# Compact Ring-LWE Cryptoprocessor

Sujoy Sinha Roy<sup>1</sup>, Frederik Vercauteren<sup>1</sup>, Nele Mentens<sup>1</sup>, Donald Donglong Chen<sup>2</sup> and Ingrid Verbauwhede<sup>1</sup>

<sup>1</sup>ESAT/COSIC and iMinds, KU Leuven Kasteelpark Arenberg 10, B-3001 Leuven-Heverlee, Belgium Email: {firstname.lastname}@esat.kuleuven.be
<sup>2</sup>Department of Electronic Engineering, City University of Hong Kong Tat Chee Avenue, Kowloon, Hong Kong SAR Email: donald.chen@my.cityu.edu.hk

Abstract. In this paper we propose an efficient and compact processor for a ring-LWE based encryption scheme. We present three optimizations for the Number Theoretic Transform (NTT) used for polynomial multiplication: we avoid preprocessing in the negative wrapped convolution by merging it with the main algorithm, we reduce the fixed computation cost of the twiddle factors and propose an advanced memory access scheme. These optimization techniques reduce both the cycle and memory requirements. Finally, we also propose an optimization of the ring-LWE encryption system that reduces the number of NTT operations from five to four resulting in a 20% speed-up. We use these computational optimizations along with several architectural optimizations to design an instruction-set ring-LWE cryptoprocessor. For dimension 256, our processor performs encryption/decryption operations in  $20/9 \ \mu s$  on a Virtex 6 FPGA and only requires 1349 LUTs, 860 FFs, 1 DSP-MULT and 2 BRAMs. Similarly for dimension 512, the processor takes  $48/21 \ \mu s$  for performing encryption/decryption operations and only requires 1536 LUTs, 953 FFs, 1 DSP-MULT and 3 BRAMs. Our processors are therefore more than three times smaller than the current state of the art hardware implementations, whilst running somewhat faster.

**Keywords.** Lattice-based cryptography, ring-LWE, Polynomial multiplication, Number Theoretic Transform, Hardware implementation

### 1 Introduction

Lattice-based cryptography is considered a prime candidate for quantum-secure public key cryptography due to its wide applicability [27] and its security proofs that are based on worst-case hardness of well known lattice problems. The *learning with errors* (LWE) problem [26] and its ring variant known as ring-LWE [17] have been used as a solid foundation for several cryptographic schemes. The significant progress in the theory of lattice-based cryptography [19, 20, 25] has recently been followed by practical implementations [1, 7, 9, 22, 23, 28].

The ring-LWE based cryptosystems operate in a polynomial ring  $R_q = \mathbb{Z}_q[\mathbf{x}]/\langle f(x) \rangle$ , where one typically chooses  $f(x) = x^n + 1$  with n a power of two, and q a prime with  $q \equiv 1 \mod 2n$ . An implementation thus requires the basic operations in such a ring  $R_q$ , with multiplication taking up the bulk of the resources both in area and time. An efficient polynomial multiplier architecture therefore is a pre-requisite for the deployment of ring-LWE based cryptography in real world systems.

The most important hardware implementations of polynomial multipliers for the rings  $R_q$  are [1,9,22,23]. In [9], a fully parallel butterfly structure is used for the polynomial multiplier resulting in a huge area consumption. For instance, even for medium security, their ring-LWE cryptoprocessor does not fit on the largest FPGA of the Virtex 6 family. In [22], a sequential polynomial multiplier architecture is designed to use the FPGA resources in an efficient way. The multiplier uses a dedicated ROM to store all the twiddle factors which are required during the NTT computation. In [23] the authors integrated the polynomial multiplier [22] in a complete ring-LWE based encryption system and propose several system level optimizations such as a better message encoding scheme and compression technique for the ciphertext. The work [1] tries to reduce the area of the polynomial multiplier by computing the twiddle factors whenever required, but as we will show, this could be improved substantially by re-arranging the loops inside the NTT computation. Furthermore, the paper does not include an implementation of a complete ring-LWE cryptoprocessor.

**Our contributions:** In this paper we present a complete ring-LWE based encryption processor that uses the Number Theoretic Transform (NTT) algorithm for polynomial multiplication. The architecture is designed to have small area and memory requirement, but is also optimized to keep the number of cycles small. In particular, we make the following contributions:

- 1. During the NTT computation, the intermediate coefficients are multiplied by the twiddle factors that are computed using repeated multiplications. In [22] a precomputed table (ROM) is used to avoid this fixed computation cost. The more compact implementation in [1] does not use ROM and computes the twiddle factors by performing repeated multiplications. In this paper we reduce the number of multiplications by re-arranging the nested loops in the NTT computation.
- 2. The implementations [1,22] use negative wrapped convolution to reduce the number of evaluations in both the forward and backward NTT computations. However, the use of the negative wrapped convolution has a pre- and post-computation overhead. In this paper we basically avoid the pre-computation which reduces the cost of the forward NTT.
- 3. The intermediate coefficients are stored in memory (RAM) during the NTT computation. Access to the RAM is a bottleneck for speeding-up the NTT computation. In the implementations [1,22], FPGA-RAM slices are placed in parallel to avoid this bottleneck. In this paper we propose an efficient memory access scheme which reduces the number of RAM accesses, optimizes the number of block RAMs and still achieves maximum utilization of the computational blocks.
- 4. The Knuth-Yao sampler [28] is slow due to the costly bit scanning operation. We reduce the cycle count using fast table lookup operations. We also optimize the area of the Knuth-Yao [28] sampler by reducing the width of the ROM. For the standard deviation 3.33 the area-optimized sampler consumes only 32 slices and is thus more compact and faster than the Bernoulli sampler in [24].
- 5. The proposed optimization techniques are applied to design a compact architecture for the NTT computation. We also implement pipelines in the architecture targeting high-speed applications. The pipeline technique derives an optimal pipeline depth for the architecture to achieve the fastest computation time.

6. Finally, we optimize one of the most popular ring-LWE encryption schemes by reducing the number of NTT computations from five to four, thereby achieving a nearly 20% reduction in the computation cost.

The above optimizations result in a very compact architecture that is three times smaller than the current state of the art implementation [23] and even runs somewhat faster.

The remainder of the paper is organized as follows: In Section 2 we provide a brief mathematical background on ring-LWE and the NTT. Section 3 contains our optimization techniques of the NTT and Section 4 presents the actual architecture of our optimized NTT algorithm. A pipelined architecture is given in Section 5. In Section 6, we propose an optimization of an existing ring-LWE encryption scheme and propose an efficient architecture for the complete ring-LWE encryption system. Finally, Section 7 reports on the experimental results of this implementation.

### 2 Background

In this section we present a brief mathematical overview of the ring-LWE problem, the encryption scheme we will be optimizing and the NTT.

#### 2.1 The LWE and ring-LWE Problem

The *learning with errors* (LWE) problem is a machine learning problem that is equivalent to worst-case lattice problems as shown by Regev [26] in 2005. Since then, the LWE problem has become popular as a basis for developing quantum secure lattice-based cryptosystems.

The LWE problem is parametrized by a dimension  $n \ge 1$ , an integer modulus  $q \ge 2$ and an error distribution, typically a discrete Gaussian distribution  $\mathcal{X}_{\sigma}$  over the integers with deviation  $\sigma$  and mean 0. The probability of sampling an integer  $z \in \mathbb{Z}$  in the Gaussian distribution  $\mathcal{X}_{\sigma}$  is given by  $\rho_{\sigma}(z)/\rho_{\sigma}(\mathbb{Z})$  where  $\rho_{\sigma}(z) = \exp\left(\frac{-z^2}{2\sigma^2}\right)$  and  $\rho_{\sigma}(\mathbb{Z}) =$  $\sum_{z=-\infty}^{+\infty} \rho_{\sigma}(z)$ . Note that some authors use the parameter  $s = \sqrt{2\pi\sigma}$  to define the Gaussian distribution or even denote the parameter s by  $\sigma$  to add to the confusion.

For a uniformly chosen  $\mathbf{s} \in \mathbb{Z}_q^n$ , the LWE distribution  $A_{s,\mathcal{X}}$  over  $\mathbb{Z}_q^n \times \mathbb{Z}_q$  consists of tuples  $(\mathbf{a}, t)$  where  $\mathbf{a}$  is chosen uniformly from  $\mathbb{Z}_q^n$  and  $t = \langle \mathbf{a}, \mathbf{s} \rangle + e \mod q \in \mathbb{Z}_q$  and e is sampled from the error distribution  $\mathcal{X}$ . The *search* version of the LWE problem asks to find  $\mathbf{s}$  given a polynomial number of pairs  $(\mathbf{a}, t)$  sampled from the LWE distribution  $A_{s,\mathcal{X}}$ . In the *decision* version of the LWE problem, the solver needs to distinguish with non-negligible advantage between a polynomial number of samples drawn from  $A_{s,\mathcal{X}}$  for hardness proofs of the search and decision LWE problems, interested readers are referred to [15].

The initial LWE encryption system in [26] is based on matrix operations which are quite inefficient and result in large key sizes. To achieve computational efficiency and to reduce the key size, an algebraic variant of the LWE called *ring*-LWE [17] uses special structured ideal lattices. Such lattices correspond to ideals in rings  $\mathbb{Z}[\mathbf{x}]/\langle f \rangle$ , where f is an irreducible polynomial of degree n. For efficiency reasons, the ring is often taken as  $R_q = \mathbb{Z}_q[\mathbf{x}]/\langle f \rangle$  with  $f(x) = x^n + 1$ , where n is a power of two and the prime q is taken as  $q \equiv 1 \mod 2n$ . The ring-LWE distribution on  $R_q \times R_q$  consists of tuples (a, t)

with  $a \in R_q$  chosen uniformly random and  $t = as + e \in R_q$ , where  $s \in R_q$  is a fixed secret element and e has small coefficients sampled from the discrete Gaussian above. The resulting distribution on  $R_q$  will also be denoted  $\mathcal{X}_{\sigma}$ .

The ring-LWE based encryption scheme that we will use was introduced in the full version of [17] and uses a global polynomial  $a \in R_q$ . Key generation, encryption and decryption are as follows:

- 1. KeyGen(a): Choose two polynomials  $r_1, r_2 \in R_q$  from  $\mathcal{X}_{\sigma}$  and compute  $p = r_1 a \cdot r_2 \in R_q$ . The public key is (a, p) and the private key is  $r_2$ . The polynomial  $r_1$  is simply noise and is no longer required after key generation.
- 2. Enc(a, p, m): The message m is first encoded to  $\overline{m} \in R_q$ . Three polynomials  $e_1, e_2, e_3 \in R_q$  are sampled from  $\mathcal{X}_{\sigma}$ . The ciphertext then consists of two polynomials  $c_1 = a \cdot e_1 + e_2$  and  $c_2 = p \cdot e_1 + e_3 + \overline{m} \in R_q$ .
- 3.  $Dec(c_1, c_2, r_2)$ : Compute  $m' = c_1 \cdot r_2 + c_2 \in R_q$  and recover the original message m from m' using a decoder.

One of the simplest encoding functions maps a binary message m to the polynomial  $\bar{m} \in R_q$  such that its *i*-th coefficient is (q-1)/2 iff the *i*-th bit of m is 1 and 0 otherwise. The corresponding decoding function then simply reduces the coefficients  $m'_i$  of m' in the interval (-q/2, q/2) and decodes to 1 when  $|m'_i| > q/4$  and 0 otherwise.

#### 2.2 Parameter sets

To enable fair comparison with the state of the art [23], we have chosen to instantiate the cryptoprocessor for the same parameter sets (n, q, s) (recall  $s = \sqrt{2\pi\sigma}$ ), namely  $P_1 = (256, 7681, 11.32)$  and  $P_2 = (512, 12289, 12.18)$ . Note that the choice of primes is not optimal for fast modular reduction. To estimate the security level offered by these two parameter sets we follow the security analysis in [16] and [14] which improves upon [15, 29]. Apart from the dimension n, the hardness of the ring-LWE problem mainly depends on the ratio  $q/\sigma$ , where clearly the problem becomes easier for larger ratios. Although neither parameter set was analyzed in [16], parameter set  $P_1$  is similar to the set (256, 4093, 8.35) from [16] which requires  $2^{105}$  seconds to break, or still over  $2^{128}$ elementary operations. For parameter set  $P_2$  we expect it to offer a high security level consistent with AES-256 (following [9]).

We limit the Gaussian sampler in our implementation to  $12\sigma$  to obtain a negligible statistical distance ( $< 2^{-90}$ ) from the true discrete Gaussian distribution. Although one can normally sample the secret  $r_2 \in R_q$  also from the distribution  $\mathcal{X}_{\sigma}$ , we restrict  $r_2$  to have binary coefficients.

#### 2.3 The Number Theoretic Transform

There are many efficient algorithms in the literature to perform polynomial multiplication and a survey of fast multiplication algorithms can be found in [2]. In this section we review the Number Theoretic Transform (NTT) which corresponds to a Fast Fourier Transform (FTT) where the roots of unity are taken from a finite ring instead of the complex numbers.

Algorithm 1: Iterative NTT

```
Input: Polynomial a(x) \in \mathbb{Z}_q[\mathbf{x}] of degree n-1 and n-th primitive root \omega_n \in \mathbb{Z}_q of unity Output: Polynomial A(x) \in \mathbb{Z}_q[\mathbf{x}] = \text{NTT}(a)
      begin
              A \leftarrow BitReverse(a);
  2
 3
             for m = 2 to n by m = 2m do

\omega_m \leftarrow \omega_n^{n/m};
  \mathbf{4}
  5
                    \omega \leftarrow 1 ~;
                    for j = 0 to m/2 - 1 do
  6
  7
                            for k = 0 to n - 1 by m do
                                   t \leftarrow \omega \cdot A[k+j+m/2];
  8
 9
                                   u \leftarrow A[k + j] \ ;
10
                                   A[k+j] \leftarrow u+t ;
11
                                   A[k+j+m/2] \leftarrow u-t ;
12
                           end
13
                           \omega \leftarrow \omega \cdot \omega_m
\mathbf{14}
                    end
15
             end
16 end
```

The FFT and NTT: Recall that the n-point FFT (with  $n = 2^k$ ) is an efficient method to evaluate a polynomial  $a(x) = \sum_{j=0}^{n-1} a_j x^j \in \mathbb{Z}[x]$  in the *n*-th roots of unity  $\omega_n^i$  for  $i = 0, \ldots, n-1$  where  $\omega_n$  denotes a primitive *n*-th root of unity. More precisely, on input the coefficients  $[a_0, \ldots, a_{n-1}]$  and  $\omega_n$ , the FFT computes  $FFT([a_j], \omega_n) = [a(\omega_n^0), a(\omega_n^1), \ldots, a(\omega_n^{n-1})]$  in  $\theta(n \log n)$  time. Due to the orthogonality relations between the *n*-th roots of unity, we can compute the inverse FFT simply as  $\frac{1}{n}FFT(\cdot, \omega_n^{-1})$ .

The NTT replaces the complex roots of unity by roots of unity in a finite ring  $\mathbb{Z}_q$ . Since we require elements of order n, q is chosen to be a prime with  $q \equiv 1 \mod 2n$ . Note furthermore that the NTT immediately leads to a fast multiplication algorithm in the ring  $S_q = \mathbb{Z}_q[x]/(x^n - 1)$ : indeed, given two polynomials  $a, b \in S_q$  we can easily compute their (reduced) product  $c = a \cdot b \in S_q$  by computing

$$c = NTT_{\omega_n}^{-1} \left( NTT_{\omega_n}(a) * NTT_{\omega_n}(b) \right), \tag{1}$$

where \* denotes point-wise multiplication.

The NTT computation is usually described as recursive, but in practice we use an in-place iterative version taken from [4] that is given in Algorithm 1. For the inverse NTT, an additional scaling of the resulting coefficients by  $n^{-1}$  is performed. The factors  $\omega$  used in line 8 are called the *twiddle factors*.

**Multiplication in**  $R_q$ : Recall that we will use  $R_q = Z_q[\mathbf{x}]/\langle f \rangle$  with  $f = x^n + 1$  and  $n = 2^k$ . Since  $f(x)|x^{2n} - 1$  we could use the 2*n*-point NTT to compute the multiplication in  $R_q$  at the expense of three 2*n*-point NTT computations and a reduction by trivially embedding the ring  $R_q$  into  $S_q$ , i.e. expanding the coefficient vector of a polynomial  $a \in R_q$  by adding *n* extra zero coefficients. However, we can do much better by exploiting the special relation between the roots of  $x^n + 1$  and  $x^{2n} - 1$  using a technique known as the *negative wrapped convolution*.

Indeed, using the same evaluation-interpolation strategy used above for the ordinary NTT, we conclude that we can efficiently multiply two polynomials  $a, b \in R_q$  if we can quickly evaluate them in the roots of f. These roots are simply  $\omega_{2n}^{2j+1}$  for  $j = 0, \ldots, n-1$  (since the even exponents give the roots of  $x^n - 1$ ) and as such can be written as  $\omega_{2n} \cdot \omega_n^j$ . These evaluations can thus be computed efficiently using a *classical* 

*n*-point NTT (instead of a 2*n*-point NTT) on the scaled polynomials  $a'(x) = a(\omega_{2n} \cdot x)$ and  $b'(x) = a(\omega_{2n} \cdot x)$ . The point-wise multiplication gives the evaluations of  $c(x) = a(x)b(x) \mod f(x)$  in the roots of f, and the classical inverse *n*-point NTT thus results in the coefficients of the scaled polynomial  $c'(x) = c(\omega_{2n} \cdot x)$ . To recover the coefficients  $c_i$  of c(x), we therefore simply have to compute  $c_i = c'_i \cdot \omega_{2n}^{-i}$ . Note that the scaling operation by  $n^{-1}$  can be combined with the multiplications of  $c'_i$  by  $\omega_{2n}^{-i}$ .

### **3** Optimization of the NTT Computation

In this section we optimize the NTT and compare with the recent hardware implementations of polynomial multipliers [1, 22, 23]. First, the fixed cost involved in computing the powers of  $\omega_n$  is reduced, then the pre-computation overhead in the forward negativewrapped convolution is optimized, and finally an efficient memory access scheme is proposed that reduces the number of memory accesses during the NTT and also minimizes the number of block RAMs in the hardware architecture.

#### 3.1 Optimizing the Fixed Computation Cost

In line 13 of Algorithm 1 the computation of the twiddle factor  $\omega \leftarrow \omega \cdot \omega_m$  is performed in the *j*-loop. This computation can be considered as a fixed cost. However in [1, 22] the *j*-loop and the *k*-loop are interchanged, such that  $\omega$  is updated in the innermost loop which is much more frequent than in Algorithm 1. To avoid the computation of the twiddle factors, in [22] all the twiddle factors are kept in a pre-computed look-up table (ROM) and are accessed whenever required. As the twiddle factors are not computed on-the-fly, the order of the two innermost loops does not result in an additional cost. However in [1] a more compact polynomial multiplier architecture is designed without using any look-up table and the twiddle factors are simply computed on-the-fly during the NTT computation. Hence in [1], the interchanged loops cause substantial additional computational overhead. In this paper our target is to design a very compact polynomial multiplier. Hence we do not use any look-up table for the twiddle factors and follow Algorithm 1 to avoid the extra computation of [1].

#### 3.2 Optimizing the Forward NTT Computation Cost

Here we revisit the forward negative-wrapped convolution technique used in [1, 22, 23]. Recall that the negative-wrapped convolution corresponds to a classical *n*-point NTT on the scaled polynomials  $a'(x) = a(\omega_{2n} \cdot x)$  and  $b'(x) = (\omega_{2n} \cdot x)$ . Instead of first precomputing these scaled polynomials and then performing a classical NTT, it suffices to note that we can integrate the scaling and the NTT computation. Indeed, it suffices to change the initialization of the twiddle factors in line 5 of Algorithm 1: instead of initializing  $\omega$  to 1, we can simply set  $\omega = \omega_{2m}$ . The rest of the algorithm remains exactly the same, and no pre-computation is necessary. Note that this optimization only applies to the NTT itself and not to the inverse NTT.

#### 3.3 Optimizing the Memory Access Scheme

The NTT computation requires memory to store the input and intermediate coefficients. When the number of coefficients is large, RAM is most suitable for hardware implementation [1, 22, 23]. In the innermost loop (lines 8-to-11) of Algorithm 1, two coefficients A[k+j] and A[k+j+m/2] are first read from memory and then arithmetic operations (one multiplication, one addition and one subtraction) are performed. The new A[k+j]and A[k+j+m/2] are then written back to memory. During one iteration of the innermost loop, the arithmetic circuits are thus used only once, while the memory is read or written twice. This leads to idle cycles in the arithmetic circuits. The polynomial multiplier in [22] uses two parallel memory blocks to provide a continuous flow of coefficients to the arithmetic circuits. However this approach could result in under-utilization of the RAM blocks if the coefficient size is much smaller than the word size (for example in the ring-LWE cryptosystem [17]). In the literature there are many papers on efficient memory management schemes using segmentation and efficient address generation (see [18]) for the classical FFT algorithm. Another well known approach is the constant geometry FFT (or NTT) which always maintains a constant index difference between the processed coefficients [21]. However the constant geometry algorithm is not in-place and hence not suitable for resource constrained platforms. In [1] memory usage is improved by keeping two coefficients A[k] and B[k] of the two input polynomials A and B in the same memory location. We propose a memory access scheme which is designed to minimize the number of block RAM slices and to achieve maximum utilization of computational circuits present in the NTT architecture.

Since the two coefficients A[k+j] and A[k+j+m/2] are processed together in Algorithm 1, we keep the two coefficients as a pair in one memory location.

Let us analyze two consecutive iterations of the m-loop (line 3 in Algorithm 1) for  $m = m_1$  and  $m = m_2$  where  $m_2 = 2m_1$ . In the  $m_1$ -loop, for some  $j_1$  and  $k_1$  (maintaining the loop bounds in Algorithm 1) the coefficients  $(A[k_1 + j_1], A[k_1 + j_1 + m_1/2])$  are processed as a pair. Then k increments to  $k_1 + m_1$  and the processed coefficient pair is  $(A[k_1+m_1+j_1], A[k_1+m_1+j_1+m_1/2])$ . Now from Algorithm 1 we see that the coefficient  $A[k_1 + j_1]$  will again be processed in the  $m_2$ -loop with coefficient  $A[k_1 + j_1 + m_2/2]$ . Since  $m_2 = 2m_1$ , the coefficient  $A[k_1 + j_1 + m_2/2]$  is the coefficient  $A[k_1 + j_1 + m_2/2]$ . Since  $m_2 = 2m_1$ , the coefficient  $A[k_1 + j_1 + m_2/2]$  is the coefficient  $A[k_1 + j_1 + m_1]$  which is updated in the  $m_1$ -loop for  $k = k_1 + m_1$ . Hence during the  $m_1$ -loop if we swap the updated coefficients for  $k = k_1$  and  $k = k_1 + m_1$  and store  $(A[k_1 + j_1], A[k_1 + j_1 + m_1])$  and  $(A[k_1 + j_1 + m_1/2], A[k_1 + j_1 + 3m_1/2])$  as the coefficient pairs in memory, then the coefficients in a pair have a difference of  $m_2/2$  in their index and thus are ready for the  $m_2$ -loop. The operations during the two consecutive iterations  $k = k_1$  and  $k = k_1 + m_1$  during  $m = m_1$  are shown in Algorithm 2 in lines 8-15. During the operations  $u_1, t_1, u_2$  and  $t_2$  are used as temporary storage registers.

A complete description of the efficient memory access scheme is given in Algorithm 2. In this algorithm for all values of m < n, two coefficient pairs are processed in the innermost loop and a swap of the updated coefficients is performed before writing back to memory. For m = n, no swap operation is required as this is the final iteration of the *m*-loop. The coefficient pairs generated by Algorithm 2 can be re-arranged easily for another (say inverse) NTT operation by performing address-wise bit-reverse-swap operation. Appendix A describes the memory access scheme using an example.

### 4 The NTT Processor Organization

In this section we present an architecture for performing the forward and backward NTT using the proposed optimization techniques. Our NTT processor (Figure 1) consists of three main components: the arithmetic unit, the memory block and the control-address unit.

Algorithm 2: Iterative NTT : Memory Efficient Version

**Input:** Polynomial  $a(x) \in \mathbb{Z}_q[\mathbf{x}]$  of degree n-1 and *n*-th primitive root  $\omega_n \in \mathbb{Z}_q$  of unity **Output:** Polynomial  $A(x) \in \mathbb{Z}_q[\mathbf{x}] = \text{NTT}(a)$ 1 begin A  $\leftarrow$  BitReverse(a); /\* Coefficients are stored in the memory as proper pairs \*/ for m = 2 to n/2 by m = 2m do  $\omega_m \leftarrow m$ -th primitiveroot(1);  $\omega \leftarrow squareroot(\omega_m)$  or 1 /\* Depending on forward or backward NTT \*/; for j = 0 to m/2 - 1 do  ${f 2}{f 3}$ 4 5 6 7 8 9  $\begin{aligned} \mathbf{for} &= 0 \ to \ m/2 - 1 \ by \ m \ \mathbf{do} \\ &(t_1, u_1) \leftarrow (A[k+j+m/2], A[k+j]) \ /* \ \text{From MEMORY[k+j] }*/ \ ; \\ &(t_2, u_2) \leftarrow (A[k+m+j+m/2], A[k+m+j]) \ /* \ \text{MEMORY[k+j+m/2] }*/ \ ; \end{aligned}$ 10  $t_1 \leftarrow \omega \cdot t_1$ ;  $t_2 \leftarrow \omega \cdot t_2$ ; 11  $\begin{array}{l} & (A[k+j+m/2], A[k+j]) \leftarrow (u_1 - t_1, u_1 + t_1) ; \\ & (A[k+m+j+m/2], A[k+m+j]) \leftarrow (u_2 - t_2, u_2 + t_2) ; \\ & MEMORY[k+j] \leftarrow (A[k+j+m], A[k+j]) ; \\ & MEMORY[k+j+m/2] \leftarrow (A[k+j+3m/2], A[k+j+m/2]) ; \end{array}$ 12 $\mathbf{13}$  $14\\15$ 16 end 17 $\omega \leftarrow \omega \cdot \omega_n \ ;$  $\mathbf{18}$  $\mathbf{end}$ 19  $\mathbf{end}$  $\mathbf{20}$  $m \leftarrow$ n21 22  $k \leftarrow 0$ ;  $\omega \leftarrow squareroot(\omega_m)$  or 1 /\* Depending on forward or backward NTT \*/; for j = 0 to m/2 - 1 do  $(t_1, u_1) \leftarrow (A[j + m/2], A[j])$  /\* From MEMORY[j] \*/;  $t_1 \leftarrow \cdots t_r$ ; 23  $\mathbf{24}$  $\mathbf{25}$  $t_1 \leftarrow \omega \cdot t_1$ ;  $\begin{array}{l} (A[j+m/2], A[j]) \leftarrow (u_1 - t_1, u_1 + t_1) ; \\ MEMORY[j] \leftarrow (A[j+m/2], A[j]) ; \end{array}$  $\frac{26}{27}$ 28  $\omega \leftarrow \omega \cdot \omega_m \ ;$  $\mathbf{end}$ 29 **30** end



Fig. 1. Hardware Architecture for NTT

**The Memory Block** is implemented as a simple dual port RAM. To accommodate two coefficients, the word size is  $2\lceil \log q \rceil$  where q is the prime modulus. In FPGAs, a RAM can be implemented as a *distributed* or as a *block* RAM. When the amount of data is large, block RAM is the ideal choice.

The Arithmetic Unit (NTT-ALU) is designed to support Algorithm 2 along with other operations such as polynomial addition, point-wise multiplication and rearrangement of the coefficients. This NTT-ALU is interfaced with the memory block and the control-address unit. The central part of the NTT-ALU consists of a modular multiplier and addition/subtraction circuits.

Now we describe how the different components of the NTT-ALU are used during the butterfly steps (excluding the last loop for m = n). First, the memory location (k + j)is fetched and then the fetched data  $(t_1, u_1)$  is stored in the input register pair  $(H_1, L_1)$ . The same also happens for the memory location (k + j + m/2) in the next cycle. The multiplier computes  $\omega \cdot H_1$  and the result is added to or subtracted from  $L_1$  using the adder and subtracter circuits to compute  $(u_1 + \omega t_1)$  and  $(u_1 - \omega t_1)$  respectively. In the next cycle the register pair  $(R_1, R_4)$  is updated with  $(u_1 - \omega t_1, u_1 + \omega t_1)$ . Another clock transition shifts the contents of  $(R_1, R_4)$  to  $(R_2, R_5)$ . In this cycle the pair  $(R_1, R_4)$  is updated with  $(u_2 - \omega t_2, u_2 + \omega t_2)$  as the computation involving  $(u_2, t_2)$  from the location (k+j+m/2) lags by one cycle. Now the memory location (k+j) is updated with the register pair  $(R_4, R_5)$  containing  $(u_2 + \omega t_2, u_1 + \omega t_1)$ . Finally, in the next cycle the memory location (k + j + m/2) is updated with  $(u_2 - \omega t_2, u_1 - \omega t_1)$  using the register pair  $(R_2, R_3)$ . The execution of the last m-loop is similar to the intermediate loops, without any data swap between the output registers. The register pair  $(R_2, R_5)$  is used for updating the memory locations. In Figure 1, the additional registers  $(H_2, H_3)$  and  $L_2$ ) and multiplexers are used for supporting operations such as addition, point-wise multiplication and rearrangement of polynomials. The Small-ROM block contains the fixed values  $\omega_m$ ,  $\omega_{2n}$ , their inverses and  $n^{-1}$ . This ROM has depth of order  $\log(n)$ .

The Control-and-Address Unit consists of three counters for m, j and k in Algorithm 2 and comparators to check the terminal conditions during the execution of any loop. The read address is computed from m, j and k and then delayed using registers to generate the write address. The control-and-address unit also generates the write enable signal for the RAM and the control signals for the NTT-ALU.

### 5 Pipelining the NTT Processor

The maximum frequency of the NTT-ALU is determined by the critical path (red dashed line in Figure 1) that passes through the modular multiplier and the adder (or sub-tracter) circuits. To increase the operating frequency of the processor, we implement efficient pipelines based on the following two observations.

**Observation 1:** During the execution of any *m*-loop in Algorithm 2, the computations (multiplication, addition and subtraction) involving a coefficient pair have no data dependency on other coefficient pairs. Such a data-flow structure is suitable for pipeline processing as different computations can be pipelined without inserting bubbles in the datapath.

Assume that the modular multiplier has  $d_m$  pipeline stages and that the output is latched in a buffer. In the  $(d_m + 1)$ th cycle after the initialisation of  $\omega \cdot t_1$ , the buffer is updated with the result  $\omega \cdot t_1$ . Now we need to compute  $u_1 + \omega \cdot t_1$  and  $u_1 - \omega \cdot t_1$  using the adder and subtracter circuits. Hence we delay the data  $u_1$  by  $d_m$  cycles so that it appears as an input to the adder and subtracter circuits in the  $(d_m + 1)$ th cycle. This delay operation is performed with the help of a shift register  $L_1, \ldots, L_{d_m+1}$  as shown in Figure 2.

**Observation 2:** Every increment of j in Algorithm 2 requires a new  $\omega$  (line 17). If the multiplier has  $d_m$  pipeline stages, then the register- $\omega$  in Figure 1 is updated with the new value of  $\omega$  in the  $(d_m + 2)$ th cycle. Since this new  $\omega$  is used by the next butterfly operations, the data dependency results in an interruption in the chain of butterfly operations for  $d_m + 1$  cycles. In any *m*-loop, the total number of such *interruption cycles* is  $(m/2 - 1) \cdot (d_m + 1)$ .

To reduce the number of interruption cycles, we use a small look-up table to store a few twiddle factors. Let the look-up table (red dashed rectangle in Figure 2) have l registers containing the twiddle factors  $(\omega, \ldots, \omega \omega_m^{l-1})$ . This look-up table is used to provide the twiddle factors during the butterfly operations for say j = j' to j = j' + l - 1. The next time j increments, new twiddle factors are required for the butterfly operations. We multiply the look-up table with  $\omega_m^l$  to compute the next l twiddle factors  $(\omega \omega_m^l, \ldots \omega \omega_m^{2l-1})$ . The multiplications are independent of each other and hence can be processed in a pipeline. The butterfly operations are resumed after  $\omega \omega_m^l$  is loaded in the look-up table. Thus using a small-look-up table of size l we reduce the number of interruption cycles to  $(\frac{m}{2l} - 1) \cdot (d_m + 1)$ . In our architecture we use l = 4; a larger value of l will reduce the number of interruption cycles, but will cost additional registers.



Fig. 2. Pipelined Hardware Architecture for NTT

**Optimal Pipeline Strategy for Speed :** During the execution of any *m*-loop in Algorithm 2, the number of butterfly operations is n/2. In the pipelined NTT-ALU, the cycle requirement for the n/2 butterfly operations is slightly larger than n/2 due to an initial overhead. The state machine jumps to the  $\omega$  calculation state  $\frac{m}{2l} - 1$  times resulting in  $(\frac{m}{2l} - 1) \cdot (d_m + 1)$  interruption cycles. Hence the total number of cycles spent in executing any *m*-loop can be approximated as shown below:

$$Cycles_m \approx \frac{n}{2} + (\frac{m}{2l} - 1) \cdot (d_m + 1)$$

Let us assume that the delay of the critical path with no pipeline stages is  $D_{comb}$ . When the critical path is split in balanced-delay stages using pipelines, the resulting delay  $(D_s)$ can be approximated as  $\frac{D_{comb}}{(d_m+d_a)}$ , where  $d_m$  and  $d_a$  are the number of pipeline stages in the modular multiplier and the modular adder (subtracter) respectively. Since the delay of the modular adder is small compared to the modular multiplier, we have  $d_a \ll d_m$ . Now the computation time for the *m*-loop is approximated as

$$T_m \approx \frac{D_{comb}}{(d_m + d_a)} \left[ \frac{n}{2} + \left( \frac{m}{2l} - 1 \right) \cdot (d_m + 1) \right] \approx D_s \frac{n}{2} + C_m \,.$$

Here  $C_m$  is constant (assuming  $d_a \ll d_m$ ) for a fixed value of m. From the above equation we find that the minimum computation time can be achieved when  $D_s$  is minimum. Hence we pipeline the datapath to achieve minimum  $D_s$ . The DSP based coefficient multiplier is optimally pipelined using the Xilinx IPCore tool, while the modular reduction block is suitably pipelined by placing registers between the cascaded adder and subtracter circuits.

### 6 The ring-LWE Encryption Scheme

The ring-LWE encryption scheme in [23] optimizes computation cost by keeping the fixed polynomials in the NTT domain. The message encryption and decryption operations require three and two NTT computations respectively. In this paper we reduce the number of NTT operations for decryption from two to *one*. The proposed ring-LWE encryption scheme is described below:

- 1. KeyGen(a): Choose a polynomial  $r_1 \in R_q$  from  $\mathcal{X}_{\sigma}$ , choose another polynomial  $r_2$  with binary coefficients and then compute  $p = r_1 a \cdot r_2 \in R_q$ . The NTT is performed on the three polynomials a, p and  $r_2$  to generate  $\tilde{a}, \tilde{p}$  and  $\tilde{r_2}$ . The public key is  $(\tilde{a}, \tilde{p})$  and the private key is  $\tilde{r_2}$ .
- 2.  $Enc(\tilde{a}, \tilde{p}, m)$ : The message m is first encoded to  $\bar{m} \in R_q$ . Three polynomials  $e_1, e_2, e_3 \in R_q$  are sampled from  $\mathcal{X}_{\sigma}$ . The ciphertext is then computed as:

$$\tilde{e_1} \leftarrow NTT(e_1); \quad \tilde{e_2} \leftarrow NTT(e_2) (\tilde{c_1}, \tilde{c_2}) \leftarrow \left( \tilde{a} * \tilde{e_1} + \tilde{e_2}; \ \tilde{p} * \tilde{e_1} + NTT(e_3 + \bar{m}) \right)$$

3.  $Dec(\tilde{c_1}, \tilde{c_2}, \tilde{r_2})$ : Compute m' as  $m' = INTT(\tilde{c_1} * \tilde{r_2} + \tilde{c_2}) \in R_q$  and recover the original message m from m' using a decoder.

The scheme requires both encryption and decryption to use a common primitive root of unity.



Fig. 3. Ring-LWE Cryptoprocessor

#### 6.1 Hardware Architecture for the Ring-LWE Encryption Scheme

Figure 3 shows a hardware architecture for the ring-LWE encryption system. The basic building blocks used in the architecture are: the memory file, the arithmetic unit, the discrete Gaussian sampler and the control-address generation unit. The arithmetic unit is the NTT-ALU that we described in the previous section. Here we briefly describe the memory file and the discrete Gaussian sampler.

The Memory File is designed to support the maximum memory requirement that occurs during the encryption of the message. Six memory blocks  $M_0$  to  $M_5$  are available in the memory file and are used to store  $\bar{a}$ ,  $\bar{p}$ ,  $e_1$ ,  $e_2$ ,  $e_3$  and  $\bar{m}$  respectively. The memory blocks have width  $2\lceil \log q \rceil$  bits and depth n/2. All six memory blocks share a common read and a write address and have a common data-input line, while their data-outputs are selected through a multiplexer. Any of the memory blocks in the memory file can be chosen for read and write operation. Due to the common addressing of the memory blocks, the memory file supports one read and one write operation in every cycle.

The Discrete Gaussian Sampler is based on the compact Knuth-Yao sampler [13] architecture proposed in [28] and have sufficiently large precision and tail-bound to satisfy a maximum statistical distance of  $2^{-90}$  to a true discrete Gaussian distribution for both s = 11.32 and s = 12.18. Though the sampler in [28] is very compact it is also quite slow due to sequential scanning of the probability bits. We improve the cycle requirement of the sampler using two look-up tables. The first lookup table directly maps eight parallel random bits into a sample value or an intermediate distance in the 8th column of the probability matrix [28]. A successful look-up operation returns a sample and the sign of the sample is determined by the 9th random bit. If the first look-up operation fails, then another lookup is performed in the next 5 columns to get a sample value or an intermediate distance in the 13th column of the probability matrix. When the second lookup operation fails (probability<0.0016) then bit-scan based Knuth-Yao

random walk [28] is started with the initial distance obtained from the second lookup operation.

**The Cycle Count** for the encryption and decryption operations can be minimized in the following way. During the encryption operation, first the three error polynomials  $e_1$ ,  $e_2$  and  $e_3$  are generated by invoking the discrete Gaussian sampler 3n times. Next the encoded message  $\bar{m}$  is added to  $e_3$  and then three consecutive forward NTT operations are performed on  $e_1$ ,  $e_2$  and  $(e_3 + \bar{m})$ . Finally the ciphertext  $\tilde{c}_1$ ,  $\tilde{c}_2$  is obtained using two coefficient-wise multiplications followed by two polynomial additions and two rearrangement operations. The decryption operation requires one coefficient-wise multiplication, one polynomial addition and finally one inverse NTT operation.

During the encryption operation, 3n samples are generated to construct the three error polynomials. Our fast Knuth-Yao sampler architecture requires 805 and 1644 cycles for the dimensions 256 and 512 respectively on average to generate the three error polynomials. The polynomial addition and point-wise multiplication operations require n cycles each with a small overhead. The consecutive processing of I forward NTTs share a fixed computation cost  $fc_{fwd}$  and require in total  $fc_{fwd} + I \times \frac{n}{2} \log(n)$  cycles. Similarly I consecutive inverse NTTs are processed in  $fc_{inv} + I \times \frac{n}{2} \log(n) + I \times n$  cycles. One interesting point is that the fixed cost  $fc_{inv}$  is larger than  $fc_{fwd}$  as it includes the computation of  $\omega_{2n}^i/N$  (Section 2.3) for  $i = (0 \dots n-1)$ . This observation has been used to optimize the overall ring-LWE based encryption scheme in Section 6. The additional  $I \times n$  cycles during the inverse NTTs are required to multiply the coefficients by the scaling factors. The rearrangement of polynomial coefficients after an NTT operation requires less than n cycles. From the above cycle counts for each primitive operations, we see that the encryption and decryption operations require total  $fc_{fwd} + \frac{3}{2}n\log(n) + 10n$ and  $fc_{inv} + \frac{n}{2}\log(n) + 3n$  cycles respectively along with additional overhead. Our ring-LWE architecture has the fixed computation costs  $fc_{fwd} = 667$  and  $fc_{inv} = 1048$  cycles for n = 256; and  $fc_{fwd} = 1139$  and  $fc_{inv} = 1959$  cycles for n = 512.

## 7 Experimental Results

We have implemented the proposed ring-LWE cryptosystem on the Xilinx Virtex 6 FPGA for the parameter sets (n, q, s): (256,7681,11.32) and (512,12289,12.18). The area and performance results are obtained from the Xilinx ISE12.2 tool after place and route analysis and are shown in Table 1. In the table we also compare our results with other reported hardware implementations of the ring-LWE encryption scheme.

Our implementations are both fast and small thanks to the proposed computational optimizations and resource efficient design style. The cycle counts shown in the table do not include the cycles for data loading or reading operations. Our Knuth-Yao samplers have less than  $2^{-90}$  statistical distances from the corresponding true discrete Gaussian distributions and consume around 164 LUTs and have delay less than 2.5ns (with optimization goal for speed). Such a small delay makes the sampler suitable for integration in the pipelined ring-LWE processor under a single clock domain. We use nine parallel true random bit generators [8, 6] to generate the random bits for the sampler. The set of true random bit generators consumes 378 LUTs and 9 FFs.

The first hardware implementation of the ring-LWE encryption scheme in [9] uses a heavily parallel architecture to minimize the number of clock cycles for the NTT computation. Due to the many parallel computational blocks, the architecture is very large (0.29 million LUTs and 0.14 million FFs for n = 256) and does not even fit on

Implementation	Parameters	Device	LUTs/FFs/	Freq	$\mathbf{Cycles}/\mathbf{Time}(\mu s)$	
Algorithm			DSPs/BRAM18	$(\mathbf{MHz})$	Encryption	Decryption
Our RLWE	(256, 7681, 11.32)	V6LX75T	1349/860/1/2	313	6.3k/20.1	2.8k/9.1
Our RLWE	(512, 12289, 12.18)		1536/953/1/3	278	13.3k/47.9	5.8k/21
RLWE [23]	(256, 7681, 11.32)	V6LX75T	4549/3624/1/12	262	6.8k/26.2	4.4k/16.8
RLWE	(512, 12289, 12.18)	V6LX75T	5595/4760/1/14	251	13.7k/54.8	8.8k/35.4
RLWE-Enc[24]	(256, 4096, 8.35)	S6LX9	317/238/95/1	144	136k/946	-
RLWE-Dec			112/87/32/1	189	-	66k/351
ECC[3]	Binary-233	V5LX85T	18097/-/5644/0	156	1.9k/12.3	1.9k/12.3
NTRU[12]	NTRU-251	XCV1600E	27292/5160/14352/0	62.3	-/1.54	-/1.41

Table 1. Performance and Comparison

the largest FPGA of the Virtex 6 family. Performance results such as cycle count and frequency are not reported in their paper. The architecture uses a Gaussian distributed array for sampling of the error coefficients up to a tail-bound of  $\pm 2s$ .

The implementation in [23] is small and fast due to its resource-efficient design style. A high operating frequency is achieved using pipelines in the architecture. The architecture uses a ROM that keeps all the twiddle factors required during the NTT operation. This approach reduces the fixed computation  $\cot(fc)$  but consumes block RAM slices in FPGAs. Additionally, the parallel RAM blocks in the NTT processor result in a larger memory requirement compared to our design. The discrete Gaussian sampler is based on the inversion sampling method [5] and has a maximum statistical distance of  $2^{-22}$  to a true discrete Gaussian distribution. Since the inversion sampling requires many random bits to output a sample value, an AES core is used as a pseudorandom number generator. The AES core itself consumes an additional 803 LUTs and 341 FFs compared to our true random number generator. Another reason behind the larger area consumption of [23] compared to our architecture is due to the fact that the architecture supports different parameter sets at synthesis time. Our ring-LWE processor is also designed to achieve scalability for various parameter sets. In our architecture the control block remains the same; while only the data-width and the modular reduction block changes for different parameter sets. Hence our architecture is also configurable by generating the HDL codes for various parameter sets using a C program.

Although our architecture does not use a dedicated ROM for storing the twiddle factors, it still achieves slightly smaller cycle count and faster computation time compared to [23]. The encryption scheme in [23] computes one forward and two inverse NTTs; while our encryption scheme computes only forward NTTs and hence does not require the 4n cycles for the scaling operation. Additionally our negative convolution method is free from the precomputation that takes n cycles in [23]. Hence we save 5n cycles in total during the NTT operations in an encryption operation. Since the fixed computation cost  $fc_{fwd}$  is smaller than 5n, we gain in cycle count for the encryption operation. The decryption operation in our case is trivially faster than [23] as only one NTT is performed. We also reduce the area and memory requirement significantly compared to [9, 23]. This reduction is achieved by our resource-efficient design decisions such as 1) absence of a dedicated ROM for the twiddle factors, 2) an efficient RAM access and storage scheme, 3) use of one modular multiplier, 4) use of a smaller and faster (lowdelay) discrete Gaussian sampler, and finally 5) the resource sharing between different computations.

A very recent paper [24] proposes ring-LWE encryption and decryption architectures targeting small area at the cost of performance. The implementation uses a quadraticcomplexity multiplier instead of a complicated NTT based polynomial multiplier. Additionally the special modulus also saves some amount of area as the modular reduction is free of cost. However if we consider a similar quadratic-complexity multiplication based architecture in the dimension n = 512, then the cycle requirement will be nearly 40 times compared to our NTT-based ring-LWE processor. Our target was to use FPGA resources more efficiently without affecting the performance and to achieve similar speed as [23]. The paper [24] also designs a compact Bernoulli sampler that consumes 37 slices for the standard deviation 3.33 and is thus smaller in area compared to the Knuth-Yao sampler in [28]. The Bernoulli sampler requires on average 96 random bits and 144 cycles to output a sample. In the contrast the Knuth-Yao sampler [28] requires on average 5 random bits and 17 cycles per sample and is thus faster than the Bernoulli sampler. In this paper we have reduced the area consumption of the Knuth-Yao sampler [28] by reducing the width of the ROM and the scan-register from 32 bits to 12 bits and by simplifying the control unit. These area optimizations do not affect the cycle requirement of the sampler, but result in an area of only 32 slices for the overall sampler. The area optimized Knuth-Yao sampler is both smaller and faster compared to the Bernoulli sampler in [24].

We also compare our results with other cryptosystems such as ECC and NTRU. The ECC processor [3] over the NIST recommended binary field  $GF(2^{233})$  requires 12.3  $\mu s$  to compute one scalar multiplication and is faster than our ring-LWE processor. However the ECC processor is designed to achieve high speed and hence consumes very large area compared to our ring-LWE processor. The NTRU scheme [12] is much faster than our ring-LWE processor due to its less complicated arithmetic. However the parameters chosen for the implementation in [12] have security around 64 bits [11]. Though secure parameter sets for the NTRU based encryption have been proposed in [10], no hardware implementation for the secure parameter sets is available in the literature.

## 8 Conclusion

This paper proposed several optimizations for implementing a ring-LWE based encryption system. The first set of optimizations improved the NTT by reducing the computation cost of the twiddle factors, avoiding the pre-computation during the forward NTT, and deriving an efficient memory access scheme that increases the utilization of the arithmetic components and the memory blocks. A further optimization reduced the number of NTTs required in the encryption scheme from five to four. The proposed optimizations are implemented in an efficient cryptoprocessor for the ring-LWE encryption system that not only is three times smaller in area and memory than any other reported implementations, but also even faster. These features make the architecture suitable for resource constrained platforms. Furthermore, the paper investigated architectural acceleration to meet the high speed requirement for real-time applications and proposes an optimal pipeline strategy that results in a very fast computation time whilst using minimum area and memory. Although the paper focuses on implementation of the ring-LWE based encryption system, we finally remark that the proposed optimization techniques for the NTT computation are applicable for other lattice based cryptosystems where similar polynomial multiplications are performed.

### Acknowledgment

This work was supported by the Research Council KU Leuven: TENSE (GOA/11/007), by iMinds, by the European Union Seventh Framework Programme (FP7/2007-2013) under grant agreement n. 609611 (PRACTICE), by the Flemish Government, FWO G.0550.12N, by the Hercules Foundation AKUL/11/19. The first author is supported by the Erasmus Mundus PhD Scholarship. We are thankful to Bohan Yang and Vladimir Rozic for useful technical discussions related to memory utilisation and random number generation.

#### References

- 1. A. Aysu, C. Patterson, and P. Schaumont. Low-cost and Area-efficient FPGA Implementations of Lattice-based Cryptography. In *HOST*, pages 81–86. IEEE, 2013.
- D. Bernstein. Fast Multiplication and its Applications. Algorithmic Number Theory, 44:325–384, 2008.
- C. Rebeiro, S. Sinha Roy, and D. Mukhopadhyay. Pushing the Limits of High-Speed GF(2<sup>m</sup>) Elliptic Curve Scalar Multiplication on FPGAs. In Cryptographic Hardware and Embedded Systems, volume 7428 of Lecture Notes in Computer Science, pages 494–511. Springer Berlin Heidelberg, 2012.
- 4. T. Cormen, C. Leiserson, and R. Rivest. *Introduction To Algorithms*. http://staff.ustc.edu.cn/~csli/graduate/algorithms/book6/toc.htm.
- 5. L. Devroye. Non-Uniform Random Variate Generation. Springer-Verlag, New York, 1986.
- M. Dichtl and J. D. Golic. High-Speed True Random Number Generation with Logic Gates Only. In Cryptographic Hardware and Embedded Systems - CHES 2007, volume 4727 of LNCS, pages 45–62. Springer Berlin, 2007.
- T. Frederiksen. A Practical Implementation of Regev's LWE-based Cryptosystem. In http://daimi.au.dk/ jot2re/lwe/resources/, 2010.
- J. D. Golic. New Methods for Digital Generation and Postprocessing of Random Data. IEEE Transactions on Computers, 55(10):1217–1229, 2006.
- N. Göttert, T. Feller, M. Schneider, J. Buchmann, and S. Huss. On the Design of Hardware Building Blocks for Modern Lattice-Based Encryption Schemes. In *Cryptographic Hardware* and *Embedded Systems CHES 2012*, volume 7428 of *LNCS*, pages 512–529. Springer Berlin, 2012.
- P. Hirschhorn, J. Hoffstein, N. Howgrave-graham, and W. Whyte. Choosing NTRUEncrypt parameters in light of combined lattice reduction and MITM approaches. In *In Proc. ACNS* 2009, LNCS 5536, pages 437–455. Springer-Verlag, 2009.
- N. Howgrave Graham. A Hybrid Lattice-Reduction and Meet-in-the-Middle Attack Against NTRU. In Advances in Cryptology - CRYPTO 2007, volume 4622 of Lecture Notes in Computer Science, pages 150–169. Springer Berlin Heidelberg, 2007.
- A. Kamal and A. Youssef. An FPGA implementation of the NTRUEncrypt cryptosystem. In *Microelectronics (ICM)*, 2009 International Conference on, pages 209–212, Dec 2009.
- D. E. Knuth and A. C. Yao. The Complexity of Non-Uniform Random Number Generation. Algorithms and Complexity, pages 357–428, 1976.
- T. Lepoint and M. Naehrig. A Comparison of the Homomorphic Encryption Schemes FV and YASHE. IACR Cryptology ePrint Archive, 2014:62, 2014.
- R. Lindner and C. Peikert. Better Key Sizes (and Attacks) for LWE-based Encryption. CT-RSA 2011, pages 319–339, 2011.
- M. Liu and P. Q. Nguyen. Solving BDD by Enumeration: An Update. In Proceedings of the 13th International Conference on Topics in Cryptology, CT-RSA'13, pages 293–309, Berlin, Heidelberg, 2013. Springer-Verlag.

- V. Lyubashevsky, C. Peikert, and O. Regev. On Ideal Lattices and Learning with Errors over Rings. In Advances in Cryptology EUROCRYPT 2010, volume 6110 of Lecture Notes in Computer Science, pages 1–23. Springer Berlin Heidelberg, 2010.
- Y. Ma and L. Wanhammar. A Hardware Efficient Control of Memory Addressing for High-performance FFT Processors. Signal Processing, IEEE Transactions on, 48(3):917– 921, Mar 2000.
- 19. D. Micciancio. Lattices in Cryptography and Cryptanalysis. 2002.
- P. Q. Nguyen and J. Stern. The Two Faces of Lattices in Cryptology. In Cryptography and Lattices, International Conference (CaLC 2001), volume 2146 of LNCS, pages 146–180. Springer-Verlag, Berlin, 2001.
- 21. J. Pollard. The Fast Fourier Transform in a Finite Field. *Mathematics of Computation*, 25:365374, 1971.
- 22. T. Pöppelmann and T. Güneysu. Towards Efficient Arithmetic for Lattice-Based Cryptography on Reconfigurable Hardware. In A. Hevia and G. Neven, editors, *Progress in Cryptology LATINCRYPT 2012*, volume 7533 of *LNCS*, pages 139–158. Springer Berlin, 2012.
- T. Pöppelmann and T. Güneysu. Towards Practical Lattice-Based Public-Key Encryption on Reconfigurable Hardware. In *Selected Areas in Cryptography – SAC 2013*, Lecture Notes in Computer Science, pages 68–85. Springer Berlin Heidelberg, 2014.
- 24. T. Pöppelmann and T. Güneysu. Area Optimization of Lightweight Lattice-Based Encryption on Reconfigurable Hardware. In *Proc. of the IEEE International Symposium on Circuits and Systems (ISCAS-14)*, 2014, Preprint.
- O. Regev. Quantum Computation and Lattice Problems. SIAM J. Comput., 33(3):738–760, Mar. 2004.
- O. Regev. On Lattices, Learning with Errors, Random Linear Codes, and Cryptography. In Proceedings of the thirty-seventh annual ACM symposium on Theory of computing, STOC '05, pages 84–93, New York, NY, USA, 2005. ACM.
- O. Regev. Lattice-Based Cryptography. In C. Dwork, editor, Advances in Cryptology -CRYPTO 2006, volume 4117 of LNCS, pages 131–141. Springer Berlin, 2006.
- S. Sinha Roy, F. Vercauteren, and I. Verbauwhede. High Precision Discrete Gaussian Sampling on FPGAs. In *Selected Areas in Cryptography – SAC 2013*, Lecture Notes in Computer Science, pages 383–401. Springer Berlin Heidelberg, 2014.
- J. van de Pol and N. P. Smart. Estimating key sizes for high dimensional lattice-based systems. In M. Stam, editor, *IMA Int. Conf.*, volume 8308 of *Lecture Notes in Computer Science*, pages 290–303. Springer, 2013.

# Appendix A

Table 2 shows the memory contents during the execution of Algorithm 2 for n = 16. The column-heading represents (m, j, k) during the iterations. The end loop in line 19 of Algorithm 2 for m = 16 performs no swap and is shown in the table using  $\star$  symbol.

# Appendix B

Our ring-LWE cryptoprocessor has one instruction-register, one iteration-register, one read-memory-index-queue and one write-memory-index-queue (Figure 4). The read and write memory-index-queues are loaded with the memory indexes. Since our ring-LWE cryptoprocessor has six memory blocks M0 to M5, the indexes are in the range 0 to 5. The instruction is stored in the *Instruction* register and the number (I) of consecutive NTT operations is kept in the *Iteration* register. The following instructions are supported by the processor.

Address	Initial	(2,0,0)	(2,0,6)	(4,0,0)	(4,0,4)	(4,1,4)	(8,3,0)	(16,7,0)*
0	$A_1 A_0$	$A_2 A_0$	$A_2 A_0$	$A_4 A_0$	$A_4 A_0$	$A_4 A_0$	$A_8 A_0$	$A_8 A_0$
1	$A_3 A_2$	$A_3 A_1$	$A_3 A_1$			$A_5 A_1$	$A_9 A_1$	$A_9 A_1$
2	$A_5 A_4$		$A_6 A_4$	$A_6 A_2$	$A_6 A_2$	$A_6 A_2$	$A_{10} A_2$	$A_{10} A_2$
3	$A_7 A_6$		$A_7 A_5$			$A_7 A_3$	$A_{11} A_3$	$A_{11} A_3$
4	$A_9 A_8$		$A_{10} A_{8}$		$A_{12} A_8$	$A_{12} A_8$	$A_{12} A_4$	$A_{12} A_4$
5	$A_{11} A_{10}$		$A_{11} A_9$			$A_{13} A_{9}$	$A_{13} A_5$	$A_{13} A_5$
6	$A_{13} A_{12}$		$A_{14} A_{12}$		$A_{14} A_{10}$	$A_{14} A_{10}$	$A_{14} A_6$	$A_{14} A_6$
7	$A_{15} A_{14}$		$A_{15} A_{13}$			$A_{15} A_{11}$	$A_{15} A_7$	$A_{15} A_7$

Table 2. Memory content during the steps in a 16-point NTT



Fig. 4. Instruction Execution Hardware

- 1. LOAD : A memory block indexed by WtQ0 is loaded with n coefficients. Since two coefficients are processed in a cycle, the instruction takes  $n/2 + \epsilon$  cycles.
- 2. ENCODE-LOAD : A memory block indexed by WtQ0 is loaded with an encoded message. The input message bits are first encoded using the encoder and then loaded in the memory block as proper coefficient-pairs. This instruction requires  $n+\epsilon$  cycles.
- 3. GAUSSIAN-LOAD : A memory block indexed by WtQ0 is loaded with n samples. The cycle count for this operation depends on the standard deviation and n.
- 4. FNTT/INTT : Is used to perform inplace forward or inverse NTT. The number of consecutive NTTs is stored in the iteration-register and the indexes of the memory blocks are kept in the read-memory-index-queue
- 5. ADD/CMULT : Two memory blocks indexed by RdQ0 and RdQ1 are added or coefficient-wise multiplied. The result is stored in the memory block indexed by WtQ0. These two instructions require  $n + \epsilon$  cycles.
- 6. REARRANGE : Performs rearrangement of coefficient pairs in a memory block indexed by RdQ0. This instruction requires less than n cycles.
- 7. READ : The contents of a memory block indexed by RdQ0 are read. This instruction requires  $n/2 + \epsilon$  cycles.