}}_i}{q} + \mathbf{e}_i \quad (11)$$

$$\text{with: } \|\mathbf{e}_i\|_\infty \leq \underbrace{6\delta_{\mathcal{R}} t \beta^2 p \left( \frac{d-1}{p-1} - \frac{1-p^{-d}}{(p-1)^2} \right)}_{B_e} + \frac{k}{2}. \quad (12)$$

*Proof.* In Appendix C.  $\square$

The three polynomials of degree  $d - 1$  in  $Y$  described in Lemma 3 approximate (2) with an error  $\mathbf{e}_i$ . They may be interpreted as encrypting  $[\mathbf{m} \cdot \mathbf{m}']_t$ . The next proposition states the inherent noise of this new ciphertext, depending on the input noises  $\mathbf{v}$  and  $\mathbf{v}'$ .

**Proposition 1.** For  $i \in \{0, 1, 2\}$ , let  $\hat{\mathbf{c}}_i = \left\lfloor \frac{\tilde{\mathbf{c}}_{i,d-1}}{p} \right\rfloor - \boldsymbol{\alpha}_{i,d-1} + \sum_{j=d}^{2d-2} \tilde{\mathbf{c}}_{i,j} \cdot p^{j-d}$  as in (11), then we have:

$$\hat{\mathbf{c}}_0 + \hat{\mathbf{c}}_1 \cdot \mathbf{s} + \hat{\mathbf{c}}_2 \cdot \mathbf{s}^2 \equiv \Delta \cdot [\mathbf{m} \cdot \mathbf{m}']_t + \hat{\mathbf{v}} \pmod{q}$$

$$\begin{aligned} \text{with: } \|\hat{\mathbf{v}}\|_\infty &\leq \delta_{\mathcal{R}} t (r_\infty + 1/2) (\|\mathbf{v}\|_\infty + \|\mathbf{v}'\|_\infty) \\ &\quad + \delta_{\mathcal{R}} \min(\|\mathbf{v}\|_\infty, \|\mathbf{v}'\|_\infty) / 2 + t^2 \delta_{\mathcal{R}} (1 + r_\infty) \\ &\quad + B_e (1 + \delta_{\mathcal{R}} B_{key} + \delta_{\mathcal{R}}^2 B_{key}^2) + 1/2(t + 1) \end{aligned} \quad (13)$$

and where  $B_e$  is as in (12) and  $r_\infty$  is as in (9).

*Proof.* In Appendix D.  $\square$

**Remark 1.** In the case where an error has occurred in the HPS basis extension, before the multiplication, then we would operate on  $\mathbf{c}_i + q\mathbf{u}$  with  $\|\mathbf{u}\|_\infty = 1$  (or  $\mathbf{c}'_i + q\mathbf{u}'$  with  $\|\mathbf{u}'\|_\infty = 1$ ). In this case  $r_\infty$  would increase to  $3(\delta_{\mathcal{R}} B_{key} + 1)/2 + 1$ . As explained in [5, Section 4.5], this would only impact the size of  $\delta_{\mathcal{R}} t (r_\infty + 1/2) (\|\mathbf{v}\|_\infty + \|\mathbf{v}'\|_\infty)$  by less than 3% in the worst case. Considering that this error occurs with probability smaller than  $4n \cdot 2^{-50} \leq 2^{-32}$  on each multiplication (for  $n \leq 2^{16}$ ), its effect on the noise growth in practice is negligible.

In practice, after having computed  $(\tilde{\mathbf{c}}_0(Y), \tilde{\mathbf{c}}_1(Y), \tilde{\mathbf{c}}_2(Y))$ , the carries of  $\tilde{\mathbf{c}}_{0,d-1}$  and  $\tilde{\mathbf{c}}_{1,d-1}$  are computed and stored apart. These are then added after the relinearisation step to the final output to reduce the overall computational complexity. The computation of the two carries requires  $2(2k + 1)$  NTT $^{-1}$  to leave the NTT representation, and  $2 \times \mathcal{C}_{1\_carry}$ . Moreover, the carry generated by  $\tilde{\mathbf{c}}_{2,d-1}$  is only computed and propagated during the relinearisation step. Therefore, at this point the three elements ciphertext has the  $d - 1$  MSDs of  $\tilde{\mathbf{c}}_0$  and  $\tilde{\mathbf{c}}_1$  and the  $d$  MSDs of  $\tilde{\mathbf{c}}_2$  in NTT representation; and the two carries stored apart in coefficient representation.

**Computational cost.** Overall, assuming that each vector-matrix product is performed with a naive algorithm, the homomorphic multiplication (before relinearisation) requires:

$$\begin{aligned} C_{\text{mult.}} &= n(4k^2 + 12k + 4) \text{EMMs} + n(4k + 4) \text{FPOs} \\ &\quad + 4d(2k + 1) \text{NTTs} + 2(2k + 1) \text{NTTs}^{-1} \\ &\quad + 3nd(2k + 1)(d + 3)/2 \text{EMMs} + 2 \cdot \mathcal{C}_{1\_carry} \end{aligned}$$

#### D. Relinearisation

At the beginning of relinearisation, the  $d$  MSDs of  $\tilde{c}_2$  are converted from the NTT to the canonical domain at a cost of  $d(2k+1)$  NTT<sup>-1</sup>. Then, carry propagation is applied and the Least Significant Digit (LSD) is discarded, producing  $\hat{c}_2$ , as per Proposition 1. By construction, the digits  $\mathbf{a} = ([\mathbf{a}_{0,\mathcal{P}}, \mathbf{a}_{0,\mathcal{B}}], \dots, [\mathbf{a}_{d-1,\mathcal{P}}, \mathbf{a}_{d-1,\mathcal{B}}])_{\text{HPR}}$  that result from propagating carries satisfy:

$$[Y^i f_{\mathcal{P} \cup \mathcal{B}}^{-1}([\mathbf{a}_{i,\mathcal{P}}, \mathbf{a}_{i,\mathcal{B}}])]_q = [Y^i \text{FBE}(\mathbf{a}_{i,\mathcal{P}}, \mathcal{P}, \{q\})]_q \quad (14)$$

where  $f_{\mathcal{P} \cup \mathcal{B}}^{-1}$  corresponds to reverting the RNS representation with respect to the  $(\mathcal{P}, \mathcal{B})$  basis (see (3)). Concretely, this means that once the carry propagation has been done, the residues of  $\hat{c}_2$  in base  $\mathcal{B}$  are no longer required for the relinearisation process.  $[\hat{c}_2 \cdot s^2]_q$  can be rewritten as:

$$[\hat{c}_2 \cdot s^2]_q = \left[ \sum_{i=0}^{d-1} \text{FBE}(\hat{c}_{2,i}, \mathcal{P}, \{q\}) \cdot s^2 Y^i \right]_q = \left[ \sum_{i=0}^{d-1} \sum_{j=1}^k \left[ \hat{c}_{2,i} \frac{p_j}{p} \right]_{p_j} \frac{p}{p_j} Y^i s^2 \right]_q. \quad (15)$$

Thus, we can define  $\mathcal{D}(\hat{c}_2)$  and  $\mathfrak{P}(s^2)$  as the vectors with entries  $\mathbf{r}_{i,j} = \left[ \hat{c}_{2,i} \frac{p_j}{p} \right]_{p_j}$  and  $\psi_{i,j} = \left[ \frac{p}{p_j} Y^i s^2 \right]_q$ , respectively, for  $0 \leq i \leq d-1$  and  $1 \leq j \leq k$ . Clearly, the  $\|\mathbf{r}_{i,j}\|_\infty < p_j/2$  are small, and  $[\langle \mathcal{D}(\hat{c}_2), \mathfrak{P}(s^2) \rangle]_q = [\hat{c}_2 \cdot s^2]_q$ .

The relinearisation-key ( $\overline{\text{rlk}}$ ) is composed of pseudo-encryptions of the  $\psi_{i,j}$  values. To produce the pseudo-encryption with respect to  $\psi_{i,j}$ , a Ring Learning with Errors (RLWE) sample  $([-(\mathbf{a}_{i,j} \cdot \mathbf{s} + e_{i,j})]_q)$ ,  $\mathbf{a}_{i,j}$  is drawn, and  $\frac{p}{p_j} \cdot s^2$  is added to the  $i$ -th digit of the first element. Then, the whole relinearisation key is stored in NTT form.

**Remark 2.** Note that, since  $\|s^2\|_\infty \leq \delta_{\mathcal{R}} B_{\text{key}}^2 < p_j$  for moduli  $p_j$  of practical size,  $\frac{p}{p_j} \cdot s^2$  is representable as a single HPR digit. Thus, the RLWE samples used to mask the  $\psi_{i,j} = \frac{p}{p_j} \cdot Y^i \cdot s^2$ , for a fixed  $i$  and  $1 \leq j \leq k$ , do not need to have a degree higher than  $i$  in  $Y$ . The  $d-1-i$  MSDs of the  $k$  elements of the relinearisation-key associated with these  $\psi_{i,j}$  can be set to zero. In practice, RLWE samples are drawn with respect to a modulus  $p^i \leq p^d$ , but with an error of same standard deviation  $\sigma_{\text{err}}$ . This technique allows us to not only reduce the size of  $\overline{\text{rlk}}$ , but also to further reduce the practical complexity of the relinearisation.

When a small polynomial  $\mathbf{r}_{i,j}$  has been extracted from a residue modulo  $p_j$ , it is easily extended to the base  $(\mathcal{P}, \mathcal{B}_{\text{sk}})$  through a simple copy-paste operation. Then, since its representation is a single digit, the computation of its NTT form only requires  $(2k+1)$  NTTs. In the next step, the polynomials  $\mathbf{r}_{i,j}$  are multiplied by the corresponding  $\overline{\text{rlk}}$  elements and the results are accumulated in  $\hat{c}_0, \hat{c}_1$ .

The relinearisation step causes an increase in the noise associated with the ciphertext. More concretely, for a homomorphic multiplication result  $(\hat{c}'_0, \hat{c}'_1)$ , one gets:

$$\begin{aligned} [\hat{c}'_0 + \hat{c}'_1 \mathbf{s}]_q &= [\hat{c}_0 + \hat{c}_1 \cdot \mathbf{s} + \langle \mathcal{D}(\hat{c}_2), \mathfrak{P}(s^2) - (\vec{e} + \vec{a} \cdot \mathbf{s}) \rangle + \langle \mathcal{D}(\hat{c}_2), \vec{a} \rangle \cdot \mathbf{s}]_q \\ &= [\hat{c}_0 + \hat{c}_1 \cdot \mathbf{s} + \hat{c}_2 \cdot s^2 - \langle \mathcal{D}(\hat{c}_2), \vec{e} \rangle]_q \\ &= [\hat{c}_0 + \hat{c}_1 \cdot \mathbf{s} + \hat{c}_2 \cdot s^2 + e_{\text{relin}}]_q. \end{aligned}$$

The extra noise caused by this step is bounded by:

$$\|e_{\text{relin}}\|_\infty \leq \delta_{\mathcal{R}} \cdot d \cdot B_{\text{err}} \cdot (p_1 + \dots + p_k). \quad (16)$$

**Computational cost.** After the carry propagation of the  $d$  MSDs of  $\tilde{c}_2$ , we have to compute the NTTs of 1-digit polynomials  $\mathbf{r}_{i,j}$  in base  $(\mathcal{P}, \mathcal{B}_{\text{sk}})$  for a total of  $dk(2k+1)$  NTTs. Then, for each  $j$ , the single-digit polynomials  $\mathbf{r}_{i,j}$  are multiplied by two relinearisation-key elements of degree  $i$ , resulting in a total of  $nk d(d+1)(2k+1)$  elementary products. The results of these products are added to  $\hat{c}_0$  and  $\hat{c}_1$ . Once all intermediate products have been accumulated,  $2d(2k+1)$  NTT<sup>-1</sup>s are used to obtain the result in the canonical domain. Finally, carry propagation is applied to  $\hat{c}'_0$  and  $\hat{c}'_1$ . Thus, the total cost of the relinearisation step is:

$$\begin{aligned} C_{\text{relin}} &= dk(2k+1) \text{ NTTs} + 3d(2k+1) \text{ NTTs}^{-1} \\ &\quad + nk d((d+1)(2k+1) + 1) \text{ EMMs} \\ &\quad + (3d-2)C_{1\_carry} \end{aligned}$$

**Remark 3.** The bit size of the noise grows linearly with the multiplicative depth  $L$ . Consequently, when the depth  $L$  reaches around  $L_{\text{max}}/d$ , the noise affects more than one digit. At this point one can apply a coarse grain decomposition to decrease the practical complexity of the relinearisation. The sum with index  $j$  in (15) can be computed by the decomposition function prior to the product with the relinearisation key. With this technique, the decomposition polynomials become  $\mathbf{r}_i = \sum_{j=1}^k \left[ \hat{c}_{2,i} \frac{p_j}{p} \right]_{p_j} \frac{p}{p_j}$ , and the new relinearisation key contains  $\psi_i = [Y^i s^2]_q$ . Since  $\mathbf{r}_i$  is no longer a single residue, its conversion toward base  $\mathcal{B}$  is not a simple copy-paste. However, the carry propagation performed prior to the relinearisation already involves such a base conversion. So,  $\mathbf{r}_i$  is obtained in the full RNS base at no extra cost. Moreover, one may disregard an increasingly larger number of LSDs during relinearisation as noise grows. Assuming that only the  $d'$  MSDs are taken into account for  $\hat{c}_{2,i}$  (in the case where  $L \simeq (d-d') \frac{L_{\text{max}}}{d}$ ), the cost boils down to:

$$\begin{aligned} C_{\text{relin},d'} &= d'(2k+1) \text{ NTTs} + 3d'(2k+1) \text{ NTTs}^{-1} \\ &\quad + nd'(d'+1)(2k+1) \text{ EMMs} \\ &\quad + (3d'-2)C_{1\_carry} \end{aligned}$$

#### IV. RELATED ART

A first adaptation of FV to the RNS was achieved in [4]. Due to the reduced complexity, speedups from 2 up to 4 and from 5 up to 20 were achieved for the homomorphic multiplication and decryption operations, respectively, in comparison to an implementation based on a generic multiprecision library. Later, [5] proposed modifications to [4], namely by computing part of the RNS basis extensions with floating point arithmetic.

It has been noticed in [9] that for similar cryptographic parameters, [4] offers a lower multiplicative depth than [5]. However as pointed out in [10], this observation results from a misimplementation of the techniques of [4] in the SEAL and PALISADE libraries. Thus, since there is no clear advantage of using one version over the other, herein, we make use of techniques from both [4] and [5].

## V. COMPLEXITY ANALYSIS

The complexities of our HPR-based variant of the FV scheme and of [4], [5] are presented in Table I which, for simplicity, associates both forward and inverse NTTs with the NTTs column.  $K_{RNS}$  corresponds to the size of the RNS bases of [4], [5]. Both [4], [5] make use of secondary bases of size  $K_{RNS} + 1$ . In HPR, each polynomial coefficient is associated with  $d$  digits, each of them represented in an RNS basis of size  $k$ , and a secondary RNS basis with  $k + 1$  moduli. If both systems target the same security level and use moduli of similar bitwidth, the value of  $k$  will be about  $d$  times smaller than  $K_{RNS}$ , since  $q_1 \dots q_{K_{RNS}} \approx (p_1 \dots p_k)^d$ .

### A. Decryption

In all schemes, the computational cost is dominated by the amount of NTTs. A radix-2 NTT requires  $\frac{n}{2} \log_2 n$  EMMs [11]. In practice, [4], [5] require  $K_{RNS}$  forward NTTs and  $K_{RNS}$  inverse NTTs; whereas  $k$  forward NTTs and  $k$  inverse NTTs are required for the HPR-based decryption. As a consequence, the cost of decryption in HPR is nearly  $d$  times smaller than that of [4], [5]. For security reasons,  $K_{RNS}$  and  $K = d \times k$  are in  $\mathcal{O}(n)$ . Hence, since NTTs have a quasi-linear complexity, the asymptotic complexity of the decryption is reduced from:

$$\mathcal{C}(\text{Dec}_{\text{RNS}}) \in \mathcal{O}(n^2 \log n) \text{ to } \mathcal{C}(\text{Dec}_{\text{HPR}}) \in \mathcal{O}(n^{1+\varepsilon} \log n)$$

when  $k = K^\varepsilon$ , with  $0 \leq \varepsilon \leq 1$  ( $\varepsilon = 1$  being the RNS variant, and  $\varepsilon = 0$  the radix- $p_1$  positional variant). As a consequence, when the number of digits  $d$  increases,  $k$  and  $\varepsilon$  decrease, which results in a more efficient decryption algorithm. In particular for  $d = K$ , and thus  $k = 1$ , we obtain a quasi-linear complexity.

### B. Homomorphic Multiplication

Using the same method than previously, since  $K_{RNS}$  and  $K \in \mathcal{O}(n)$ , the asymptotic complexity of [4], [5] is:

$$\mathcal{C}(\text{Mult}_{\text{RNS}}) \in \mathcal{O}(n^3 \log n)$$

while for our HPR variant it is reduced to:

$$\mathcal{C}(\text{Mult}_{\text{HPR}}) \in \mathcal{O}(n^{2+\varepsilon} \log n + n^3) = \mathcal{O}(n^3).$$

Since asymptotically the choice of  $\varepsilon$  does not impact the complexity of homomorphic multiplication, it would seem best to take  $\varepsilon = 0$  so as to minimise the cost of the decryption procedure. However, homomorphic multiplication is executed more often and is the practical bottleneck of these schemes. Hence, one should rather minimise its computational cost by considering the optimal  $k_n$ , or equivalently  $d_n$ , in dimension  $n$ , that minimises the costs presented in Table I.

**Remark 4.** When Remark 3 is considered, the asymptotic complexity of relinearisation can be reduced to

$$\mathcal{O}(dkn \log(n) + nd^2k) = \mathcal{O}(n^2 \log(n) + n^{3-\varepsilon}) = \mathcal{O}(n^{3-\varepsilon}).$$

## VI. EXPERIMENTAL EVALUATION

The NTT-based Fast Lattice (NFL) library [7] was adapted to handle moduli with the bitwidths featured in Table II. Our HPR-based variant and [4], [5] were implemented in C++, using the modified NFL library [7]<sup>1</sup>. They were compiled with gcc 8.3.0 and executed on a single thread on an Intel Core i7-6700K processor running at 4.0GHz with 32GB of main memory, operated by CentOS 7.3. Experimental results for the proposed HPR-based scheme and [4], [5], for the cryptographic parameters presented in Table II, are depicted in Table III.

Asymptotically, decryption has a speed-up that is linear with  $n$ . This behaviour is more noticeable for  $n \geq 2^{14}$  in Table III because therein  $k = 1$  is featured for  $n = 2^{13}$  and  $k = 3$  for  $n \geq 2^{14}$ , since  $k$  was chosen to minimise the complexity of homomorphic multiplication. Maximum speed-ups of 8.6 and 8.0 are achieved when the proposed HPR-based scheme is compared with [4] and [5], respectively, for  $n = 2^{16}$ .

Table III suggests that the theoretically predicted  $\mathcal{O}(\log n)$  speed-up for homomorphic multiplication is only achieved in practice for large parameters. While the novel HPR-based representation brings performance improvements to the multiplication part of the procedure (i.e. without relinearisation), that are noticeable for most parameters, the complexity improvements pertaining to relinearisation only take effect for  $n \geq 2^{16}$ . This suggests that the proposed HPR-based scheme is preferable to [4], [5] for applications relying on sums of products, wherein relinearisation is only applied once, after computing the homomorphic sum of the optimised products, or when large parameters ( $n \geq 2^{16}$ ) are considered. Maximum speed-ups regarding the whole procedure of 2.0 and 1.6 are achieved when the proposed scheme is compared with [5] and [4], respectively.

Moreover, we have measured the depth reached in practice for each scheme by squaring several ciphertexts until they did not decrypt correctly and keeping the smallest value as the practical multiplicative depth. We have not observed any significant difference between the HPR and the RNS variants, which confirms the conclusions of [10].

## CONCLUSION

Whereas state-of-the-art implementations of FV make use of the RNS, the FV scheme requires operations that are inefficient to implement in this representation. In particular, division and rounding require large basis extensions, and the representation of small values is ample. Herein, an alternative approach to accelerate FV was proposed, supported on HPR. Divisions are approximated by considering only the MSDs of HPR values, reducing the asymptotic complexity of decryption from

<sup>1</sup>The code is available at <https://gitlab.com/fvhp/hare>

Scheme	DECRYPTION			MULTIPLICATION & RELINEARISATION		
	NTTs	EMMs	FPOs	NTTs	EMMs	FPOs
[4]	$2dk$	$3ndk$	–	$(dk)^2 + 16dk + 7$	$n(12(dk)^2 + 35dk + 17)$	–
[5]	$2dk$	$ndk$	$n(1 + dk)$	$(dk)^2 + 16dk + 7$	$n(12(dk)^2 + 35dk + 14)$	$n(10dk + 13)$
HPR	$2k$	$3nk$	–	$2dk^2 + (23d + 4)k + 7d + 2$	$n((2d^2 + 8d + 4)k^2 + (5d^2 + 28d + 8)k + 2d^2 + 10d)$	$n(4k + 4)$

Table I: Computational costs of decryption and homomorphic multiplication for the proposed HPR-based scheme, [4] and [5].

Par.	$n$	$\log_2 q$	$K_{RNS}$	$\log_2 q_i$ (RNS)	$K = k \times d$	$\log_2 p_i$ (HPR)
A	$2^{13}$	275	5	55	$5 = 1 \times 5$	55
B	$2^{14}$	549	9	61	$9 = 3 \times 3$	61
C	$2^{15}$	1098	18	61	$18 = 3 \times 6$	61
D	$2^{16}$	2196	36	61	$36 = 3 \times 12$	61

Table II: Parameters achieving at least 100 bits of security according to the estimator of [12] with  $\sigma_{err} = 8.0$ . Values of  $k$  and  $d$  were computed by minimizing the cost of homomorphic multiplication presented in Table I, assuming that 1 NTT =  $n \log_2(n)$  EMMs and 1 FPO = 1 EMM.

Par.	Scheme	Dec.	Mult.	Rel.	Mult. + Relin
A	HPR	<b>0.75</b>	32.0	24.7	56.7
A	[4]	2.64	<b>28.3</b>	13.3	<b>41.6</b>
A	[5]	2.36	36.8	<b>13.2</b>	50.0
B	HPR	<b>6.46</b>	<b>132</b>	135	267
B	[4]	12.8	157	97	<b>254</b>
B	[5]	12.1	199	<b>95</b>	294
C	HPR	<b>14.7</b>	<b>705</b>	892	<b>1,597</b>
C	[4]	64.5	840	<b>775</b>	1,615
C	[5]	58.5	1,164	784	1,948
D	HPR	<b>33.2</b>	<b>3,925</b>	<b>5,688</b>	<b>9,613</b>
D	[4]	285	8,227	7,118	15,345
D	[5]	267	12,298	7,018	19,316

Table III: Execution time (ms) of decryption and homomorphic multiplication for the HPR-based variant, [4] and [5].

$\mathcal{O}(n^2 \log n)$  to  $\mathcal{O}(n \log n)$ . Since the homomorphic multiplication operation also depends on division and rounding, it also benefits from the use of HPR. Moreover, since in HPR small values can be represented as a single digit, the amount of NTTs one needs to compute during relinearisation is significantly reduced. As NTTs no longer dominate the asymptotic complexity of homomorphic multiplication, its complexity is reduced from  $\mathcal{O}(n^3 \log n)$  to  $\mathcal{O}(n^3)$ , when compared to RNS-based approaches. Experimental speedups of up to 8.0 and 8.6 for the decryption operation and of up to 2.0 and 1.6 for the homomorphic multiplication are achieved when compared with two state-of-the-art schemes. The experimental results also confirm the asymptotic gain since the speed-ups increase as  $n$  grows.

#### ACKNOWLEDGEMENTS

This work was partially supported by Portuguese funds through Fundação para a Ciência e a Tecnologia (FCT) with reference UIDB/50021/2020 and by the Ph.D. grant with reference SFRH/BD/103791/2014; by the ANR grant ARR-AND 15-CE39-0002-01; through the Pessoa/Hubert Curien programme with reference 4335 (FCT)/40832XC (CAMPUS-FRANCE); by EU's Horizon 2020 research and innovation programme under grant agreement No. 779391 (FutureTPM).

This work was also partially supported by CyberSecurity Research Flanders with reference VR20192203, by the Office of the Director of National Intelligence (ODNI) and by Intelligence Advanced Research Projects Activity (IARPA) via Contract No. 2019-1902070006. J. Eynard was initially supported by University of Toulon (France) for this project.

Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the ODNI or IARPA.

The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein.

#### REFERENCES

- [1] C. Gentry, "A fully homomorphic encryption scheme," Ph.D. dissertation, Stanford, CA, USA, 2009, aAI3382729.
- [2] J. Fan and F. Vercauteren, "Somewhat Practical Fully Homomorphic Encryption," Cryptology ePrint Archive, Report 2012/144, 2012, <http://eprint.iacr.org/>.
- [3] L. Sousa, S. Antao, and P. Martins, "Combining residue arithmetic to design efficient cryptographic circuits and systems," *IEEE Circuits and Systems Magazine*, vol. 16, no. 4, pp. 6–32, Fourthquarter 2016.
- [4] J.-C. Bajard, J. Eynard, M. A. Hasan, and V. Zucca, "A Full RNS Variant of FV Like Somewhat Homomorphic Encryption Schemes," in *Selected Areas in Cryptography – SAC 2016*, R. Avanzi and H. Heys, Eds. Cham: Springer International Publishing, 2017, pp. 423–442.
- [5] S. Halevi, Y. Polyakov, and V. Shoup, "An Improved RNS Variant of the BFV Homomorphic Encryption Scheme," in *Topics in Cryptology – CT-RSA 2019*, M. Matsui, Ed. Cham: Springer International Publishing, 2019, pp. 83–105.
- [6] K. Bigou and A. Tisserand, "Hybrid Position-Residues Number System," in *ARITH: 23rd Symposium on Computer Arithmetic*, J. Hormigo, S. Oberman, and N. Revol, Eds. Santa Clara, CA, United States: IEEE, Jul. 2016. [Online]. Available: <https://hal.inria.fr/hal-01314232>
- [7] C. Aguilar-Melchor, J. Barrier, S. Guelton, A. Guinet, M.-O. Killijian, and T. Lepoint, *Topics in Cryptology - CT-RSA 2016: The Cryptographers' Track at the RSA Conference 2016, San Francisco, CA, USA, February 29 - March 4, 2016, Proceedings*. Cham: Springer International Publishing, 2016, ch. NFLlib: NTT-Based Fast Lattice Library, pp. 341–356. [Online]. Available: [http://dx.doi.org/10.1007/978-3-319-29485-8\\_20](http://dx.doi.org/10.1007/978-3-319-29485-8_20)
- [8] A. P. Shenoy and R. Kumaresan, "Fast Base Extension Using a Redundant Modulus in RNS," *IEEE Trans. Comput.*, vol. 38, no. 2, pp. 292–297, Feb. 1989. [Online]. Available: <http://dx.doi.org/10.1109/12.16508>
- [9] A. Qaisar Ahmad Al Badawi, Y. Polyakov, K. M. M. Aung, B. Veeravalli, and K. Rohloff, "Implementation and Performance Evaluation of RNS Variants of the BFV Homomorphic Encryption Scheme," *IEEE Transactions on Emerging Topics in Computing*, pp. 1–1, 2019.
- [10] J.-C. Bajard, J. Eynard, P. Martins, L. Sousa, and V. Zucca, "Note on the noise growth of the RNS variants of the BFV scheme," Cryptology ePrint Archive, Report 2019/1266, 2019, <https://eprint.iacr.org/2019/1266>.
- [11] T. Pöppelmann and T. Güneysu, "Towards Efficient Arithmetic for Lattice-Based Cryptography on Reconfigurable Hardware," in *Progress in Cryptology – LATINCRYPT 2012*, A. Hevia and G. Neven, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 139–158.
- [12] M. R. Albrecht, R. Player, and S. Scott, "On the concrete hardness of Learning with Errors." *Journal of Mathematical Cryptology*, vol. 9, pp. 169–203, October 2015.



- [13] J. Bos, K. Lauter, J. Loftus, and M. Naehrig, "Improved Security for a Ring-Based Fully Homomorphic Encryption Scheme," in *Cryptography and Coding*, ser. Lecture Notes in Computer Science, M. Stam, Ed. Springer Berlin Heidelberg, 2013, vol. 8308, pp. 45–64. [Online]. Available: [http://dx.doi.org/10.1007/978-3-642-45239-0\\_4](http://dx.doi.org/10.1007/978-3-642-45239-0_4)

APPENDIX A  
PROOF OF LEMMA 1

A valid encryption of  $[m]_t$  satisfies  $c_0 + c_1 \cdot s = \Delta[m]_t + v + uq$ . Thus we have that the approximation in (4) satisfies:

$$[c_{0,d-1} + c_{1,d-1}s]_p = \frac{\Delta[m]_t + v + qu}{p^{d-1}} - \sum_{i=0}^{d-2} (c_{0,i} + c_{1,i} \cdot s)p^{i-d+1} + \epsilon p$$

with  $e = -\sum_{i=0}^{d-2} (c_{0,i} + c_{1,i} \cdot s)p^{i-d+1}$ . The norm of  $e$  can thus be bounded as

$$\|e\|_\infty \leq \sum_{i=0}^{d-2} (\beta p + \delta_{\mathcal{R}} \beta p B_{\text{key}}) p^{i-d+1} \leq (\beta p + \delta_{\mathcal{R}} \beta p B_{\text{key}}) \frac{p^{d-1}-1}{p^{d-1}} \leq \frac{p}{p-1} \beta (1 + \delta_{\mathcal{R}} B_{\text{key}}). \quad \square$$

APPENDIX B  
PROOF OF LEMMA 2

First, we show that if

$$\|\tilde{e}\|_\infty < \frac{1}{2} \left(1 - \frac{k}{\gamma}\right) \quad (17)$$

is satisfied then  $[\tilde{e}] = 0$  and  $\left\| \left[ \gamma \left[ p\tilde{e} \right]_p / p \right] - \alpha \right\|_\infty < \gamma/2$ . In particular,  $[\tilde{e}] = 0$  since  $k, \gamma > 0$  and hence  $\|\tilde{e}\|_\infty < 1/2$ . In addition,  $[p\tilde{e}]_p = p\tilde{e}$  since  $\|p\tilde{e}\|_\infty < p/2$ . Thus,

$$\left\| \left[ \gamma \left[ p\tilde{e} \right]_p / p \right] - \alpha \right\|_\infty = \left\| \frac{\gamma p\tilde{e} - [\gamma p\tilde{e}]_p}{p} - \alpha \right\|_\infty < \frac{\gamma p(1-k/\gamma)+p}{2p} + \frac{k-1}{2} \leq \gamma/2.$$

Finally, it will be shown that (8) implies (17). This is achieved by exploiting (6) to bound  $\|\tilde{e}\|_\infty$ , and as consequence  $\|v\|_\infty$ :

$$\|\tilde{e}\|_\infty \leq \frac{t}{q} \|v\|_\infty + \frac{t|q|_t}{2q} + \frac{t}{p-1} \beta (1 + \delta_{\mathcal{R}} B_{\text{key}}) < \frac{1}{2} \left(1 - \frac{k}{\gamma}\right).$$

Through a manipulation of the previous expression, (8) is finally reached:

$$\frac{t}{q} \|v\|_\infty + \frac{t|q|_t}{2q} + \frac{t}{p-1} \beta (1 + \delta_{\mathcal{R}} B_{\text{key}}) < \frac{1}{2} \left(1 - \frac{k}{\gamma}\right) \Leftrightarrow \|v\|_\infty < \frac{q}{2t} \left(1 - \frac{k}{\gamma} - \frac{2t\beta}{p-1} (1 + \delta_{\mathcal{R}} B_{\text{key}})\right) - \frac{|q|_t}{2}. \quad \square$$

APPENDIX C  
PROOF OF LEMMA 3

For  $i \in \{0, 1, 2\}$  we have by definition:  $\tilde{c}_i = \sum_{j=0}^{2d-2} \tilde{c}_{i,j} \cdot p^j$ . Moreover, given that the digits of  $c_l$  and  $c'_l$ , for  $l = 0, 1$ , have their norm smaller than  $\beta p$  then:

$$\begin{cases} \text{for } i \in \{0, 2\} \text{ and for } j \in \{0, \dots, d-2\}, \\ \|\tilde{c}_{i,j}\|_\infty \leq \delta_{\mathcal{R}} t(j+1)\beta^2 p^2; \\ \text{for } j \in \{0, \dots, d-2\}, \|\tilde{c}_{1,j}\|_\infty \leq 6\delta_{\mathcal{R}} t(j+1)\beta^2 p^2. \end{cases} \quad (18)$$

Now by using the second bound for every  $\tilde{c}_{i,j}$  we can write:

$$\begin{aligned} & \|\tilde{c}_{i,0} + \dots + \tilde{c}_{i,d-2} \cdot p^{d-2} + ([\tilde{c}_{i,d-1}]_p + p\alpha_{i,d-1}) \cdot p^{d-1}\|_\infty \\ & \leq 6\delta_{\mathcal{R}} t \beta^2 p^2 (1 + 2p + \dots + (d-1)p^{d-2}) + \left(\frac{1}{2} + \frac{k-1}{2}\right) p^d \\ & \leq 6\delta_{\mathcal{R}} t \beta^2 p^2 \left(\frac{(d-1)p^{d-1}}{p-1} - \frac{p^d-1}{p(p-1)^2}\right) + \frac{k}{2} p^d. \end{aligned}$$

Recalling that  $p^d = q$ , we have:

$$\begin{aligned} & \left( \left[ \frac{\tilde{c}_{i,d-1}}{p} \right] - \alpha_{i,d-1} + \tilde{c}_{i,d} + \tilde{c}_{i,d+1} \cdot p + \dots + \tilde{c}_{i,2d-2} \cdot p^{d-2} \right) \\ & = \frac{\tilde{c}_i}{q} - \frac{\tilde{c}_{i,0} + \dots + \tilde{c}_{i,d-2} \cdot p^{d-2}}{p^d} - \frac{[\tilde{c}_{i,d-1}]_p \cdot p^{d-1} + \alpha_{i,d-1} \cdot p^d}{p^d} \\ & \text{and then we deduce (11) with:} \\ & e_i = -\frac{\tilde{c}_{0,0} + \dots + \tilde{c}_{0,d-2} \cdot p^{d-2}}{p^d} - \frac{[\tilde{c}_{0,d-1}]_p \cdot p^{d-1} + \alpha_{0,d-1} \cdot p^d}{p^d}. \end{aligned}$$

In particular, we obtain (12):  $\|e_i\|_\infty \leq 6\delta_{\mathcal{R}} t \beta^2 p \left(\frac{d-1}{p-1} - \frac{1-p^{-d}}{(p-1)^2}\right) + \frac{k}{2}. \quad \square$

APPENDIX D  
PROOF OF PROPOSITION 1

Knowing that  $\Delta t = q - |q|_t$ ,  $[m]_t \cdot [m']_t = [m \cdot m']_t + tr_m$  and  $v \cdot v' = \Delta r_v + [v \cdot v']_\Delta$ , we can write:

$$\begin{aligned} & \frac{t}{q} (c_0 + c_1 \cdot s) \cdot (c'_0 + c'_1 \cdot s) = \frac{t}{q} (\Delta \cdot [m]_t + v + r_q) \cdot (\Delta \cdot [m']_t + v' + r'_q) \\ & = \frac{t}{q} (\Delta^2 [m]_t \cdot [m']_t + \Delta([m]_t \cdot v' + [m']_t \cdot v) + \Delta q([m]_t \cdot r' + [m']_t \cdot r) + v \cdot v' + q(v \cdot r' + v' \cdot r) + q^2 r \cdot r') \\ & = \frac{q-|q|_t}{q} (\Delta([m \cdot m']_t + tr_m) + [m]_t \cdot v' + [m']_t \cdot v + r_v) \\ & + t(v \cdot r' + v' \cdot r) + (q-|q|_t)([m]_t \cdot r' + [m']_t \cdot r) + \frac{t}{q} [v \cdot v']_\Delta + tq r \cdot r' \\ & = \Delta[m \cdot m']_t + [m]_t \cdot v' + [m']_t \cdot v + r_v + t(v \cdot r' + v' \cdot r) \\ & + \frac{t}{q} [v \cdot v']_\Delta - \frac{|q|_t}{q} (\Delta[m \cdot m']_t + [m]_t \cdot v' + [m']_t \cdot v + r_v) \\ & + tq r \cdot r' + (q-|q|_t)([m]_t \cdot r' + [m']_t \cdot r) + qr_m - 2|q|_t r_m \\ & + \frac{|q|_t^2}{q} r_m = \Delta[m \cdot m']_t + v_{\text{mult}} + q([m]_t \cdot r' + [m']_t \cdot r + r_m + tr \cdot r'), \end{aligned}$$

$$\begin{aligned} & \text{with: } v_{\text{mult}} = ([m]_t \cdot v' + [m']_t \cdot v + r_v) \left(1 - \frac{|q|_t}{q}\right) \\ & + t(v \cdot r' + v' \cdot r) + \frac{t}{q} [v \cdot v']_\Delta - \frac{|q|_t}{q} \Delta[m \cdot m']_t \\ & - |q|_t ([m]_t \cdot r' + [m']_t \cdot r) - r_m |q|_t \left(2 - \frac{|q|_t}{q}\right) \end{aligned}$$

As shown in [13], we have:  $\|r_m\|_\infty \leq \frac{\delta_{\mathcal{R}} t}{2}$  and  $\|r_v\|_\infty \leq \frac{\delta_{\mathcal{R}}}{2} \min(\|v\|_\infty, \|v'\|_\infty)$ . Therefore:

$$\begin{aligned} & \|v_{\text{mult}}\|_\infty \\ & \leq \frac{t}{2} \delta_{\mathcal{R}} (\|v'\|_\infty + \|v\|_\infty) + \frac{\delta_{\mathcal{R}}}{2} \min(\|v\|_\infty, \|v'\|_\infty) \\ & + t\delta_{\mathcal{R}} r_\infty (\|v'\|_\infty + \|v\|_\infty) + \frac{t}{q} \frac{\Delta}{2} + \frac{t^2}{2q} \Delta + t^2 \delta_{\mathcal{R}} r_\infty + t^2 \delta_{\mathcal{R}} \\ & \leq \delta_{\mathcal{R}} t \left(r_\infty + \frac{1}{2}\right) (\|v\|_\infty + \|v'\|_\infty) \\ & + \frac{\delta_{\mathcal{R}}}{2} \min(\|v\|_\infty, \|v'\|_\infty) + \frac{1}{2} + t \left(\frac{1}{2} + t\delta_{\mathcal{R}} (1 + r_\infty)\right). \end{aligned}$$

However in the conditions herein considered, we need also to take into account the errors  $e_i$  coming from (11), thus:

$$\begin{aligned} & \hat{c}_0 + \hat{c}_1 \cdot s + \hat{c}_2 \cdot s^2 \\ & = \frac{t}{q} c_0 \cdot c'_0 + e_0 + \left(\frac{t}{q} (c_0 \cdot c'_1 + c_1 \cdot c_0) + e_1\right) \cdot s \\ & + \left(\frac{t}{q} c_1 \cdot c'_1 + e_2\right) \cdot s^2 \\ & = \frac{t}{q} (c_0 + c_1 \cdot s) \cdot (c'_0 + c'_1 \cdot s) + e_0 + e_1 \cdot s + e_2 \cdot s^2. \end{aligned}$$

Therefore,  $\hat{v} = v_{\text{mult}} + e_0 + e_1 \cdot s + e_2 \cdot s^2$ , with  $\|e_0 + e_1 \cdot s + e_2 \cdot s^2\|_\infty \leq B_e \left(1 + \delta_{\mathcal{R}} B_{\text{key}} + \delta_{\mathcal{R}}^2 B_{\text{key}}^2\right)$ . We obtain (13) by summing the last two bounds.  $\square$