

LETTER

A Pure Hardware Implementation of CRYSTALS-KYBER PQC Algorithm through Resource Reuse

Yiming Huang¹, Miaoqing Huang², Zhongkui Lei^{1a)}, and Jiaxuan Wu³

Abstract This paper presents a pure hardware implementation of CRYSTALS-KYBER algorithm on Xilinx FPGAs. CRYSTALS-KYBER is one of 26 candidate algorithms in Round 2 of NIST Post-Quantum Cryptography (PQC) standardization process. The proposed design focuses on maximizing resource utilization by reusing most of the functional modules in the encapsulation and decapsulation processes of the algorithm. For instance, the hash module integrates several different hash functions in one module. Efficient parallel and pipelined computations are applied in the NTT module. Through the analysis of simulation and synthesis results, it is found that the proposed work has the advantages of higher frequencies and lower execution times. The scheme operates at 155 MHz and 192 MHz frequencies on Xilinx Artix-7 and Virtex-7 FPGAs, respectively. Compared with the performance of an embedded Cortex-M4 processor, the hardware implementation can achieve a maximum speedup of 129 times for encryption/decryption.

key words: CRYSTALS-KYBER, Cryptography, field-programmable gate arrays (FPGAs), PQC

Classification: Integrated circuits

1. Introduction

It has been reported [1] that Post-Quantum Cryptography (PQC) algorithm are going to replace the classic public-key cryptography algorithms such as RSA, which will become vulnerable once quantum computers become mature. Therefore, NIST started the process to standardize PQC algorithms. In December 2017, 69 algorithms were released in Round One of the process. Among these 69 algorithms, 26 algorithms (including 17 encryption/key-encapsulation (PKE/KEM) and 9 signature schemes) advanced to Round Two in January 2019. CRYSTALS-KYBER is one of 17 PKE/KEM schemes in Round Two.

Most Post-Quantum Cryptography algorithms in Round Two of the NIST standardization process contain a large amount of complicated calculations and repeated math operations. Pure software implementations are not good at parallel computing. An encryption or decryption period on microprocessors could take a massive number of clock cy-

cles. Therefore, it is critical to use hardware designs to accelerate the calculation processes and assess diverse PQC algorithms on hardware platforms such as FPGAs.

Recently, several PQC schemes have been implemented in pure hardware. Among these implementations, [2] presented the hardware implementation of NewHope-Simple algorithm on Xilinx Artix-7 FPGAs. In [3] the Rainbow digital signature algorithm was implemented on FPGAs. The implementation reduced about half of the number of multiplications, and it can reconfigure different security levels. Besides, a software and hardware co-design on Xilinx Zynq FPGAs was presented in [4], including three Lattice-based PQC algorithms: FrodoKEM, Round5, and Saber. It used hardware logic to accelerate most of the computation. The hardware logic is connected to the ARM processor on the Zynq platform through the AXI bus. As for the PQC CRYSTALS-KYBER (Kyber) [5] algorithm version 2.0, the work in [6] presented an implementation on ARM Cortex-M4 embedded processor. It was able to achieve 18% performance speedup while using a tiny memory footprint.

The main difficulty of a pure hardware implementation of the Kyber PQC algorithm is the large amount of math operations, including Number Theoretic Transforms (NTT), division, and shifting computation through several matrices. We mainly used two strategies in our hardware implementation. (1) We identified the operations that contribute to the most computation in the whole process and tried to accelerate them. (2) We used BRAM (block random access memory) on FPGAs to reduce the overall cost. One challenge was to coordinate dozens of bottom-level modules (including encryption/decryption modules, pre-encryption/pre-decryption modules, etc.) to realize the key encapsulation/decapsulation mechanism of the Kyber algorithm.

This paper focuses on the hardware implementations of Kyber algorithm on two different Xilinx FPGAs, Artix-7 XC7A200T on AC701 board, and Virtex-7 XC7VX485T on VC707 board. We cover three different cryptography security levels (i.e., 512, 768, and 1024) in our implementations. AC701 is the primary target board if it contains enough hardware resources. Otherwise the target board would be VC707. Using the Verilog hardware description language, we designed the major operations of the Kyber algorithm.

¹Nanjing University of Aeronautics and Astronautics, China

²University of Arkansas, USA

³ShanghaiTech University, China

a) leizhongkui@nuaa.edu.cn

DOI: 10.1587/elex.17.20200234

Received July 07, 2020

Accepted July 16, 2020

Publicized August 14, 2020

The remainder of this paper is organized as follows. We analyze the mathematical logic of the Kyber algorithm and make an appropriate direction of implementations in Section 2. In Section 3 we demonstrate the specific scheme of overall architecture and the design of key modules. In Section 4, we illustrate main performance results and comparison with other implementation. Lastly, we give the concluding remarks in Section 5.

2. Implementation Analysis of Kyber Algorithm

Considering Kyber is an IND-CCA2-secure KEM whose mathematic basis focuses on the learning-with-errors problem in module lattices (MLWE problem [7]) presented in [8]. The whole Kyber encryption or decryption is packed by regularization data and IND-CPA cryptography with a slightly tweaked Fujisaki-Okamoto (FO) transform [9, 10] in IND-CCA2 KEM cryptography. The mathematical fundamental is Ring-LWE introduced in [11, 12], and the mathematical carrier is the polynomial rings. Recent work [13] implemented arithmetic in the polynomial ring with algorithms of Karatsuba [14] and ToomCook [15, 16]. For Kyber, the original ring $Z[X]/(X^n+1)$ is denoted by R where $n = 2^{n'}-1$ such that X^n+1 is the $2^{n'}-1$ -th cyclotomic polynomial. The Kyber algorithm picks the polynomial ring $R_q = Z_q[X]/(X^n+1)$ where q is the modulo parameter [17]. Specifically, Kyber cryptography packs $n = 256$ and $q = 3,329$ among different security levels. All parameter sets for Kyber related FPGA implementations are demonstrated in Table I.

Table I. Parameter sets for Kyber Implementation.

Algorithm	Level	Parameters (n/k/q)	Public key/Secret key/Ciphertext size p/s/c (in Bytes)
Kyber512	1	256/2/3,329	800/1,632/736
Kyber768	3	256/3/3,329	1,184/2,400/1,088
Kyber1024	5	256/4/3,329	1,568/3,168/1,568

Before performing encryption or decryption, the Kyber algorithm preconditions data to regular form, including compression, decompression, sampling, rejection, encoding, and decoding. For compression and decompression function, the former transfers an element data $x \in Z_q$ to an integer in $(d < \log_2(q))$ and vice versa. The superimposed error matrix in Kyber is sampled from a centered binomial distribution (CBD) [18] so that each of coefficient from polynomial $f \in R_q$ is sampled. After that, they are moduled by rejection function. The encoding works in serializing polynomials to byte arrays and decoding translates byte stream to polynomials vectors. We did not include key generation operation in our implementation. This work implements the encryption and decryption processes of the Kyber algorithm. We slightly modify the Kyber algorithm so that the demonstration looks more clearly and it becomes easier for FPGA implementation. Figure 1 and Figure 2 show the Kyber encryption KEM and decryption KEM approaches, respectively.

Specifically, the IND-CPA Kyber encryption in encryp-

tion KEM Flowchart Step 5 involves several modules including SHAKE-128, SHAKE-256, centered binomial distribution, NTT, pointwise-multiplied accumulation (PACC), inverse NTT, modulo polynomials, compress, decompress, encode, and decode. The IND-CPA Kyber decryption in Figure 2 Step 3 includes compress, decompress, NTT, inverse NTT, encode, and decode modules. From above, it can be found that the encryption and decryption share a lot of modules together. In other words, we can reuse a lot of modules when we implement both encryption and decryption of the Kyber algorithm.

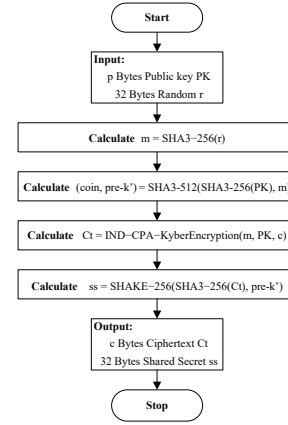


Fig. 1. The flowchart of Kyber Encryption KEM.

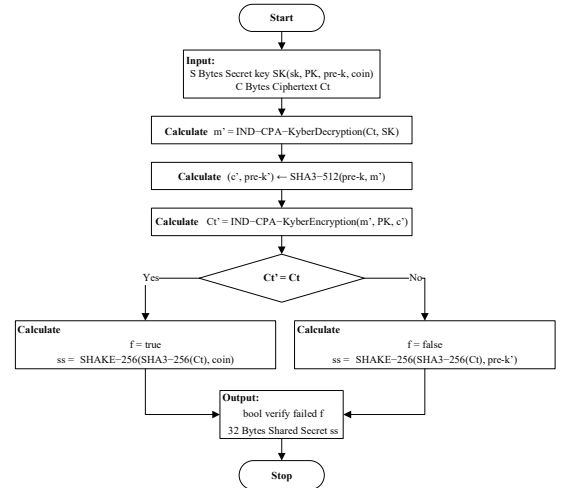


Fig. 2. The flowchart of Kyber Decryption KEM.

The reference software implementation of Kyber algorithm defines two projects, one for encryption and the other for decryption. A lot of subroutines are called in both projects. In software, calling a subroutine does not cost much physical resources. However, if we create a new instance of a module each time we call it in hardware implementation, the final implementation will occupy a massive amount of resource. For Kyber algorithm, the encryption operation and decryption operation typically do not run at the same time. Therefore, we decided to share those modules that are used in both operations. For instance, various hash functions are

3. The FPGA Design Scheme

There are a lot of intermediate data produced by modules in both the encryption and decryption processes. The amount of data would increase dramatically as the security level grows. For example, the number of hash polynomials is 9 with security level 1. It would increase to 25 with

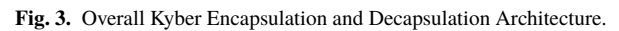


Fig. 4. Hash module structure.

A^T character is mainly generated by the same seed. Due to one-way hash function, different permutation initial-state inputs have different permutation outputs. The differences of the initial state are small discrepancies of the dimension number of A^T and the succession number of noises. When permutation is done, the noise character and A^T character are fed into the Sampling module and the Rejection module in parallel. The former module samples the coefficient to be CBD satisfied. The latter module makes data within the boundary of q , in addition to transferring character array to a matrix.

In this hash design, we use combinational logic to implement Keccak permutation process. On microprocessor, the Keccak permutation process may take hundreds of clock cycles [23]. Our hardware implementation only takes one clock cycle. Then different hash functions can be built on top of the Keccak permutation process. For example, the SHAKE-256 based noise polynomial only needs one round of Keccak permutation. The SHAKE-128 based A^T generation needs four rounds to Keccak permutation. Given different hash functions and input permutation states, the controller will schedule different rounds of Keccak permutation accordingly.

3.3 NTT Module

Many PQC algorithms consist of a lot of multiplications. Number Theoretic Transform (NTT) is typically used to improve the multiplication capacity [24, 25, 26]. Recent research on NTT applications involves with the Intel processors by Seiler [27], Lyubashevsky [28], and ARM Cortex-M4 in [29]. Kyber algorithm also applies NTT to accelerate multiplication performance. A typical NTT polynomial $f = \sum_{i=0}^{n-1} f_i X^i \in R_q$ with negacyclic is denoted as Equation 1.

$$NTT(f) = \hat{f} = \sum_{i=0}^{n-1} \hat{f}_i X^i, \quad (1)$$

$$\hat{f}_i = \sum_{j=0}^{n-1} \psi^j f_j \omega^{ij} \mod q$$

By using this equation, the multiplication between two polynomials $f, g \in R_q$ can be simplified as $f \bullet g = NTT^{-1}(NTT(f) \circ NTT(g))$. The negacyclic NTT recursion with the changing of ψ transfers the coefficient of polynomial from real number domain to NTT domain. After the point-wise multiplication of accumulation in NTT domain is done, the inverse NTT would recover the data then. According to the latest Kyber algorithm specification, one of the biggest changes in version 2.0 [5] is the NTT refinement compared with version 1.0 [30]. In version 2.0 of Kyber, the q is fixed with \mathbb{Z}_q , which includes 256^{th} roots of unity, not 512^{th} . Thus, the NTT of a polynomial $f \in R_q$ is a vector of

128 polynomials as Equation 2. It is used to handle the CBD-sampled error matrices and public key through encryption and decryption.

$$NTT(f) = \hat{f} = \sum_{i=0}^{127} (\hat{f}_{2i} + \hat{f}_{2i+1} X) \quad (2)$$

Meanwhile, according to Montgomery Reduction [31], the solution $a \bullet \omega^{-1} \pmod{q}$ is suitable for NTT transform with standard order by inputs and bitreversed order by outputs. A pre-multiplication for $\psi^j \bullet f_j$ and a Montgomery Reduction computation in Kyber have three multiplications and one subtraction, which can be pipelined in the hardware implementation. At the first clock cycle, the $(j + len)^{th}$ coefficient is pushed into Montgomery Reduction and the previous $(j - 3)^{th}$ coefficient needs to be written back to BRAM. Then, the address of coefficient changes to j , whose original data would be captured by registers marked as 1-bit `pp_state`. In the meantime, the $(j + len - 3)^{th}$ coefficient, which is subtracted by the original j^{th} coefficient, is written back to BRAM. At the third and fourth clock cycles, the coefficient with new address would be pushed into Montgomery Reduction again. The result of previous 2 computation cycles would be written back to BRAM. During the NTT computation cycles, four terms of coefficients are isolated. This isolation assures the success of the pipelining design. The structure of the NTT module is shown in Figure 5. Major operations of the NTT module are as follows.

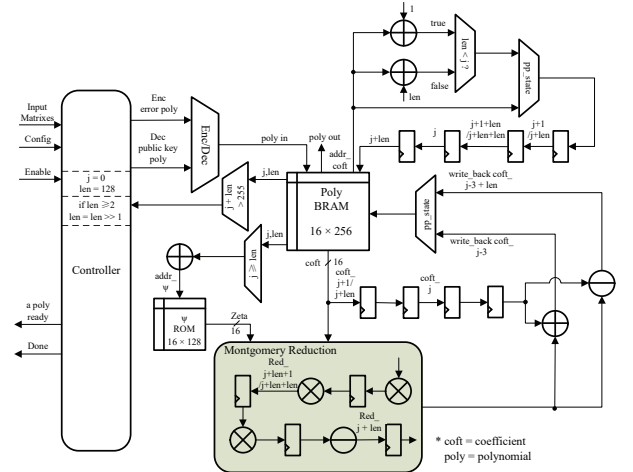


Fig. 5. NTT module structure.

- When it is enabled, the Poly BRAM would store the pending polynomial from the enc/dec mux. The initial values of `addr_ψ`, `j`, `len` would be 1, 0, 128, respectively. The first coefficient_address in cycle is denoted by $j + len$.
- Push current coefficient and ψ into Montgomery Reduction module.

Table II. Preference comparison(clock cycles) in three implementations.

Algorithm	Implementation Platform	Encapsulation [cycles]	Decapsulation [cycles]	Encapsulation Time Reduction Percentage [%]	Decapsulation Time Reduction Percentage [%]
Kyber512	This Work	49,015	68,815	—	—
	Haswell [5]	161,440	190,206	69.6	63.8
	Cortex-M4 [6]	634,000	597,000	92.3	88.5
Kyber768	This Work	77,481	102,113	—	—
	Haswell [5]	272,254	315,976	71.5	67.7
	Cortex-M4 [6]	946,000	1,167,000	91.8	91.2
Kyber1024	This Work	107,054	135,553	—	—
	Haswell [5]	396,928	451,096	73.0	70.0
	Cortex-M4 [6]	1,525,000	1,732,000	93.0	92.2

Table III. Resource Utilization and Timing Consumption.

Algorithm	Process	FPGA	Clock Freq. [MHz]	LUTs	Slices	DSPs	BRAMs	Div.IPs	Cortex-M4 Time [ms]	This Work Time [ms]	Speedup
Kyber512	Enc.	Aritix7 AC701	155	80,322	141,825	54	200.5	2	26.417	0.316	83.6
	Dec.	Aritix7 AC701	155	88,901	152,875	354	202	3	24.875	0.444	56.0
Kyber768	Enc.	Aritix7 AC701	155	97,085	153,867	36	200.5	2	46.375	0.500	92.75
	Dec.	Aritix7 AC701	155	110,260	167,293	292	202	3	44.125	0.659	67.0
Kyber1024	Enc.	Virtex7 VC707	192	119,189	162,636	36	200.5	2	72.167	0.558	129.3
	Dec.	Virtex7 VC707	192	132,918	172,489	548	202	3	68.876	0.706	97.6

- Update coefficient_address. Following the current Montgomery Reduce coefficient address, the pending addition/subtraction coefficient address would be subtracted by len to form an address pair. Then, the new Montgomery Reduction coefficient address would increase 1 (if $len < j$) or len (otherwise). In the meantime, we update the values of $addr_ψ$ and len as follows.
 - If $j \geq len$, $addr_ψ = addr_ψ + 1$; otherwise, $addr_ψ = addr_ψ$.
 - If $j + len > 255$, len does a bit right shift until len equal 2; otherwise, $len = len$.
- When addition/subtraction coefficient is read, a temporary register stores the data marked as pp_state . Meanwhile, the previous coefficient, which adds or subtracts the Montgomery result would be written back to BRAM. For example, after initial four clock cycles, $(j - 3)^{th}$, $(j - 3 + len)^{th}$, $(j - 2)^{th}$ and $(j - 2 + len)^{th}$ coefficients have all been written back to BRAM.
- The Montgomery Reduction recursion would stop when the last coefficient is pushed into the BRAM. The j and len would be 253 and 2, respectively, so that $j + len > 255$. After the final reduction is done, a polynomial NTT process is finished.

The polynomial would be handled one by one with the ready flag is set true. When all polynomials are processed, the whole NTT process is completed. With this NTT implementation design, for improving performance, we pipeline the Montgomery Reduction process. It would increase the efficiency more than three times compared with the non-pipelined implementation. The non-pipelined NTT design would take 6,550 clock cycles while the pipelined NTT module only takes 1,834 clock cycles. This implementa-

tion also solves the storage issue of polynomials by using BRAM inside the module.

4. Result

Table II compares the performance of this work with the performances of the original reference implementations on Intel Core i7-4770K (Haswell) by C language [6] and on Cortex-M4 in 24MHz [5].

With a variety of performance optimizations in hardware implementations, the amount of total clock cycles for both encryption and decryption of this proposed design reduces notably compared with Cortex-M4 implementation as well as Haswell implementation. Besides the typical techniques such as parallel execution and pipelining in hardware, innovative design techniques, such as the integration of multiple hash functions in a single module, reusing most of the functional modules during encryption and decryption, and using variable input/output widths of BRAM IPs, lead to this remarkable performance improvement.

The target platforms of hardware designs are Xilinx AC701 or VC707 FPGA boards. The detailed results, including resource utilization, single encapsulation, and decapsulation process timing and performance comparison with [5], are presented in Table III.

The results above show that the maximum clock frequencies can reach 155 MHz and 192 MHz, respectively, on AC701 and VC707 after synthesis and implementation. As expected, for both encryption and decryption, the time consumption on hardware drops significantly compared with the software implementation on Coretex-M4. The highest speedup could reach 129.3 times and the average speedup

could be 87.7 times. Also, due to functional modularization design in the top module and maximum BRAM utilization through the whole design, we are about to accommodate both encryption and decryption on the same device. Meanwhile, the use of DSPs has a certain contribution to performance speedup for operations such as additions and multiplications to accelerate the entire scheme.

5. Conclusion

In this paper, a pure hardware implementation scheme on FPGA for the CRYSTALS-KYBER Post-Quantum Cryptography algorithm is presented. This scheme implements both encryption and decryption, and uses top-down modular design approach, in which BRAM is adopted to interface communicating components. In the whole design, pipelining and parallel execution are extensively used in all modules for improving the performance. At the same time, we tried to save hardware resources by reusing bottom modules in designing upper-level modules. For example, the hash module can support multiple hash functions sharing a set of basic functional components. Compared with the software implementations running on desktop and embedded processors, the hardware implementation can achieve a speedup as high as 129 times while fitting on a single FPGA device.

References

- [1] NIST: Post-Quantum Cryptography Standardization <https://csrc.nist.gov/Projects/post-quantum-cryptography>.
- [2] Oder, Tobias and Güneysu, Tim: "Implementing the NewHope-Simple key exchange on low-cost FPGAs," International Conference on Cryptology and Information Security in Latin America (2017) 128-142 (DOI: 10.1007/978-3-030-25283-0_7).
- [3] Ahmed Ferozpur and Kris Gaj: "High-speed FPGA Implementation of the NIST Round 1 Rainbow Signature Scheme," 2018 International Conference on ReConfigurable Computing and FPGAs (ReConFig) (2018) (DOI: 10.1109/RECONF.2018.8641734).
- [4] Dang Viet B, *et al.*: "Implementing and benchmarking three lattice-based post-quantum cryptography algorithms using software/hardware codesign," 2019 International Conference on Field-Programmable Technology (ICFPT) (2019) 206-214 (DOI: 10.1109/ICFPT47387.2019.00032).
- [5] Schwabe P, *et al.*: "CRYSTALS-Kyber-Algorithm Specifications And Supporting Documentation," NIST Technical Report (2019).
- [6] B. Leon, *et al.*: "Memory-efficient high-speed implementation of Kyber on Cortex-M4," International Conference on Cryptology in Africa (2019) 209-228 (DOI: 10.1007/978-3-030-23696-0_11).
- [7] Langlois, Adeline and Stehlé, Damien: "Worst-case to average-case reductions for module lattices," Designs, Codes and Cryptography **75** (2015) 565-599 (DOI: 10.1007/s10623-014-9938-4).
- [8] B.Joppe, *et al.*: "CRYSTALS-Kyber: a CCA-secure module-lattice-based KEM," 2018 IEEE European Symposium on Security and Privacy (EuroS&P) (2018) 353-367.
- [9] Fujisaki, Eiichi and Okamoto, Tatsuaki: "Secure integration of asymmetric and symmetric encryption schemes," Annual International Cryptology Conference (1999) 537-554 (DOI: 10.1007/s00145-011-9114-1).
- [10] D. Hofheinz, *et al.*: "A Modular Analysis of the Fujisaki-Okamoto Transformation," TCC 2017 (2017) 341-371 (DOI: 10.1007/978-3-319-70500-2_12).
- [11] O. Regev: "On Lattices, Learning with Errors, Random Linear Codes, and Cryptography," Journal of the ACM **56** (2009) 1-40 (DOI: 10.1145/1568318.1568324).
- [12] V. Lyubashevsky, *et al.*: "On ideal lattices and learning with errors over rings," Journal of the ACM **60** (2013) 1-35 (DOI: 10.1145/2535925).
- [13] Matthias J. Kannwischer, *et al.*: "Faster Multiplication in $\mathbb{Z}_2^m[x]$ on Cortex-M4 to Speed up NIST PQC Candidates," (2019) (DOI: 10.1007/978-3-030-21568-2_14).
- [14] Karatsuba: "Multiplication of multidigit numbers on automata," Doklady Akad Nauk Sssr **145** (1963) 595
- [15] A. Cook Stephen, *et al.*: "On the minimum computation time of functions," Ph.D Dissertation, Harvard University, Boston (1966).
- [16] A. L. Toom: "The complexity of a scheme of functional elements realizing the multiplication of integers," Doklady Akademii Nauk Sssr **3** (1963) 496-498 (DOI: 10.1016/j.actao.2009.04.001).
- [17] P. Chris: "Public-key cryptosystems from the worst-case shortest vector problem," ACM on Theory of computing (2009) 333-342 (DOI: 10.1145/1536414.1536461).
- [18] Z. Brakerski: "Classical Hardness of Learning with Errors," Proceedings of the Annual ACM Symposium on Theory of Computing (2013) (DOI: 10.1145/2488608.2488680).
- [19] K. John, *et al.*: "SHA-3 Derived Functions: cSHAKE, KMAC, TupleHash and ParallelHash," NIST Special Publications 800-185 (2016).
- [20] Sha, NIST: "standard: Permutation-based hash and extendable-output functions," DOI, 3AD **3** (2015).
- [21] Daniel J. Bernstein, *et al.*: Tweetable FIPS 202 (2015) <https://keccak.team>.
- [22] Guido Bertoni, *et al.*: "Keccak specifications," Submission to the NIST SHA-3 competition (2011).
- [23] Adam Langley: Maybe Skip SHA-3 (2017) <https://www.imperialviolet.org/2017/05/31/skipsha3.html>.
- [24] L. Vadim, *et al.*: "SWIFFT: A modest proposal for FFT hashing," International Workshop on Fast Software Encryption (2008) 54-72 (DOI: 10.1007/978-3-540-71039-4_4).
- [25] T. Poppelmann, *et al.*: "Towards Practical Lattice-Based Public-Key Encryption on Reconfigurable Hardware," International Conference on Selected Areas in Cryptography (2013) (DOI: 10.1007/978-3-662-43414-7_4).
- [26] R. Sujoy Sinha, *et al.*: "Compact ring-LWE cryptoprocessor," International Workshop on Cryptographic Hardware and Embedded Systems (2014) 371-391 (DOI: 10.1007/978-3-662-44709-3_21).
- [27] Gregor Seiler: "Faster AVX2 optimized NTT multiplication for Ring-LWE lattice cryptography," Cryptology ePrint Archive, Report **39** (2018).
- [28] V. Lyubashevsky, *et al.*: "NTTRU: Truly Fast NTRU Using NTT," Transactions on Cryptographic Hardware and Embedded Systems **3** (2019) (DOI: 10.13154/tches.v2019.i3.180-201).
- [29] Alkim, E., *et al.*: "A new hope on ARM Cortex-M," 6th Security, Privacy, and Advanced Cryptography Engineering (2016) (DOI: 10.1007/978-3-319-49445-6_19).
- [30] B. Joppe W, *et al.*: "CRYSTALS-Kyber: a CCA-secure module-lattice-based KEM," Cryptology Archive (2017).
- [31] Montgomery, *et al.*: "Modular multiplication without trial division," Mathematics of Computation **44** (1985) 519 (DOI: 10.2307/2007970).