

HPKA: A High-Performance CRYSTALS-Kyber Accelerator Exploring Efficient Pipelining

Ziying Ni, *Student Member, IEEE*, Ayesha Khalid, *Senior Member, IEEE*, Dur-e-Shahwar Kundi, Máire O'Neill, *Senior Member, IEEE*, and Weiqiang Liu, *Senior Member, IEEE*

Abstract—CRYSTALS-Kyber (Kyber) was recently chosen as the first quantum resistant Key Encapsulation Mechanism (KEM) scheme for standardisation, after three rounds of the National Institute of Standards and Technology (NIST) initiated PQC competition which began in 2016 and search of the best quantum resistant KEMs and digital signatures. Kyber is based on the Module-Learning with Errors (M-LWE) class of Lattice-based Cryptography, that is known to manifest efficiently on FPGAs. This work explores several architectural optimizations and proposes a high-performance and area-time (AT) product efficient hardware accelerator for Kyber. The proposed architectural optimizations include inter-module and intra-module pipelining, that are designed and balanced via FIFO based buffering to ensure maximum parallelisation. The implementation results show that compared to state-of-the-art designs, the proposed architecture delivers 25-51% speedups for Kyber's three different security levels on Artix-7 and Zynq UltraScale+ devices, and a 50-75% reduction in DSPs at comparable security level. Consequently, the proposed design achieves higher AT product efficiencies of 19-33%.

Index Terms—Post-quantum Cryptography (PQC), Lattice-based Cryptography (LBC), Module-Learning with Errors (M-LWE), CRYSTALS-Kyber, Hardware Accelerator.



1 INTRODUCTION

The advent of quantum computers threatens the security of all existing cryptosystems. A quantum algorithm, called Shor's algorithm [1], is capable of completely breaking all currently deployed Public-key Cryptography (PKC), including RSA [2] and Elliptic Curve Cryptography (ECC) [3]. In addition, Grover's search quantum algorithm [4] reduces the complexity of the search space of a brute force attack on symmetric-key encryption schemes (e.g., AES [5]) and hashing (e.g., SHA-3 [6]) to half. The National Institute of Standards and Technology (NIST) announced a formal global call in 2016, to standardize new Post-quantum Cryptography (PQC) based Public-key Encryption (PKE) and digital signature schemes [7]. In 2017, 69 proposals were selected in Round 1 of the NIST PQC and four candidates and five alternatives were shortlisted in Round 3 in July 2020. Three out of these four candidate algorithms were lattice-based cryptographic (LBC) schemes, namely CRYSTALS-Kyber [8], SABER [9], and NTRU [10]. NIST intended to standardize no more than one of these lattice-based Public-

key Encryption and Key-encapsulation algorithms and on July 5, 2022, NIST announced the first group of winners from its six-year competition. CRYSTALS-Kyber was announced as the first PQC algorithm to be standardized as a Key-encapsulation Mechanism (KEM) [11].

CRYSTALS-Kyber (hereafter called Kyber) [8] KEM is based on the module learning with errors (M-LWE) problem, which is a lattice based hard problem. M-LWE is an "algebraic" LWE with a tight formal mathematical security reduction of the ring-LWE (R-LWE) problem [12]. Schemes based on the M-LWE problem have a more elaborate algebraic structure and consequently, higher security than R-LWE schemes while achieving higher performance than LWE schemes. In the M-LWE scheme, a parameter k is introduced to restrict the dimensions of the public-key matrix A , but all elements of the matrix are on the ring $\mathbb{Z}_q[x]/(x^n + 1)$. Unlike the lattice-based scheme SABER [9], the polynomial operations in Kyber can be computed using the Number Theoretic Transform (NTT), which allows Kyber to gain a high throughput performance. For round 2 of the NIST PQC submission, the Kyber team adopted a technique to reduce the parameter q of Kyber from 7,681 to 3,329, further reducing the complexity of the modular reduction and area resources.

High speed implementations of Kyber, as the first standardized KEM scheme, are desired for high-performance applications. However, the current Kyber architecture's reuse of the main computational modules, influenced by the storage resources for large amounts of polynomial data, limits further speedups. To the best of the authors' knowledge, this is the first work to present an inter- and intra-module pipelined hardware accelerator for the IND-Chosen-ciphertext Attack-2 (IND-CCA2) secure Kyber KEM scheme that speeds up the operation time. The accelerator comprises Key Generation, Encapsulation, and Decapsu-

Manuscript received X XX, 2023; revised X XX, 2023. This work is supported by the National Natural Science Foundation of China (62022041) and the Fundamental Research Funds for the Central Universities (NP2022103), and grants from the Engineering and Physical Sciences Research Council (EPSRC) Quantum Communications Hub (EP/T001011/1).

Z. Ni is with the Centre for Secure Information Technologies (CSIT), Queens University Belfast, UK and College of Electronic and Information Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing. E-mail: zni03@qub.ac.uk.

A. Khalid and M. O'Neill are with the Centre for Secure Information Technologies (CSIT), Queens University Belfast, UK. E-mail: {a.khalid}@qub.ac.uk, m.oneill@ecit.qub.ac.uk.

D.-S. Kundi is with the PQShield, E-mail: {dur-e-shahwar.kundi}@pqshield.com.

W. Liu is with College of Electronic and Information Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing, China. E-mail: liuweiqiang@nuaa.edu.cn.

lation modules for the three NIST specific security levels 1/3/5. The proposed architecture has high speed and high Area-Time (AT) product efficiency. The major contributions of this work are summarized as follows:

Our Kyber accelerator *aggressively exploits architectural parallelisation* via optimal inter-module and intra-module pipelines. To balance the pipeline, buffering is provided at the interface of several modules. The computation order of the modules is rearranged to facilitate an optimal usage of pipeline.

A fully pipelined *Radix-2-Multipath Delay Commutator (MDC)-NTT* core is presented that multiplexes the resources for both the NTT and inverse-NTT (INTT) computations. By using different delay units, the bit reversal operation in the NTT/INTT calculation is completely eliminated. Due to pipelining, a single NTT/INTT computation requires only 128 clock cycles, once the pipeline is full.

Resource utilization is reduced in terms of DSPs and BRAMs by several strategies. The hardware for NTT/INTT is shared. To balance the pipeline, buffering is done via FIFOs (restricting first in first out data access) instead of BRAMs that allow more flexible access but are more expensive in resource consumption. The data input order of the FIFOs for the NTT/INTT module has been organized so as not to require the use of any BRAM in our proposed architecture.

Our proposed Kyber accelerator surpasses all previously reported FPGA based implementations with comparable security levels in terms of execution time and design efficiency (i.e., Area-Time (AT) product). Compared to the state-of-the-art design, the proposed architecture has a speedup of 1.25-1.44 at the three security levels on an Artix-7 and 1.26-1.51 on a Zynq-UltraSale+. In terms of hardware efficiency, the proposed architecture improves the AT efficiency by 33.3%, 26.9%, and 18.9% for the three different security levels, respectively.

Intra module architectural pipelining in hardware designs is not a novel acceleration technique. However, considering simultaneously the intra and the inter module parallelization maximization and consequently balancing and orchestrating the whole data buffering and data flow based on that is not considered in previous research for ultra high speed PQC hardware (especially Kyber).

This paper is organized as follows. Section II introduces the Kyber protocol and NTT, and Section III presents the proposed overall architecture including various modules. A fast pipelined scheduling scheme and memory approaches for the Kyber hardware architecture are presented in Section IV. Implementation results and a comparison with previous designs are provided in Section V, and Section VI concludes our work.

2 PRELIMINARIES

In this section, the Kyber KEM is explained, describing the main NTT construct in Kyber.

2.1 Kyber.v3 (NIST PQC Round 3)

Kyber KEM is the first lattice based PQC algorithm chosen by NIST for standardisation. The relative balance between performance and security can be directly adjusted by tweaking the size of the matrix k ; the choice of k varies between 2, 3, or 4 for security levels 1 (Kyber512), 3 (Kyber768) and 5 (Kyber1024), respectively. The noise parameter η is adjusted according to the security level. The IND-CCA2 secure Kyber KEM, submitted to NIST PQC Round 3 referred to as Kyber.CCA consists of three main steps: key generation (Kyber.CCA.KeyGen), key encapsulation (Kyber.CCA.Enc), and key decapsulation (Kyber.CCA.Dec). The prime used in Kyber, p is changed from 7, 681 to 3, 329, which enables the polynomial multiplication in Kyber to be accelerated using NTT. The Kyber.CCA implementation is built on top of the Kyber.CPA, using the Fujisaki-Okamoto transform [13]. Kyber.CPA comprises three components: key generation (Kyber.CPA.KeyGen), encryption (Kyber.CPA.Enc), and decryption (Kyber.CPA.Dec). A functional description of the three constituent functions of Kyber.CPA are described as follows, and for more details of Kyber, the reader is kindly referred to [8].

Kyber.CPA.KeyGen(): Before a communication can be established between Alice and Bob, the Key generation function generates a public key pk and secret key sk . First, a matrix A is generated directly in the NTT domain by uniform sampling; vectors s and e are generated by binomial distribution sampling. Then, the polynomial $\hat{t} = \hat{A} \text{NTT}(s) + \text{NTT}(e)$ is computed, and the key pair is generated by $sk = (\text{NTT}(s))$ and $pk = (\hat{t}, \rho)$, where ρ is from a random number.

Kyber.CPA.Enc(pk, m, r): First, the matrix A^T is generated in the NTT domain by uniform sampling, other vectors generated are r, e_1 and polynomial e_2 by seed r . The vector u and v are then computed by $u = \text{NTT}^{-1}(\hat{A}^T \hat{r}) + e_1$ and $v = \text{NTT}^{-1}(\hat{t}^T \hat{r}) + e_2 + \text{Decompress}_q(\text{Decode}_1(m), 1)$. Next, u and v are compressed and encoded to produce c_1 and c_2 . Finally, the Ciphertext $c = (c_1 || c_2)$ is returned.

Kyber.CPA.Dec(sk, c): The vectors u and v from the input c are decompressed. Then the message m° is calculated from $m^\circ = \text{Encode}_1(\text{Compress}_q(v \text{NTT}^{-1}(\hat{S}^T \text{NTT}(u), 1)))$. The m° is returned as the output of the function.

Algo. 1, 2, and 3 show the Kyber.CCA.KeyGen, Kyber.CCA.Enc and Kyber.CCA.Dec functions in Kyber.CCA, respectively. Kyber.CCA.KeyGen() calls Kyber.CPA.KeyGen() function to generate public key pk and sk° . The secret key sk is then generated by $sk = (sk^\circ || \text{H}(pk) || z)$, where H represents SHA3-256. Kyber.CCA.Enc() receives the public key pk and calls the Kyber.CPA.Enc(pk, m, r) function to generate the ciphertext c . The shared secret ss is then generated by SHAKE-256. The output includes the ciphertext c and shared secret ss . Kyber.CCA.Dec() takes the ciphertext c and the secret key sk as inputs. First the Kyber.CPA.Enc(sk, c) function is called to recover the message m , and then Kyber.CPA.Dec(pk, m°, r°) computes the new ciphertext c° . If the comparison between c and c° returns a success, the encryption is successful.

Algorithm 1 Kyber.CCA.KeyGen()

-
- 1: **Output:** Public key pk , Secret key sk
 - 2: $z = B^{32}$
 - 3: $(pk, sk^0) := \text{Kyber.CPA.KeyGen}()$
 - 4: $sk := (sk^0 || pk || H(pk) || z)$
 - 5: **return** (pk, sk)
-

Algorithm 2 Kyber.CCA.Enc(pk)

-
- 1: **Input:** Public key pk
 - 2: **Output:** Ciphertext c , Shared key ss
 - 3: $m = B^{32}$
 - 4: $m = H(m)$
 - 5: $(\bar{K}, r) := G(m || H(pk))$
 - 6: $c := \text{Kyber.CPA.Enc}(pk, m, r)$
 - 7: $ss := \text{KDF}(\bar{K} || H(c))$
 - 8: **return** (c, ss)
-

Algorithm 3 Kyber.CCA.Dec(c, sk)

-
- 1: **Input:** Ciphertext c , Secret key sk
 - 2: **Output:** Shared key ss
 - 3: $m^0 := \text{Kyber.CPA.Dec}(sk, c)$
 - 4: $(\bar{K}^0, r^0) := G(m^0 || h)$
 - 5: $c^0 := \text{Kyber.CPA.Enc}(pk, m^0, r^0)$
 - 6: **if** $c = c^0$ **then**
 - 7: $ss := \text{KDF}(\bar{K}^0 || H(c))$
 - 8: **else**
 - 9: $ss := \text{KDF}(z || H(c))$
 - 10: **end if**
 - 11: **return** ss
-

2.2 NTT in Kyber

The NTT algorithm is derived from the Fast Fourier Transform (FFT) algorithm. Compared to standard schoolbook polynomial multiplication, the complexity of the NTT algorithm is reduced from $O(n^2)$ to $O(n \log n)$. The choice of modulus in the construction of Kyber satisfies the modulus restriction for the NTT, and polynomial multiplication calculations in Kyber can be accelerated using NTT. Table 1 shows the number of times the NTT/INTT (single NTT/INTT for 256 points) is needed during the key generation, encapsulation and decapsulation functions in Kyber, under three different security levels.

For NTT transformations, polynomials are expressed in terms of a vector of coefficients, e.g., the polynomial $a(x) = a_0 + a_1x + \dots + a_{n-1}x^{n-1} + a_nx^n$ is represented as a set of n points $a(x_i) = f(x_0, y_0), (x_1, y_1), \dots, g$. The NTT computation can be substantially improved when using n special points, i.e. n powers of the twiddle factor w .

TABLE 1

Number of Calls for The NTT/INTT Operation in The Key Generation, Encapsulation and Decapsulation, For The Three Security Levels.

Functions	Number of NTT calls Security Levels: 1/3/5	Number of INTT calls Security Levels: 1/3/5
Key Generation	4/6/8	0/0/0
Encapsulation	2/3/4	3/4/5
Decapsulation	4/6/8	4/5/6

Given a polynomial with n elements, k is an integer ranging from 0 to $n-1$, w is the square of ψ , where ψ is the primitive root of unity in NTT, the NTT and INTT transformations are shown as follows:

$$\hat{a}_m = \sum_{k=0}^{n-1} a_k \psi^{(2m+1)k} = \sum_{k=0}^{n-1} (a_k \psi^k) w^{mk} \text{ mod } q \quad (1)$$

$$a_k = \frac{1}{n} \sum_{m=0}^{n-1} \hat{a}_m \psi^{-(2m+1)k} = \psi^{-k} \frac{1}{n} \sum_{m=0}^{n-1} \hat{a}_m w^{-mk} \text{ mod } q \quad (2)$$

There are some differences between the NTT defined in Kyber [8] and the classical NTT. In Kyber, the base field \mathbb{Z}_q contains the primitive 256^{th} root of unity but not the primitive 512^{th} root of unity. Let $\zeta = 17$ be the first primitive 256^{th} root of unity modulo q . Therefore, the defining polynomial $x^{256} + 1$ of R factors into 128 polynomials of degree 2 modulo q . NTT in Kyber is defined in Eq. (3), br_7 in Eq. (3) means bit reversal of the unsigned 7-bit integer i and $k \geq \text{br}_7(i)$. In addition, when performing NTT calculations, these two parity polynomials are calculated independently.

$$\hat{f}_{2i+k} = \sum_{j=0}^{127} f_{2i+k} \zeta^{(2\text{br}_7(i)+1)j} \quad (3)$$

The polynomial multiplication $\text{NTT}(f) \cdot \text{NTT}(g) = \hat{h} \text{ mod } x^2 - \zeta^{2\text{br}_7(i)+1}$, where \hat{f} , \hat{g} , and \hat{h} are polynomials in NTT representation, can be expressed as:

$$\begin{aligned} \hat{h}_{2i} + \hat{h}_{2i+1}x &= \hat{f}_{2i}\hat{g}_{2i} + \zeta^{2\text{br}_7(i)+1} \hat{f}_{2i+1}\hat{g}_{2i+1} \\ &+ x(\hat{f}_{2i}\hat{f}_{2i+1} + \hat{g}_{2i}\hat{g}_{2i+1}) \end{aligned} \quad (4)$$

2.3 Related Work

Hardware based Kyber KEM accelerator designs have primarily focused on the optimizations of its most computationally intensive constituent component, i.e., the polynomial multiplication module (implemented via the NTT module) [14]–[18]. Several recent works on Kyber focus on improving the NTT memory access and modulo multiplication units. In 2020, Chen *et al.* proposed a pipelined processor for the vector of polynomials using two-column sequential storage and bit-inverted addressing access for Kyber ($p = 7, 681$) [19]. Zhang *et al.* in 2021 proposed an effective NTT structure for the prime in Round 2 ($p = 3, 329$) [15]. In the same year, Bisheh-Niasar *et al.* proposed the K2-RED module algorithm based on the prime $p = 3, 329$ of Kyber, and using four parallel-computing NTT butterfly units to achieve a high-speed polynomial multiplication accelerator [18]. Yaman *et al.* implemented different architectures of multiplication units by increasing the number of butterfly units in the NTT based on a unified butterfly structure [16], with lightweight, balanced, and high-performance hardware architectures, using 1, 4, and 16 parallel butterfly units, respectively. Resistance to attacks on the NTT has also received attention, with some work focusing on the error detection capabilities of the NTT [20], [21]. These NTT-based error detection architectures are compact and are suitable for deeply embedded NTT architectures.

Several hardware-only complete Kyber implementations have also been reported in the literature [18], [22]–[25]. The

first full hardware implementation of Kyber was presented by Huang *et al.* in 2020 [22]. In 2021, Xing *et al.* proposed a low-cost, high-performance Kyber processor on the Artix-7 platform [23]. This architecture utilized a predefined control order table with short control codes for the scheduling of various NTT processes and used different sized FIFOs for data transmission/ reception to achieve good throughput performance with limited computational resources. Bisheh-Niasar *et al.* proposed a polynomial multiplier for Kyber in 2021, using a 2x2 reconfigurable butterfly cell with pure combinational logic referred to as K2-RED modulo subtraction cells [18]. The same reconfigurable butterfly module was used in an instruction set processor for Kyber [24], whose overall operating frequency was limited. In addition, Dang *et al.* implemented three different lattice-based designs of NIST Round 3 candidates (Kyber, SABER, NTRU) in hardware [25] and achieved a 52.5%, 65.7%, 76.2% improvement in speed of Kyber at the three different security levels compared to earlier high performance Kyber implementations [23]. In addition, some Kyber accelerator designs implemented on other devices are also worthy of reference. In 2020, Viet Ba Dang *et al.* summarized the advantages of SW/HW joint designs of the candidates in the second round of NIST PQC [26]. Software-hardware co-design offers a rapid deployment of cryptographic algorithms and can result in faster performance compared to pure software design. However, the hardware only accelerators can optimize each component within the accelerator, generally achieving superior throughput performance. In 2021, Bisheh *et al.* optimized the Kyber accelerator and NTT unit on an ASIC, which greatly reduced the calculation speed of the Kyber accelerator [27]. In the same year, the Kyber512 accelerator

designed by Sanal *et al.* on ARM Cortex-A processor carried out several optimizations for NTT, noise sampling, improve the speed of the accelerator [28].

HPKA (the proposed high-performance Kyber accelerator) is designed to improve the AT-efficiency of CRYSTALS-Kyber FPGA based hardware implementations. While countermeasures against side channel analysis (SCA) attacks, more generally including fault attacks and power attacks, are not currently incorporated in HPKA, all of the components in our design, including of both the server-side and client-side for all the three security levels, are constant time in execution, and are consequently resistant against the timing attacks.

3 THE PROPOSED KYBER HARDWARE ACCELERATOR

The proposed Kyber accelerator comprises server and client side implementations. The server-side accelerator includes the key generation and decapsulation functions while the client-side implementation consists of a sub-set of components, performing only the key encapsulation function. An overall design of the server-side Kyber accelerator is shown in Fig. 1. Our primary goal is to use pipelining to the maximum extent possible to reduce latency in all three functions. The architecture is only multiplexed when considering that separate modules would significantly reduce the hardware efficiency of the accelerator, e.g. NTT/INTT units. The accelerator comprises a controller, computational units and storage units. The storage unit consists of multiple FIFOs and ROMs (in the NTT/INTT and polynomial multiplication unit). In Fig. 1, the black, green, and purple colors indicate the data flow for Kyber.CPA.KeyGen, Kyber.CPA.Dec and Kyber.CPA.Enc and

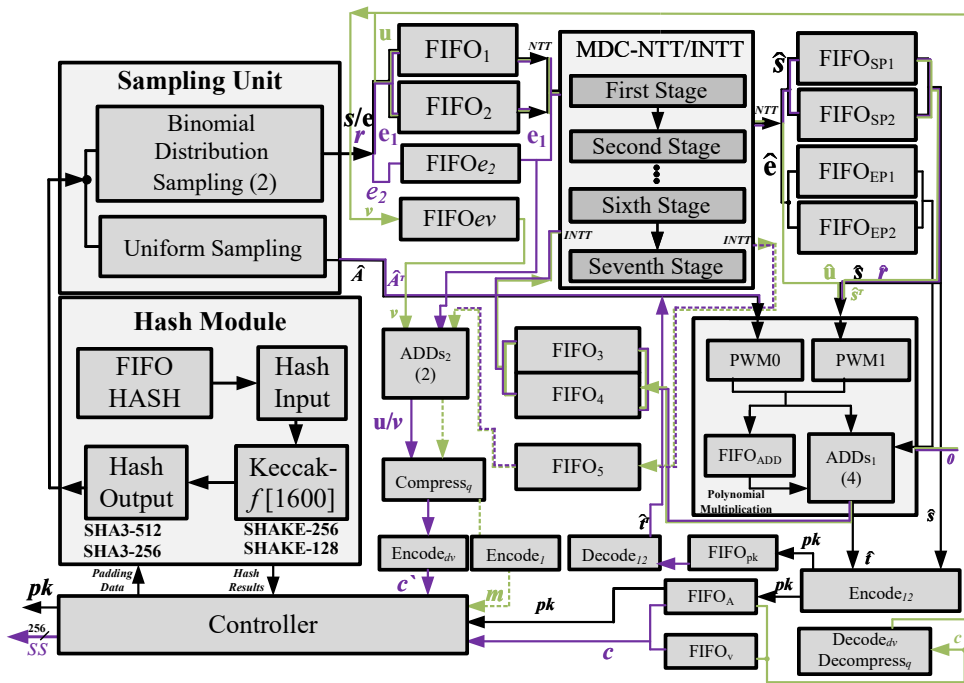


Fig. 1. The overall server-side Kyber768 architecture. The black, green, and purple colored arrows indicate the data flow for Kyber.CPA.KeyGen, Kyber.CPA.Dec and Kyber.CPA.Enc functions, respectively. The dashed line in the figure is the data flow after the solid line of the same color is complete.

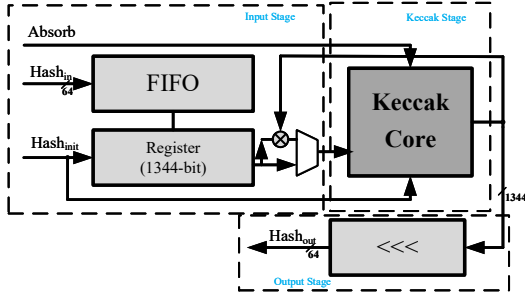


Fig. 2. Hash module comprising of input/ output stages and a Keccak core.

Kyber.CPA.Enc functions, respectively. The dashed lines indicate the data flow through INTT, which is computed after the solid line of the same color.

In our HPKA, four parallel data blocks (totalling 48-bits for 12-bit blocks) are simultaneously processed. Therefore, two Point-wise Multiplication (PWM) units and four groups of parallel adders (ADDs₁) are used after the NTT computation. A finite state machine based controller controls the execution of the hash module, until enough random samples are generated. Then the Kyber accelerator enters a pipelined state for computational units until the hash function computation is needed again. The controller assembles the padding for the hash module, based on the hash function needed for the current state and feeds it into the hash module. In addition, the controller includes four 256-bit registers for storing and distributing the results generated by the hash module. There are several differences in the accelerator’s execution for the three different security levels of Kyber. For Kyber512, the centered-3 binomial distribution (*CBD*₃) sampling module is added to the overall architecture. For Kyber1024, the compress/ decompress modules differ from the other two security levels. Various modules are shared between the key generation and decapsulation functions on the server-side for area minimization.

In the rest of this section, we present the main modules used in the Kyber accelerator, i.e., the hash module, the sampling module, the NTT/INTT module, the PWM module, and the compression and encoding modules. All descriptions are based on the Kyber768 implementation (e.g., FIFO sizes, *k* = 3 etc.), while differences for Kyber512 and Kyber1024 are mentioned.

3.1 SHA-3 Based Hash Module

The hash module generates the random distribution samples to the sampling modules, to provide the coefficients of the noise polynomial and consequently can become the potential computational bottle-neck of the design. Hence consideration is given to match the Hash module throughput with other modules. Our HPKA uses one Keccak core [29] that is implemented serially for different SHA-3 functions. The Keccak core needs significant hardware resources, e.g., 54.2% of the total LUTs used in the Kyber design [23].

Kyber uses four different SHA-3 functions, i.e., SHA3-256, SHA3-512, SHAKE128, and SHAKE256. While the Keccak core computations remain identical for these functions,

the padding method differs for all of them. The maximum size of data output in a single computation is also different, i.e., 1,344, 1,088 and 576 bits for SHA3-128, SHA3-256, and SHA3-512, respectively.

The hash module receives up to 1,344-bits of data, in 64-bit chunks (in 21 cycles). This data is fed to the Keccak core, after the controller has added the SHA-3 function appropriate padding, e.g., the SHA3-512 function takes 64-bit data inputs in 9 cycles and in the subsequent 12 cycles 64-bit ‘0’ values are padded to get 1,344-bits of data. The Hash module defaults to the squeeze stage from the second round of computation and automatically feeds the results of the first round into the Keccak core for a further 24 rounds of computations; However the Hash module also contains an absorb signal for cases when the amount of data is greater than the maximum amount of data that can be carried in a single round.

The architecture of the Hash module comprises the input/output stages and a Keccak core, as shown in Fig. 2. These three stages are independently buffered and operate in a fully pipelined manner ensuring substantial acceleration. It has a 64 × 64 FIFO to cache-in large volumes of input data. The FIFO output is fed as 64-bit words into the 1,344-bit shift register in big-endian format (in 21 cycles). The Keccak core takes the data and applies 24 rounds of iterative operations on it. The internal state of the Keccak core is 1,600 bits in length, out of which the 1,344 most significant bits (MSBs) are taken for data output or absorbing after the 24 round calculations. At the output stage of the Hash module, the data is loaded into the 256-bit registers in the controller, 64-bit per cycle. Although 36-bit or 48-bit inputs are generally used in the sampling module in the proposed architecture, going from high-bit-width data to low-bit-width data, i.e., 48/32-bit will not cause discontinuities for input to the sampling module. For the absorbing stage, in the proposed architecture, the 1,344 MSBs of the output from the Hash module are selected to be XORed with the subsequent input.

3.2 Sampling of Noise

Kyber uses two types of sampling, namely uniform distribution sampling (*Parse*) and central binomial distribution

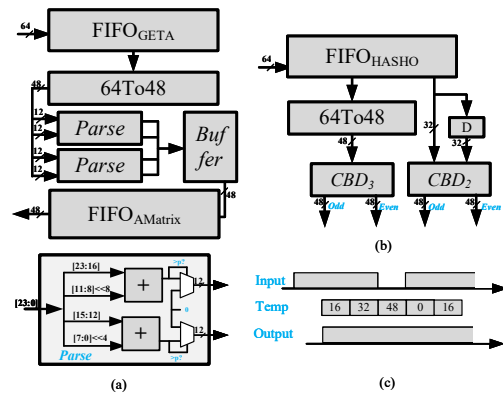


Fig. 3. Sampling modules: (a) *Parse* module, (b) *CBD* module, and (c) The output stage of *FIFO_GETA* and *FIFO_HASHO*.

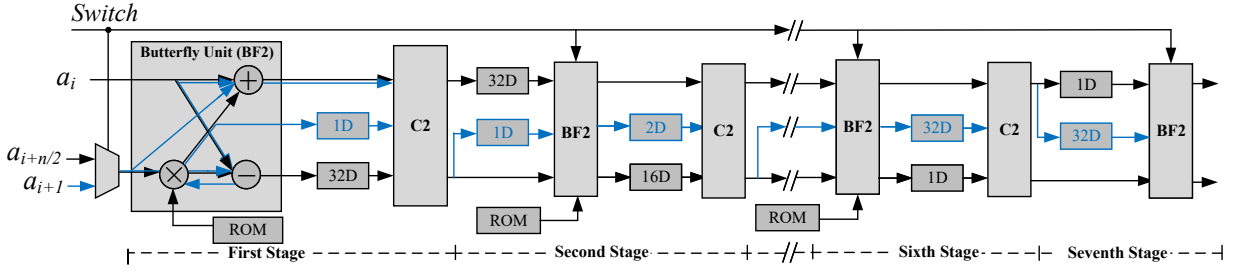


Fig. 4. The switch Radix-2 multipath delay commutator (MDC) NTT/INTT pipelined architecture with seven stages.

sampling ($CBD_\eta, \eta = 2, 3$). All three security levels of Kyber need *Parse* to get polyvector matrix \hat{A} (as well as \hat{A}^T). For Kyber512, two different *CBD* modules (CBD_2, CBD_3) are used, the other two security levels use only CBD_2 . Fig. 3(a) and (b) represent the *Parse* and CBD_η modules, respectively, where Fig. 3(c) shows the state of the FIFO when a 64-bit to 48-bit data conversion is performed.

While the hash module takes 21 clock cycles to output 1,344-bit of data (64-bit per cycle) and a single computation of Keccak takes 24 cycles, there is always meant to be a waiting interval even if pipelining is employed. The sampled data is stored in 64×32 FIFO called $FIFO_{GETA}$, converted to 48 bits, and fed into the sampling module. For the 64To48 module, as shown in Fig.3 (c), 64-bit data blocks are fed three times every four cycles. The data is internally registered to enable a consecutive stream of 48-bit outputs in 4 consecutive cycles. The results of *Parse* do not need to be stored and can be multiplied directly with the NTT results at output. As the results of *Parse* greater than 3,329 are discarded, a 48×256 FIFO is required, called $FIFO_{AMatrix}$, to filter out the *Parse* data \hat{A} . For each set of matrix \hat{A} data, it is packaged into 64 48-bit data outputs (64×4) for subsequent polynomial multiplication.

The CBD_2 and CBD_3 modules require 32 and 48 bits of data, respectively, to generate eight output results. Data output from the Hash module is stored in a 64×32 FIFO called $FIFO_{HASHO}$. The CBD_3 requires a 64To48 module before sampling while CBD_2 does not need this conversion. The $FIFO_{HASHO}$ pops out 64-bits of data once in alternate clock cycles. The upper and lower halves of this data (32 bits each) are processed in two consecutive cycles.

3.3 NTT/INTT Module

The polynomial multiplication is both a resource hungry and a computational bottle neck in a lattice based cryptography design. The Kyber parameters were tweaked to be more 'NTT-friendly' in the 3rd round submission to the NIST PQC and the use of NTT/INTT for polynomial multiplication is now part of Kyber specifications [8]. For polynomial multiplication via NTT, first the NTT of the multiplicand and multiplier polynomials is computed, a multiplication of the two vectors is carried out and then the inverse NTT (INTT) is computed to complete the polynomial multiplication. Since NTT/INTT calculations are not simultaneously performed, the same architecture is reused to minimize resource consumption. To balance area and speed and match the throughput of data in the subsequent

modules, the Radix-2 multipath delay commutator (MDC)-based architecture turns out to be the optimal choice [30].

The proposed switch-MDC-NTT (S-NTT) architecture is shown in Fig. 4. Similar to previous research [31], [32], the S-NTT consists of only seven general processing units. The first six general processing units contain a radix-2 butterfly unit (BF2), a modulo multiplier, a delay unit (D), and a two-channel commutator (C2), while the 7th general processing unit only contains a BF2 unit and a delay unit. The proposed MDC architecture employs pipelining, significantly reducing the computation time by enabling simultaneous execution of multiple consecutive NTT/INTT calculations. The BF2 modules in the processing unit are floating combinations, this is to accommodate the different order of calculations within the BF in NTT and INTT. C2 is used to select the appropriate data for the next level. For example, in NTT, the first stage output coefficient pair is (a_0, a_{64}) , then $(a_{64}, a_{65}, \dots, a_{127})$ is deferred by 32D. At the beginning of the second stage, $(a_0, a_1, \dots, a_{31})$ is deferred by 32D to obtain the coefficient pair (a_0, a_{32}) , at which point $(a_{64}, a_{65}, \dots, a_{95})$ is deferred by 32D. After the coefficient pair (a_{31}, a_{63}) , the coefficient pair changes to (a_{64}, a_{96}) . The internal architecture of the modulo multipliers follows Algorithm 4 in [23], which uses an optimized Barrett reduction [33] algorithm with multiple partitioning and addition operations. The traditional Barrett algorithm uses two constant multiplications and a shift instead of a division operation. In Kyber, it is obvious that the use of shifts for constant multiplication will be more resource and computation time efficient due to the shorter multiplication bit width. The modulo multiplier is pipelined and requires six cycles for completion, the first two perform the 12-bit multiplication and the last four perform the reduction operation.

The improved architecture eliminates the pre-processing and post-processing modulo multiplier units compared to the architecture of [30], but adds a modulo multiplier unit in the seventh stage, so that the computation time is similar to [30] (the prime in [30] is 7,681). However, despite the addition of a Div2 unit in each stage, the module contains only one adder, which is insignificant compared to the reduced number of four multipliers, so the optimised NTT/INTT unit has a lower hardware overhead. In addition, the rotation factors in the INTT are all stored as $\zeta/2$ to eliminate the Div2 unit for half of the path.

Cooley-Tukey (CT) is a common structure for NTT computation which allows the NTT to have natural order inputs and bit-reversed order outputs, while Gentleman-Sande (GS) used for INTT does the opposite. As shown in Fig. 5,

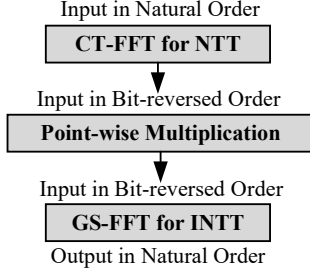


Fig. 5. NTT and INTT are computed using different flow graph-topologies to avoid bit reversal operation.

using different flow graph-topologies in NTT and INTT can completely avoid bit reversal operations in polynomial multiplication calculations. Thus, for NTT/INTT calculations, the CT and GS butterfly structures correspond to NTT and INTT calculations, respectively. The data flow of NTT/INTT is the same in the proposed architecture. In order to ensure that the output order of INTT is the natural order, the INTT in S-NTT is input in a bit-reversed order. Thus the result obtained from the polynomial multiplier can be directly input into the INTT in order. The calculation units in S-NTT, such as BF2 and modular multiplier, are all reused, but the delay unit (D) in each stage cannot be reused due to the different input order. In addition, the 12-bit multiplication in NTT/INTT uses two 6 × 12 multipliers composed of LUTs to further reduce resources. The bit-inverted order of the INTT input is used in Kyber, therefore for INTT, a independent delay unit is used, as shown in the blue line and block diagram in Fig. 4. The result of NTT/INTT is sequential data separated by parity and even, no BRAM unit is used to store data in the input and output of S-NTT. The start-up time for the very first NTT and INTT calculation (when the pipeline is not full) is 119 cycles. The cycle time for a single 256-point-wise calculation after full pipelining is 128. Hence the first NTT/INTT computation requires 119+128 cycles but the subsequent NTT/INTT computations require only 128 cycles.

3.4 Point-wise Multiplications (PWMs) and ADDs

The point-wise multiplication (PWM) in Kyber is not as straightforward as required in Ring-LWE. It requires different computations for even and odd values of the point-wise

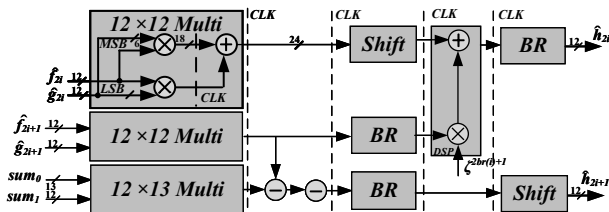


Fig. 6. The pipelined point-wise multiplication modules.

multiplication product. For Kyber, assume that polynomials \hat{f} and \hat{g} are multiplied, and the result is \hat{h} , root of unity ζ , then an optimized PWM calculation formula based on the Karatsuba algorithm from Eq.(4) was proposed in 2020 [23] as follows:

$$\hat{h}_{2i} = \hat{f}_{2i}\hat{g}_{2i} + \hat{f}_{2i+1}\hat{g}_{2i+1} \zeta^{2br(i)+1} \quad (5)$$

$$\hat{h}_{2i+1} = (\hat{f}_{2i} + \hat{f}_{2i+1})(\hat{g}_{2i} + \hat{g}_{2i+1}) - (\hat{f}_{2i}\hat{g}_{2i} + \hat{f}_{2i+1}\hat{g}_{2i+1}) \zeta^{2br(i)+1} \quad (6)$$

Thus, the result of h_{2i} comes from the sum of $\hat{f}_{2i}\hat{g}_{2i}$ and $\hat{f}_{2i+1}\hat{g}_{2i+1} \zeta^{2br(i)+1}$ and $\hat{f}_{2i+1}\hat{g}_{2i+1} \zeta^{2br(i)+1}$ comes from multiplying three 12-bit data ($\zeta^{2br(i)+1}$ can be stored in ROM by precomputation). If $\hat{f}_{2i+1}\hat{g}_{2i+1}$ is not reduced in time, it will cause the bit width of the addition to increase. Also, note that the results of $\hat{f}_{2i}\hat{g}_{2i}$ and $\hat{f}_{2i+1}\hat{g}_{2i+1}$ are directly usable in the calculation of \hat{h}_{2i+1} , and only one additional multiplication is needed when calculating \hat{h}_{2i+1} . Since the pipeline length has minimal impact on the overall design calculation time, the proposed PWM design matches the frequency of the other modules. As shown in Fig. 6, the PWM module requires eleven cycles to provide a full pipeline operation. In the first two cycles two 12-bit multiplications and one 13 × 12 bit multiplication is carried out. sum_0 and sum_1 in Fig. 6 are derived from $\hat{f}_{2i} + \hat{f}_{2i+1}$ and $\hat{g}_{2i} + \hat{g}_{2i+1}$, respectively; However sum_1 is reduced to modulo p to ensure that $sum_0 - sum_1$ does not exceed 25 bits first. During the calculation of h_{2i+1} if the result becomes negative, $24^{\text{th}}hd01000$ is added to ensure that the data in the input of Barrett reduction (BR) is positive. In the 3rd cycle, the result of \hat{h}_{2i+1} is reduced to 12-bit using the Barrett reduction (BR) module, which is the same as the BR used in the NTT module. Thereafter, \hat{h}_{2i+1} passes through a shift register until \hat{h}_{2i} is output at the same time. After $\hat{f}_{2i+1}\hat{g}_{2i+1}$ is output from the BR module, the 12-bit result is calculated with the value of ζ , added to $\hat{f}_{2i}\hat{g}_{2i}$ and the final result is reduced in one pass in the BR module.

Two PWM units are used in the overall design, allowing simultaneous calculation of the four polynomial coefficients. A 24 × 64 distributed ROM is pre-computed to provide ζ values to two channels of PWMs ($\zeta[23 : 12]$ and $\zeta[11 : 0]$). Four 12-bit adders with three input (ADD_{S1}) units are used in the design for the adder. This unit will perform the final addition for the 256 point-wise polynomial multiplication.

3.5 Storage: FIFOs and ROMs

The polynomials in lattice based cryptography schemes are large in dimension with low bit-width components, e.g., a single polynomial of Kyber has 256 elements, each of 12-bit size. High speed simultaneous access to several low bit-width data elements is critical to ensure high throughput performance but also dictates the amount of on-chip resources. The FPGA storage structures include Look-up-table RAM (LUTRAM) and Block RAMs (BRAMs). While BRAMs are dual ported, allowing fast and dual read/write access, they can be under-utilized due to their limited width-depth configurations, resulting in a waste of resources. FIFOs instead are much more resource efficient in comparison and can be custom sized as needed. However, they only access

data in the order that is pushed into the FIFO. There are distributed ROMs and two types of FIFOs in our Kyber accelerator. All pre-computed data in the design uses distributed ROM, e.g., twiddle factors in NTT/INTT and PWMs. Two types of FIFOs are used in the proposed architecture, one for buffering large amounts of data that cannot be computed on-the-fly. The depth in this type of FIFO does not need to match the total amount of data but just enough to meet the current pipeline speed. Another type of FIFO is used for the complete storage of data. Since the FIFO restricts read and write access only, for NTT computations, the order of input data is adjusted before it is input into the FIFOs. In a fully pipelined architecture, the data is continuously pushed forward. Therefore, in many cases the computation modules can also be treated as a kind of storage module, which saves a lot of storage units. Benefiting from the ability of INTT to directly use bit-reversal sequential inputs and to output data in a natural order, the proposed architecture does not use BRAM to store data at all.

3.6 Decode Modules

In Kyber KEM, data in the Kyber.CPA function must be compressed before encoding. Compression is relatively easy to implement, with the quotient being output by shifting the polynomial addition and reducing. Decompression does not require a reduction module, but all the data is multiplied by $b_3; 329=2c$, using a constant multiplier composed of LUTs to reduce DSP consumption. The encode module takes data input from the INTT output. Various encoding modules are obtained by shifting the input data and performing arithmetic operations. The bit width of the encoded output depends on the need of subsequent calculations. For example, the encoded output of v and u is 64 bits per cycle, and the m output in Kyber.CPA.Dec generates 8 bits per cycle.

used in Kyber512 and Kyber768, and it takes only eight cycles to complete a round of computation for v and u with 64 bit inputs. Fig. 7(a) shows the computation flow for the u decoding under Kyber512 and Kyber768. For the u polyvector, a single cycle produces 48 bits of data after processing 40 bits (4 groups) of data, and the 8 data blocks both generated in two consecutive cycles are separated by odd and even positions and fed into the FIFO in the same way as the sampled data. Fig. 7(a) represents the data processed in the current cycle, where the white block is the 40-bit output data, the grey block is the data that cannot be output in the current cycle and needs to be processed in the next cycle, and the red block represents the data from the previous cycle. In addition, the grey block of the Input signal represents read data blocks from the FIFO in the current cycle. Thus, every 64-bit data blocks are fetched in eight consecutive cycles, generating the 32 polyvector data of u .

The decoding in Kyber1024 is more complicated than the other two security levels due to the irregularity of individual data block (11 bits) in the decoding of the u polyvector. If the shifting method in Kyber768 is followed in Kyber1024, the generated single-round states reach more than 25. Therefore, an adaptive feedback (AF) scheme is proposed using the output control of the FIFO by the Decode module. In the AF scheme, the output of the FIFO is controlled by the decode module, whose internal counter is constantly updated with the current amount of unprocessed data to guide the output data and the FIFO read data. Each time the decode module receives 64 bits of data, the internal counter is increased by 8. As shown in Fig. 7(b), in the case of the decode u vector, when the counter is greater than 11 (88 bits), the feedback FIFO pauses to take out the data and uses the 88 bits of data to generate eight u coefficients in the subsequent two cycles. For decoding the v vector coefficients, the feedback is performed and calculated when the counter is greater than 5.

4 DATA FLOW IN THE PROPOSED KYBER ACCELERATOR

In order to ensure simultaneous execution of various modules of the Kyber accelerator, its modules are designed to support pipelining. The modules in Kyber vary greatly in terms of their input/output word sizes and speeds and need to be thoughtfully scheduled so that data flow does not halt. This is ensured by using storage units whenever needed. For example, in CBD_2 , a single 36 bit input results in eight 12-bit polynomial coefficients, while in the NTT module, only 24 bits of data are input and 24 bits are output in a single clock cycle. This section walks through the primary data flow for the main modules of Kyber.CPA.Enc (Kyber768) as an example, as shown in Fig. 8. The primary data flow consists of the Hash module, NTT/INTT module, PWM modules, and the data compression module. All cycles represent the total cycles of the current module calculation, except for the NTT/INTT, which represent the cycles needed to write the output data.

The Hash function computation time is critical to the data sampling time. For example, SHA3-256 takes 21 cycles to receive input data, 24 cycles for the internal Keccak core

Fig. 7. Decode_{dv} in different security levels. (a) Decode_{dv} of Kyber512 and Kyber768. (b) Decode_{dv} of Kyber1024.

The Decode module is more complex, with different methods for the v vector and u polyvector for different security levels in Decode_{dv}. The same Decode method is

Fig. 8. Execution order and clock cycles for main modules in Kyber.CPA.Enc of Kyber768.

computations, and up to 17 cycles (64-bit data each) can be output in one squeeze. Adding the time for data to pass through modules, a total of 70 cycles is needed for one hash calculation. However, the overhead of input and out can be made negligible by pipelining. In the Hash module, cycles for $r_{[0]}$ are computed from the start of data input, and the 70 cycles include the output stage of the previous parallel computation of the Hash module. The input data from $r_{[1]}$ is fed into the Hash module simultaneously with $r_{[0]}$ generating the output. Thus the computation cycles of the single Hash module are reduced to 53 cycles for CBD_2 sampling. To ensure that the matrix $A^{\wedge T}$ gets the 256 data points from the sampling module, the Hash module calculates SHAKE-128 by squeezing four times. Due to pipelining, the single calculation cycle is only 115 cycles.

The NTT/INTT computation is also pipelined, generating a complete set of 256 coefficients every 128 cycles. During NTT computation, the first 64 cycles of data for each output set are delayed by 64 cycles of output using shift registers to meet t^{\wedge} starting at the NTT output stage. The compress module takes the polynomial addition output after INTT calculation and generates a pair of data in a single pass, taking 512 cycles. The data flow of Fig. 8 consumes 1.9k cycles in total.

4.1 Primary Data Flow in Kyber768

The primary data flow starts with the sampling modules. A portion of the data generated by the Hash module is stored in four 256-bit registers for re-entry into the Hash module, while the rest is sampled data and is fed into the uniform Parse and binomial sampling $CBD_{2=3}$. A 64-bit data is fed into the Parse sampler every cycle, but the polyvector matrix $A^{\wedge T}$ samples are not immediately generated. The 64-bit data is first input into a 64 \times 32 $FIFO_{GETA}$ for caching before sampling. To match the 48-bit input per cycle of Parse sampling, the $FIFO_{GETA}$ outputs 64 bits of data three times in four consecutive cycles. The 64To48 module reads 192 bits of data every four cycles and splits and reassembles it into 48 bit blocks per cycle. Each sampling cycle of Parse produces four results, but not all are valid. This uncertainty in the generation of the matrix $A^{\wedge T}$ can cause difficulties in matching data in subsequent polynomial multiplications.

Fig. 9. The access process of CBD_2 and FIFOs ($FIFO_1$ and $FIFO_2$). 8 data blocks per cycle in CBD_2 output and 4 data blocks per cycle input into $FIFO_1$ and $FIFO_2$.

The 12 bits of valid results by Parse are stitched together as 48 bits and fed into a 48 \times 256 $FIFO_{AMatrix}$. Due to the late timing of the output of r^{\wedge} after the NTT calculation, the $FIFO_{AMatrix}$ waits for the complete output of r^{\wedge} from the NTT module before generating the output data to ensure that the polynomial multiplication of t^{\wedge} does not conflict with matrix $A^{\wedge T}$ in the Kyber.CPA.Enc calculation.

The binomial distribution sampler comprises CBD_2 and CBD_3 that take 32-bit and 48-bit inputs every cycle, respectively, to generate 8 results. The NTT module takes the CBD_2 output as odd/ even samples of 256 data points separately, and only a pair of data inputs is taken in each cycle. Hence the data generated by a single CBD_2 cycle (8 \times 12 bits) is much more than what the NTT module can accommodate. Therefore, the data of r at the odd positions from $CBD_{2=3}$ is shifted back by 16 cycles, after which two FIFOs ($FIFO_1$ and $FIFO_2$, 48 \times 128 FIFO) are utilized to store the sampling data. e_1 and e_2 do not need to perform a NTT calculation, the data output from CBD_2 is directly stored in $FIFO_1$, $FIFO_2$ and $FIFO_{e_2}$ (96 \times 32 FIFO).

As shown in Fig. 9, four consecutive even position data blocks in the order of 0, 2, 4, 6 are combined and written

Fig. 10. The access process of NTT results.

to FIFO₁ in the rst cycle. During the rst 16 write cycles, only FIFO₁ is used to store the even positioned data blocks. Starting from the 17th cycle of the CBD₂ output, the FIFO₁ data is replaced with the data output from the shift register, i.e., the CBD₂₌₃ odd results f 1, 3, 5, 7, and the even position data blocks from CBD₂₌₃ f 128, 130, 132, 134 are stored in FIFO₂, simultaneously. Therefore, eight data locations stored at the lowest part of these two FIFOs are f 0, 128, 2, 130, 4, 132, 6, 134. After that, the output signals (shown as greyed out in bottom of the Fig. 9, the FIFO₍₁₌₂₎ Read) are generated every four cycles to ensure that the FIFO data is read in order (f (0, 128), (2, 130), (4, 132), (6, 134)). FIFO₁ and FIFO₂ continue to take in the results of the CBD while generating output data until the FIFO₁ and FIFO₂ signal that they are empty. The data storage mechanism does not change for higher security levels but the depth of FIFO₁ and FIFO₂ is increased to match the higher data size to be cached.

The polynomial multiplication module multiplies the data, i.e., matrix \hat{A}^T and vector \hat{t}^T , by the vector \hat{r} , generated by the NTT. The operation carried out is $\hat{A}_{ji}^T = \hat{A}^T[j][i] \hat{r}[i]$ or $\hat{t}_i^T = \hat{t}^T[i] \hat{r}[i]$ ($i, j \in \{0, 1, 2\}$ in Kyber768). The polynomial \hat{t}^T is the output of the key generation function and is stored in the FIFO_{PK} (48 256 FIFO) and fed into the PWM modules simultaneously as the NTT produces its output. The $\hat{r}[i]$ generated by the NTT is fed directly into the PWM module and also stored in two FIFOs (FIFO_{SP1} and FIFO_{SP2} as well as FIFO_{EP1} and FIFO_{EP2} for Kyber.CPA.KeyGen, 24 256 FIFO). The order of the NTT outputs in continuous data separated by parity is f 0, 2, 4, 6, ..., 1, 3, 5, 7, g. Fig. 10 shows the relationship between NTT results and FIFOs during the polynomial multiplication stage. A different approach is used for the FIFO_{SP1} and FIFO_{SP2} input to ensure that the order of the PWM module input is f 0, 1, 2, 3.g. In the rst 64 cycles of the NTT operation, its two 12-bit results will be spliced and input into FIFO_{SP1}. Starting from the 65th cycle, the results of the NTT will be input into FIFO_{SP2}, thus ensuring that the order of the lowest part in these two FIFOs is f 0, 1, 2, 3g. As the security level increases, more data can be accommodated simply by increasing the depth of the FIFOs. The data in FIFO_{SP1} and FIFO_{SP2} will be used multiple

times. Therefore, the data output by current FIFOs will be stored in the same FIFOs again to reduce the use of storage resources. Since $\hat{r}[0], \hat{r}[1], \dots, \hat{r}[n]$ are computed sequentially, the output data is restored at the top of the FIFOs and does not affect the polynomial multiplication computation of the current stage. In Kyber.CCA.KeyGen, all CBD sampling results are required for the NTT calculation. Therefore, e from CBD is also stored in FIFO₁ and FIFO₂, waiting to be fed into the NTT module as s and r.

Fig. 11. The access process of nal ADDs (ADDs₁) in polynomial multiplication of Kyber768.

The nal step of the polynomial multiplication calculation should add up the different dimensional data (f $\hat{A}_{00}^T, \dots, \hat{A}_{22}^T$ g) to the same dimension, in the Kyber768 case, i.e. to compute $\hat{A}^{OT}[j] = \hat{A}_{j0}^T + \hat{A}_{j1}^T + \hat{A}_{j2}^T$ ($j \in \{0, 1, 2\}$) and in Kyber.CCA.Enc and Kyber.CCA.Dec also should calculate $\hat{t}^{OT} = \hat{t}_0^T + \hat{t}_1^T + \hat{t}_2^T$. The output of PWM is spaced, and the number of cycles between data in \hat{t}_i^T and \hat{A}_{ji}^T is not the same, although it is possible to combine $\hat{A}_{j0}^T, \hat{A}_{j1}^T, \hat{A}_{j2}^T$ using shift registers to perform additions, the area will increase dramatically as the security level increases with the addition of shift registers. Therefore, using a FIFO to access sequential results of ADDs (ADDs₁) is still the best option. FIFO_{ADD} (48 64 FIFO) is set up to write and read data from the ADDs module. Fig. 11 shows the data ow of the $\hat{A}_{00}^T + \hat{A}_{01}^T + \hat{A}_{02}^T$ computation and the states of FIFO_{ADD} in Kyber768. There are three states in FIFO_{ADD} : store-only,

store-read, and read-only. For Kyber512, only store-only and read-only states are included, while for Kyber768 and Kyber1024, the three states are all included, and the store-read state will be performed twice in Kyber1024. In the store-only state, the $FIFO_{ADD}$ is used to collect the result of the ADDs sequentially; in the store-read state, the data stored in the $FIFO_{ADD}$ is output sequentially rst, and when the ADDs produces the result, the calculation result is stored in the $FIFO_{ADD}$ at the same time; in the read-only state, the $FIFO_{ADD}$ no longer accepts input, and all the data currently stored in the FIFO is output.

The ADDs₁ in the polynomial multiplication module generates four 12-bit data per cycle, which needs to be separated and buffered in odd and even positions before the INTT calculation. Therefore, two FIFOs of the same size (24 256, $FIFO_3$ and $FIFO_4$) store the data in odd and even positions, respectively. The S-NTT module reads the data stored in $FIFO_3$ $FIFO_4$ directly in sequence and calculates them. The order of INTT results is close to the sequential order $f(0, 128), (2, 130), \dots (1, 129), g$. Firstly using a shift register to delay INTT results in 64 cycles, then the output of the shift register is combined with the output of the current INTT results into four 12-bit data blocks in the order $f(0, 1, 128, 129), (2, 3, 130, 131), g$. Two data with order $f(0, 1g)$ of these four are fed directly into ADDs₂ to complete the polynomial adder calculation and the last two data blocks with order $f(128, 129)$ are stored in a FIFO (24 64, $FIFO_5$). When S-NTT starts to output the second set of INTT results, the input of ADDs₂ is from $FIFO_5$, and the second set of INTT results will be shifted by 64 cycles.

4.2 Additional Data Flow in Key Decapsulation

Kyber.CCA.Dec contains both single Kyber.CPA.Enc and Kyber.CPA.Dec functions, while Kyber.CPA.Dec is more like a reduced Kyber main data ow. Fig. 12 shows the data ow in a single Kyber.CPA.Dec and contains mainly NTT/INTT, polynomial multiplication, and encode modules. In Kyber.CPA.Dec, when the server-side receives the cipher-text c from the client-side, it uses the decompress and decode module to process the data rst, e.g., in Kyber768, it takes 64 4 cycles to decompress and decode all the data to polyvector u and vector v . The u and v output from the decompress and decode modules perform different operations when fed into the computation unit. Polyvector u will be fed into $FIFO_1$ and $FIFO_2$, and since every two cycles form 8 sets of data blocks, the NTT starts the calculation after 32 cycles. v is stored in the FIFO (48 64 FIFO, $FIFO_{6v}$) and then output when ADDs₂ is running. When the NTT calculation is complete, results from the S-NTT module are only multiplied with s^T , which is stored in $FIFO_{SP1}$ and $FIFO_{SP2}$. Thus the NTT results do not need to be stored in FIFOs again. The s^T output from Kyber.CPA.KeyGen is still stored in $FIFO_{SP1}$ and $FFIFO_{SP2}$, and no longer uses separate resource storage. In this case, the NTT result is shifted back 64 bits using a shift register, enabling a succession of coefficients in the order $f(0, 1, 2, 3), g$ to be calculated simultaneously by the PWM unit for all four sets of data.

In addition, when the service-side architecture accepts the ciphertext c value, the c value is stored in $FIFO_A$

Fig. 12. Cycles and order in Kyber.CPA.Dec (Kyber768).

(64 128 FIFO) and $FIFO_v$ (48 256 FIFO) according to the difference between the u vector and v vector, respectively.

When Kyber.CPA.Enc generates $c^0 = (c_{1j}^0 c_2^0)$, they are compared with the data blocks from $FIFO_A$ and $FIFO_v$ respectively to define if the decapsulation results are correct.

5 IMPLEMENTATION AND RESULTS

The proposed Kyber hardware architecture has been synthesized and implemented using Xilinx Vivado 2020.1 suite targeting two different devices, e.g., Artix XC7A200 and Zynq UltraScale+ XCZU7EV. The designs proposed in this paper have passed the post-place & route (post-PAR) simulation and functional verification. The main modules of the Kyber accelerator are the same under the three different security levels, except for the increase in FIFO depth due to the need of higher amounts of data.

Table 2 shows the speed and area of our Kyber accelerator architecture, compared against the state-of-the-art architectures, for the three different security levels it offers. Since the server side needs more computational processing, it consumes more resources than the client-side. We specify the resources for the client-side and the server side respectively in Table 2, as done in [23]. In Table 2, the hardware resources (LUT, FF, DSP, BRAM) consumption is indicated for both the client/server-side, but for the number of slices column, only the slices for server side core are indicated since the server-side contains essentially all the constituent modules of Kyber. The total operation cycle count of the server and client accelerator under the same security level is constant, thus providing protection against timing attacks.

The DSPs are mainly used in the PWM module. In the S-NTT module, one pipeline cycle is added to each stage to perform the 12-bit multiplication use LUTs. Therefore, the proposed architecture uses only one DSP in each of the two PWM units, and only two DSPs are used for all security levels of the design. In addition, all storage units including distributed FIFO and distributed ROM using LUT resources and no BRAM resources are used in the proposed architecture. The accelerators at the three different security levels run at almost the same frequency as the critical paths are the same.

A comparison with related work focuses only on the hardware implementation of Kyber Round 3. Key, Enc, and Dec in Table 2 represent the Kyber.CCA.KeyGen, Kyber.CCA.Enc, and Kyber.CCA.Dec, respectively. Exploiting

TABLE 2
Post-PAR Implementation Results for our Kyber Accelerator (Area and Timing) and Comparison with The state-of-the-art Designs.

Design	LUT	FF	DSP	BRAM	Slices	ENS	F. [MHz]	Key/Enc./Dec. [K Cycles]	Key/Enc./Dec./Total [us]	² Improv. (Total Time)	AT ³ (ENS s)	FPGA Device
Kyber-512 NIST PQC Security level 1												
[24]	18,000	5,000	6	15	5,000	8,540	115	4.0/7.0/10.0	34.8/60.9/86.9/182.6	87.6%	1.56	Artix-7 XC7A100
[23]	6,785/7,412	3,981/4,644	2/2	3/3	2,126	2,914	161/167	3.8/5.1/6.7	23.4/30.5/41.3/95.2	75.9%	0.28	Artix-7 XC7A12
[18]	10,502	9,859	8	13	3,549	6,897	200	1.9/2.4/3.7	9.4/12.0/18.8/40.2	43.8%	0.28	Artix-7 XC7A100
[25] ⁴	9,347	8,186	4	6	2,337	3,913	220	2.1/3.3/4.5	9.7/14.8/20.3/44.8	49.6%	0.18	Artix-7 XC7A200
Ours	13,999/16,147	12,503/13,868	2/2	0/0	5,273	5,473	208	1.1/1.5/2.1	5.3/7.2/10.1/22.6	-	0.12	Artix-7 XC7A200
[25]	9,435	8,605	4	6	-	-	450	2.1/3.3/4.5	4.8/7.2/9.9/21.9	51.1%	-	Zynq-UltraScale+ XCZU7EV
Ours	13,300/16,004	12,533/13,890	2/2	0/0	-	-	435	1.1/1.5/2.1	2.5/3.4/4.8/10.7	-	-	Zynq-UltraScale+ XCZU7EV
Kyber-768 NIST PQC Security level 3												
[24]	16,000	6,000	9	16	4,000	8,036	115	7.0/10.0/14.0	60.9/86.9/121.7/269.5	87.3%	2.2	Artix-7 XC7A100
[23]	6,785/7,412	3,981/4,644	2/2	3/3	2,126	2,914	161/167	6.3/7.9/10.0	39.2/47.6/62.3/149.1	77.1%	0.43	Artix-7 XC7A12
[18]	11,783	10,424	12	14	3,952	7,896	200	2.7/3.2/4.8	13.3/16.3/24.0/53.6	36.4%	0.42	Artix-7 XC7A100
[25] ⁴	10,434	9,473	6	8.5	2,609	4,875	220	2.7/3.9/5.0	12.3/17.7/22.9/52.9	35.5%	0.26	Artix-7 XC7A200
Ours	15,246/16,834	12,482/13,722	2/2	0/0	5,393	5,593	208	1.7/2.4/3.0	8.2/11.5/14.4/34.1	-	0.19	Artix-7 XC7A200
[25]	10,512	10,105	6	8.5	-	-	450	2.7/3.9/5.0	6.0/8.6/11.2/25.8	36.8%	-	Zynq-UltraScale+ XCZU7EV
Ours	15,196/16,797	12,443/13,705	2/2	0/0	-	-	435	1.7/2.4/3.0	3.9/5.5/6.9/16.3	-	-	Zynq-UltraScale+ XCZU7EV
Kyber-1024 NIST PQC Security level 5												
[24]	16,000	6,000	12	17	5,000	9,532	112	10.0/14.0/18.0	86.9/121.7/156.5/365.1	86.6%	3.48	Artix-7 XC7A100
[23]	6,785/7,412	3,981/4,644	2/2	3/3	2,126	2,914	161/167	9.4/11.3/13.9	58.2/67.9/86.2/212.3	76.9%	0.62	Artix-7 XC7A12
[18]	13,347	11,639	16	16	4,585	9,321	185	3.5/4.1/6.2	17.3/20.6/31.3/69.2	29.2%	0.65	Artix-7 XC7A100
[25] ⁴	11,527	11,767	8	10.5	2,882	5,740	220	3.6/4.8/6.0	16.3/21.8/27.1/65.2	24.8%	0.37	Artix-7 XC7A200
Ours	15,666/17,827	12,476/13,979	2/2	0/0	6,022	6,222	208	2.7/3.4/4.1	13.0/16.3/19.7/49.0	-	0.30	Artix-7 XC7A200
[25]	11,598	11,606	8	10.5	-	-	450	3.6/4.8/6.0	8.0/10.6/13.2/31.8	26.4%	-	Zynq-UltraScale+ XCZU7EV
Ours	15,544/17,812	12,510/14,002	2/2	0/0	-	-	435	2.7/3.4/4.1	6.2/7.8/9.4/23.4	-	-	Zynq-UltraScale+ XCZU7EV

1. ENS (equivalent number of slices) = DSP 100 + BRAM 196 + Slices [30]

2. Time(Total) = Time(Key+ Enc+ Dec)

3. Area and time product (AT) = ENS Time (Ttotal)

4. The slices in Ref [25] are estimated by formula #LUTs/4 as one slice consists of 4 LUTs in Artix-7 [34].

a fully pipelined implementation enables simultaneous execution of several modules, resulting in a 49.6% higher speed performance for Kyber512 compared to [25] (Artix-7). Compared with [18], which also uses a high-speed NTT architecture, the speedup for Kyber512 is 44% for Key, 40% for Enc, and 46% for Dec (Artix-7). With faster devices (Zynq-UltraScale+) with larger resources, the speedup is 51.1% compared with [25] for Kyber512. As the security level increases, the speedup is reduced as the proposed architecture uses the same computational architecture. Compared with the state-of-the-art architecture [25], the total time for Key, Enc, and Dec is reduced by 35.5% and 24.8% for Kyber768 and Kyber1024, respectively (Artix-7). For Kyber1024, using the faster hardware architecture (Zynq-UltraScale+), it takes 6.2, 7.8, and 9.4us for Key, Enc, and Dec, respectively. Compared with the designs in [23]

and [24], the proposed architecture achieves a speedup of 4-7.5x for all three different security levels (Artix-7).

In terms of area, the proposed design is higher in LUTs than the previous design due to the use of more FIFO cells. However, for other on-chip resources, the proposed design uses significantly less. For Kyber512, compared with [24], [18], and [25], the number of DSP is reduced by 4, 6, and 2 blocks, and the number of BRAM is reduced by 15, 13, and 6 blocks, respectively. For a fairer comparison of resource consumption, we estimate the equivalent number of slices (ENS), as undertaken in [30]. One DSP is taken as equivalent to 100 Slices, a 36K BRAM is equivalent to 196 Slices, resulting in the ENS computation, $ENS = DSP \ 100 + BRAM \ 196 + Slices$. Compared with [24] and [18], the proposed architecture reduces ENS by 30.4-34.7% and 20.6-33.2% for the three different security

levels, respectively. Compared with [23], which implements a lightweight design, the proposed architecture uses more resources, but the speed increase outweighs the resource consumption. Compared with [25], the proposed architecture uses more LUT resources but less DSP and BRAM resources. The number of DSPs is reduced by 50.0%, 66.7%, and 75.0% for the three different security levels, respectively. To better balance the advantages of area and speed, we utilise the AT (area and time product) metric, where $AT = ENS \times \text{Time (Total)}$. As can be seen from Table 2, the proposed architecture reduces the AT in 57.1%, 54.8%, and 51.6% for the three different security levels compared with [23] and [18], respectively. Compared with [25], the proposed architecture reduces the AT by 33.3%, 26.9%, and 18.9%, where the slices in [25] is estimated by the formula $\#LUTs/4$, as a single slice in a Xilinx FPGA consists of 4 LUTs and 8 FFs [34]. Therefore, the proposed Kyber accelerator significantly improves speed and hardware efficiency compared to the state-of-the-art at all three different security levels.

While Kyber was the only KEM scheme taken up for standardized by NIST PQC [11], we compare the results of our hardware accelerator implementation with several R-LWE based Round 2 candidates and a Mod-LWR based Round 3 candidate. Compared with the Round 3 PQC candidate SABER [25], which is based on the Mod-LWR scheme, our proposed Kyber hardware architecture offers nearly double throughput performance with a comparable area at the security level 5. Compared to the Round 2 PQC candidate NewHope, our Kyber core offers nearly 2.5 times higher throughput for security level 5 [35]. Compared with the R-LWE hardware core [36], the proposed Kyber accelerator does not use BRAM units and achieves faster speed. Considering another Round 2 PQC candidate LAC, our proposed Kyber accelerator requires nearly 2.5 times fewer LUT resources [37] and speeded up by 59.7% for security level 5. In conclusion, our HPKA Kyber accelerators are marginally superior with significantly faster throughput performance, and use fewer resources than R-LWE based designs.

6 CONCLUSION

This work presents an ultra high-performance FPGA based Kyber accelerator that undertakes an optimally designed pipelined architecture for parallel execution of various modules in the design. The accelerator uses a pipelined MDC-NTT to speed up operations but to keep the area efficiency high, resource reuse is orchestrated during NTT/INTT. Multiple FIFOs are used to buffer data for pipeline balancing. We performed a hardware implementation of the proposed architecture using two different devices, the Artix-7 and the Zynq-UltraScale+. The results show that the proposed Kyber accelerator on the Artix-7 is 1.44x, 1.36x, and 1.25x faster than previous research for security levels 1/3/5, respectively. In terms of equivalent slice count, the proposed architecture reduces AT (area and time product) by 18.9-33.3% for the three different security levels. In addition, in the Zynq-UltraScale+ device, the proposed architecture achieves a speedup of 1.26-1.51x compared to the state-of-the-art designs reported to date.

REFERENCES

- [1] P. W. Shor, "Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer," *SIAM review*, vol. 41, no. 2, pp. 303–332, 1999.
- [2] R. L. Rivest, A. Shamir, and L. Adleman, "A method for obtaining digital signatures and public-key cryptosystems," *Communications of the ACM*, vol. 21, no. 2, pp. 120–126, 1978.
- [3] V. S. Miller, "Use of elliptic curves in cryptography," in *Proc. Advances in Cryptology (CRYPTO'85)*, pp. 417–426, Springer, 1985.
- [4] L. K. Grover, "A fast quantum mechanical algorithm for database search," in *Proc. 28th Annu. ACM Symp. Theory Comput.*, pp. 212–219, 1996.
- [5] J. Daemen and V. Rijmen, "Reijndael: The advanced encryption standard," *Dr. Dobbs' Journal: Software Tools for the Professional Programmer*, vol. 26, no. 3, pp. 137–139, 2001.
- [6] M. J. Dworkin, "SHA-3 standard: Permutation-based hash and extendable-output functions." 2015. [Online]. Available: <https://doi.org/10.6028/NIST.FIPS.202>.
- [7] D. Moody, "Post-quantum cryptography: NIST's plan for the future," in *Talk given at PQCrypto'16 Conference*, 2016.
- [8] R. Avanzi, J. Bos, L. Ducas, E. Kiltz, T. Lepoint, V. Lyubashevsky, J. M. Schanck, P. Schwabe, G. Seiler, and D. Stebila, "CRYSTALS-Kyber algorithm specifications and supporting documentation," *NIST PQC Round 3* 2020.
- [9] Y. Zhu, M. Zhu, B. Yang, W. Zhu, C. Deng, C. Chen, S. Wei, and L. Liu, "LWRpro: An energy-efficient configurable crypto-processor for module-LWR," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 68, no. 3, pp. 1146–1159, 2021.
- [10] A. Hülsing, J. Rijneveld, J. M. Schanck, and P. Schwabe, "NTRU-KEM-HRSS17: Algorithm specification and supporting documentation." Submission to the NIST Post-Quantum Cryptography Standardization Project, 2017.
- [11] A. Gorjan et al., "Status report on the third round of the NIST post-quantum cryptography standardization process." 2022. [Online]. Available: <https://csrc.nist.gov/publications/detail/nistir/8413/nal>.
- [12] C. Peikert and Z. Pepin, "Algebraically structured LWE, revisited," in *The 17th Theory of Cryptography Conf. (TCC)*, pp. 1–23, Springer, 2019.
- [13] E. Fujisaki and T. Okamoto, "Secure integration of asymmetric and symmetric encryption schemes," in *Proc. Advances in Cryptology (CRYPTO'99)*, pp. 537–554, Springer, 1999.
- [14] W. Guo and S. Li, "Area-efficient modular reduction structure and memory access scheme for NTT," in *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS)*, pp. 1–5, 2021.
- [15] C. Zhang, D. Liu, X. Liu, X. Zou, G. Niu, B. Liu, and Q. Jiang, "Towards efficient hardware implementation of NTT for Kyber on FPGAs," in *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS)*, pp. 1–5, 2021.
- [16] F. Yaman, A. C. Mert, E. Öztürk, and E. Savaş, "A hardware accelerator for polynomial multiplication operation of CRYSTALS-Kyber PQC scheme," in *Proc. Design, Automat. Test Eur. Conf. Exhib. (DATE)*, pp. 1020–1025, IEEE, 2021.
- [17] W. Tan, B. M. Case, A. Wang, S. Gao, and Y. Lao, "High-speed modular multiplier for lattice-based cryptosystems," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 68, no. 8, pp. 2927–2931, 2021.
- [18] M. Bisheh-Niasar, R. Azarderakhsh, and M. Mozaffari-Kermani, "High-speed NTT-based polynomial multiplication accelerator for post-quantum cryptography," in *Proc. IEEE 28th Symp. Comput. Arithmetic (ARITH)*, pp. 94–101, 2021.
- [19] Z. Chen, Y. Ma, T. Chen, J. Lin, and J. Jing, "Towards Efficient Kyber on FPGAs: A Processor for Vector of Polynomials," in *Proc. 25th Asia South Pacific Design Autom. Conf. (ASP-DAC)*, pp. 247–252, 2020.
- [20] A. Sarker, A. C. Canto, M. M. Kermani, and R. Azarderakhsh, "Error detection architectures for hardware/software co-design approaches of number-theoretic transform," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 40, pp. 2022–2032, 2022.
- [21] A. Sarker, M. Mozaffari-Kermani, and R. Azarderakhsh, "Hardware constructions for error detection of number-theoretic transform utilized in secure cryptographic architectures," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 27, no. 3, pp. 738–741, 2018.
- [22] Y. Huang, M. Huang, Z. Lei, and J. Wu, "A pure hardware implementation of CRYSTALS-Kyber PQC algorithm through resource reuse," *IEICE Electronics Express*, vol. 17, no. 17, pp. 1–6, 2020.

