

# Fast and Efficient Hardware Implementation of HQC

Sanjay Deshpande<sup>1</sup>, Chuanqi Xu<sup>1</sup>, Mamuri Nawan<sup>2</sup>, Kashif Nawaz<sup>2</sup> and Jakub Szefer<sup>1</sup>

<sup>1</sup>CASLAB, Department of Electrical Engineering, Yale University, New Haven, USA,  
[sanjay.deshpande@yale.edu](mailto:sanjay.deshpande@yale.edu), [chuanqi.xu@yale.edu](mailto:chuanqi.xu@yale.edu), [jakub.szefer@yale.edu](mailto:jakub.szefer@yale.edu)

<sup>2</sup>Cryptography Research Centre, Technology Innovation Institute, Abu Dhabi, UAE  
[mamuri@tii.ae](mailto:mamuri@tii.ae), [kashif.nawaz@tii.ae](mailto:kashif.nawaz@tii.ae)

**Abstract.** This work presents a hardware design for constant-time implementation of the HQC (Hamming Quasi-Cyclic) code-based key encapsulation mechanism. HQC has been selected for the fourth-round of NIST’s Post-Quantum Cryptography standardization process and this work presents first, hand-optimized design of HQC key generation, encapsulation, and decapsulation written in Verilog targeting implementation on FPGAs. The three modules further share a common SHAKE256 hash module to reduce area overhead. All the hardware modules are parametrizable at compile time so that designs for the different security levels can be easily generated. The architecture of the hardware modules includes novel, dual clock domain design, allowing the common SHAKE module to run at slower clock speed compared to the rest of the design, while other faster modules run at their optimal clock rate. The design currently outperforms the other hardware designs for HQC, and many of the fourth-round Post-Quantum Cryptography standardization process, with one of the best time-area products as well. For the combined HighSpeed design targeting lowest security level, we show that the HQC design can perform key generation in 0.1 ms, encapsulation in 0.14 ms, and decapsulation in 0.23 ms when synthesized for an Xilinx Artix 7 FPGA. As this work shows, code-based algorithms can be competitive with other schemes when optimized hardware is developed. The presented design will further be made available under open-source license.

**Keywords:** HQC · Hamming Quasi-Cyclic · PQC · Post-Quantum Cryptography · Key Encapsulation Mechanism · Code-Based Cryptography · FPGA · Hardware Implementation

## 1 Introduction

Since 2016 NIST has been conducting a standardization process with the goal to standardize cryptographic primitives that are secure against attacks aided by quantum computers. There are today five main families of post-quantum cryptographic algorithms: hash-based, code-based, lattice-based, multivariate, and isogeny-based cryptography. Very recently NIST has selected one algorithm for standardization in the key encapsulation mechanism (KEM) category, CRYSTALS-Kyber, and four fourth-round candidates that will continue in the process. One of the four fourth-round candidates is HQC. It is a code-based KEM based on structured codes.

As the standardization process is coming to an end after the fourth round, the performance as well as hardware implementations of the algorithms are becoming very important factor in selection of the algorithms to be standardized. The motivation for our work is to understand how well hand-optimized HQC hardware implementation can be designed

and realized on FPGAs. To date, most of the post-quantum cryptographic hardware has focused on lattice-based candidates, with code-based algorithms receiving much less attention. All existing hardware implementations for HQC are based on high-level synthesis (HLS) [AAB<sup>+</sup>20]. While HLS can be used for rapid prototyping, in our experience it cannot yet outperform Verilog or other hand optimized designs. Indeed, as we show in this work, our design outperforms the existing HQC HLS design.

In addition, our hardware design competes very well with the hardware designs for other candidates currently in the fourth round of NIST’s process: BIKE, Classic McEliece, and SIKE. The presented design has best time-area product as well as time for key generation and decapsulation compared to the hardware for these designs. We also achieve similar time-area product for encapsulation when compared to BIKE. Due to limited breakdown of data for SIKE’s hardware [MLRB20] comparison to SIKE for all aspects is more difficult, but we believe our design is better since for similar area cost, their combined encapsulation and decapsulation times are two orders of magnitude larger. Detailed comparison to related work is given in Section 4.

As this work aims to show, code-based designs can be competitive with other schemes when optimized hardware is developed. Further our design is constant-time, eliminating timing-based attacks. We believe our work shows that HQC can be a strong contender in the fourth round of NIST’s process.

## 1.1 Open-Source Design

All our hardware designs reported in this paper are fully reproducible and their source code will be released as open-source after the acceptance of this paper to a journal or a conference with proceedings.

## 1.2 Paper Outline

The remainder of the paper is structured as follows. Section 2 gives background on the HQC algorithm. Section 3 presents the hardware designs of the HQC modules, as well as, it provides the evaluation results. Section 4 summarizes related work and presents comparison of the HQC design to other existing designs. Section 5 concludes the paper.

# 2 Preliminaries

In this section, we briefly introduce HQC. We first introduce notations used in this paper. Then concatenated Reed–Muller and Reed–Solomon codes that are used to encode and decode messages in HQC are presented. In the end, HQC public key encryption (PKE) and key encapsulation mechanism (KEM) are described. We refer to the specification of HQC [AAB<sup>+</sup>20] for more detailed information.

## 2.1 Notation

In this paper, we denote  $\mathbb{F}_2$  the binary finite field, and  $\mathcal{R} = \mathbb{F}_2[X]/(X^n - 1)$  the quotient ring on which vectors and operations of HQC are defined. For any field or ring,  $\mathbb{F}_2^l$  or  $\mathcal{R}^l$  denotes the field or ring of  $l$  dimensional vectors over  $\mathbb{F}_2$  or  $\mathcal{R}$ . An element  $\mathbf{x} \in \mathcal{R}$  can be represented as either a vector  $\mathbf{x} := (x_0, x_1, \dots, x_{n-1})$  or a polynomial  $\mathbf{x} := \sum_{i=0}^{n-1} x_i X^i$ . The Hamming weight of  $\mathbf{x}$  is defined as  $\sum_{i=0}^{n-1} x_i$ , i.e., the number of non-zero coefficients in  $\mathbf{x}$ .  $\mathbf{x} \leftarrow \mathcal{R}$  denotes  $\mathbf{x}$  is chosen uniformly random from  $\mathcal{R}$ , while  $\mathbf{x} \xleftarrow{w} \mathcal{R}$  denotes  $\mathbf{x}$  is randomly chosen from  $\mathcal{R}$  and the Hamming weight of  $\mathbf{x}$  is  $w$ . Optionally, a random seed can be specified for the sampling process, and the sampling process with random seed  $\theta$  is denoted as  $\mathbf{x} \xleftarrow{w, \theta} \mathcal{R}$ . For  $\mathbf{x}, \mathbf{y} \in \mathcal{R}$ , the addition  $\mathbf{a} = \mathbf{x} + \mathbf{y} \in \mathcal{R}$  and is

defined as  $a_i = x_i + y_i \bmod 2$  or  $a_i = x_i \text{ xor } y_i$  for  $i = 0, \dots, n-1$ . The multiplication  $\mathbf{a} = \mathbf{x} \cdot \mathbf{y} \in \mathcal{R}$  and is defined as  $a_k = \sum_{i+j \equiv k \pmod n} x_i \cdot y_j \bmod 2$  for  $k = 0, \dots, n-1$ .  $\text{Encode}(\cdot)$  and  $\text{Decode}(\cdot)$  are the encode and decode function of concatenated Reed–Muller and Reed–Solomon codes, which will be introduced in Section 2.2.1.  $\mathcal{G}(\cdot), \mathcal{H}(\cdot), \mathcal{K}(\cdot)$  are hash functions with domain separation bytes 3, 4, 5 respectively. Parameters  $n, w, w_r$  depend on the security level, which can be found on Table 1.

## 2.2 HQC PKE and KEM Schemes

To introduce HQC KEM, we first introduce how to encode and decode concatenated Reed–Muller and Reed–Solomon codes, which are important parts in HQC PKE to encode and decode messages. Then we introduce HQC PKE. Finally, HQC KEM that achieves IND-CCA2 by performing the Fujisaki–Okamoto [HHK17] transformation is described.

### 2.2.1 Concatenated Reed–Muller and Reed–Solomon Codes

Concatenated Reed–Muller and Reed–Solomon Codes are used in the encode and decode function of HQC to encode and decode messages. Specifically, the encode of concatenated Reed–Muller and Reed–Solomon Codes is to first encode with Reed–Solomon code, and then the output is encoded with Reed–Muller code. The decode process is executed in the reverse order, i.e., first decode with Reed–Muller code, and then the output is decoded with Reed–Solomon code. The parameter sets for the concatenated code can be found in Table 1. Notice that shortened Reed–Solomon code and duplicated Reed–Muller code instead of the conventional codes are used in HQC.

**Encode of shortened Reed–Solomon code.** In the polynomial representation, the message can be denoted as  $u(x) = u_0 + \dots + u_{k-1}x^{k-1} \in \mathbb{F}_{2^8}[x]/(x^8 + x^4 + x^3 + x^2 + 1)$ . The codeword is given by  $c(x) = b(x) + x^{n-k}u(x)$ , where  $b(x) = x^{n-k}u(x) \bmod g(x)$ , and  $g(x)$  is the generator polynomial<sup>1</sup>, which is shown below for different security levels:

$$g_{\text{hqc-128}}(x) = 89 + 69x + 153x^2 + 116x^3 + 176x^4 + 117x^5 + 111x^6 + 75x^7 + 73x^8 + 233x^9 + 242x^{10} + 233x^{11} + 65x^{12} + 210x^{13} + 21x^{14} + 139x^{15} + 103x^{16} + 173x^{17} + 67x^{18} + 118x^{19} + 105x^{20} + 210x^{21} + 174x^{22} + 110x^{23} + 74x^{24} + 69x^{25} + 228x^{26} + 82x^{27} + 255x^{28} + 181x^{29} + x^{30}$$

$$g_{\text{hqc-192}}(x) = 45 + 216x + 239x^2 + 24x^3 + 253x^4 + 104x^5 + 27x^6 + 40x^7 + 107x^8 + 50x^9 + 163x^{10} + 210x^{11} + 227x^{12} + 134x^{13} + 224x^{14} + 158x^{15} + 119x^{16} + 13x^{17} + 158x^{18} + x^{19} + 238x^{20} + 164x^{21} + 82x^{22} + 43x^{23} + 15x^{24} + 232x^{25} + 246x^{26} + 142x^{27} + 50x^{28} + 189x^{29} + 29x^{30} + 232x^{31} + x^{32}$$

$$g_{\text{hqc-256}}(x) = 49 + 167x + 49x^2 + 39x^3 + 200x^4 + 121x^5 + 124x^6 + 91x^7 + 240x^8 + 63x^9 + 148x^{10} + 71x^{11} + 150x^{12} + 123x^{13} + 87x^{14} + 101x^{15} + 32x^{16} + 215x^{17} + 159x^{18} + 71x^{19} + 201x^{20} + 115x^{21} + 97x^{22} + 210x^{23} + 186x^{24} + 183x^{25} + 141x^{26} + 217x^{27} + 123x^{28} + 12x^{29} + 31x^{30} + 243x^{31} + 180x^{32} + 219x^{33} + 152x^{34} + 239x^{35} + 99x^{36} + 141x^{37} + 4x^{38} + 246x^{39} + 191x^{40} + 144x^{41} + 8x^{42} + 232x^{43} + 47x^{44} + 27x^{45} + 141x^{46} + 178x^{47} + 130x^{48} + 64x^{49} + 124x^{50} + 47x^{51} + 39x^{52} + 188x^{53} + 216x^{54} + 48x^{55} + 199x^{56} + 187x^{57} + x^{58}$$

The generator polynomial can also be computed by  $g(x) = \prod_{i=0}^{n-k-1} (x - \alpha^i)$ , where  $\alpha$  is the primitive element of the field.

**Decode of shortened Reed–Solomon code.** We denote the codeword to be  $v(x) = v_0 + v_1x + \dots + v_{n-1}x^{n-1}$ , the error polynomial to be  $e(x) = e_0 + e_1x + \dots + e_{n-1}x^{n-1}$ , and the received word to be  $r(x) = r_0 + r_1x + \dots + r_{n-1}x^{n-1}$ . With these definitions,  $r(x) = v(x) + e(x)$ . The primitive element  $\alpha$  of the field satisfies  $v(\alpha^i) = 0$  for  $i = 1, \dots, 2d$  (notice  $2d = n - k$ ), since  $g(\alpha^i) = 0$  and  $v(x) \bmod g(x) = 0$ . If there is no error

<sup>1</sup>Zeroth Coefficient of  $g_{\text{hqc-128}}(x)$  is updated based on the software reference implementation given at <https://pqc-hqc.org/implementation.html>

in the received word,  $r(\alpha^i) = v(\alpha^i) = 0$ , so  $e(\alpha^i) = 0$ . Otherwise, we can denote  $r(\alpha^i) = e(\alpha^i) = e_{j_1}(\alpha^i)^{j_1} + \dots + e_{j_t}(\alpha^i)^{j_t}$  where  $e(x)$  has  $t$  errors at locations  $j_1, \dots, j_t$ . Let us define:

$$\begin{aligned} S_i &= r(\alpha^i) = e(\alpha^i) = e_{j_1}(\alpha^{j_1})^i + \dots + e_{j_t}(\alpha^{j_t})^i, \quad i = 1, \dots, 2d \\ \sigma(x) &= (1 + \alpha^{j_1}x)(1 + \alpha^{j_2}x) \dots (1 + \alpha^{j_t}x) = 1 + \sigma_1x + \sigma_2x^2 + \dots + \sigma_tx^t \\ Z(x) &= 1 + (S_1 + \sigma_1)x + (S_2 + \sigma_1S_1 + \sigma_2)x^2 + \dots + (S_t + \sigma_1S_{t-1} + \sigma_2S_{t-2} + \dots + \sigma_t)x^t \end{aligned}$$

Then the error value at location  $j_l$  can be computed by:

$$e_{j_l} = \frac{Z((\alpha^{j_l})^{-1})}{\prod_{i=1, i \neq l}^t [1 + \alpha^{j_i} \cdot (\alpha^{j_l})^{-1}]}$$

The decode steps are:

1. Compute  $S_i = r(\alpha^i), i = 1, \dots, 2d$ .
2. Because  $\sigma((\alpha^{j_i})^{-1}) = 0, i = 1, \dots, t$ , the coefficients  $\sigma_i, i = 1, \dots, t$  can be calculated from the linear equation set:  $0 = \sum_{i=1}^t e_{j_i}(\alpha^{j_i})^{l+t} \sigma((\alpha^{j_i})^{-1}) = S_{l+t} + \sigma_1S_{l+t-1} + \sigma_2S_{l+t-2} + \dots + \sigma_tS_l, l = 0, \dots, t-1$ .
3. The roots of  $\sigma(x)$  can be calculated, which are  $(\alpha^{j_i})^{-1}, i = 1, \dots, t$ .
4.  $Z((\alpha^{j_l})^{-1}), l = 1, \dots, t$  can be computed, and so do  $e_{j_l}, l = 1, \dots, t$ .
5. The codeword can be computed by  $v(x) = r(x) - e(x)$ .

**Encode of duplicated Reed–Muller code.** Encode of duplicated Reed–Muller code is to directly perform a matrix vector multiplication. The generator matrix is shown below (note that numbers are big endian and in hexadecimal):

$$\mathbf{G} = \begin{pmatrix} \text{aaaaaaaa} & \text{aaaaaaaa} & \text{aaaaaaaa} & \text{aaaaaaaa} \\ \text{cccccccc} & \text{cccccccc} & \text{cccccccc} & \text{cccccccc} \\ \text{f0f0f0f0} & \text{f0f0f0f0} & \text{f0f0f0f0} & \text{f0f0f0f0} \\ \text{ff00ff00} & \text{ff00ff00} & \text{ff00ff00} & \text{ff00ff00} \\ \text{ffff0000} & \text{ffff0000} & \text{ffff0000} & \text{ffff0000} \\ \text{00000000} & \text{ffffffffff} & \text{00000000} & \text{ffffffffff} \\ \text{00000000} & \text{00000000} & \text{ffffffffff} & \text{ffffffffff} \\ \text{ffffffffff} & \text{ffffffffff} & \text{ffffffffff} & \text{ffffffffff} \end{pmatrix}$$

If the message is  $\mathbf{m} = (m_0, \dots, m_7) \in \mathbb{F}_{2^8}$ , then  $\mathbf{c} = \mathbf{mG}$ , and the codeword is given by duplicating  $\mathbf{c}$  3 or 5 times, depending on the security level.

**Decode of duplicated Reed–Muller code.** The decoding of duplicated Reed–Muller codes is done in three steps:

1. The first step is applying the function  $F$  on the received codeword. Let  $v$  be a duplicated Reed–Muller codeword with multiplicity 3, it can be seen as  $v = (a_1b_1c_1, \dots, a_{n_1}b_{n_1}c_{n_1})$  where each  $a_i, b_i, c_i$  has 128 bits size ( $a_i = (a_{i_1}, \dots, a_{i_{128}})$ ,  $b_i = (b_{i_1}, \dots, b_{i_{128}})$  and  $c_i = (c_{i_1}, \dots, c_{i_{128}})$ ). The transformation  $F$  is applied to each element in  $v$  as  $((-1)^{a_{i_1}} + (-1)^{b_{i_1}} + (-1)^{c_{i_1}}, \dots, (-1)^{a_{i_{128}}} + (-1)^{b_{i_{128}}} + (-1)^{c_{i_{128}}})$ . For multiplicity 5, it follows the same process.
2. The second step is applying Hadamard transform on the output of the previous step.
3. The third step is finding the location of the highest value on the output of Hadamard transform. When the peak is positive, we add all-one-vector. If there are two identical peaks, we take the peak with smallest value in the lowest 7 bits.

### 2.2.2 HQC PKE

**Key generation.** First a vector  $\mathbf{h}$  is sampled uniformly random, which is viewed as the vector to generate a circulant matrix and further a systematic quasi-cyclic code of index 2. More specifically, let  $\mathbf{h} = (h_0, \dots, h_{n-1})$ . Then

$$\mathbf{rot}(\mathbf{h}) = \begin{pmatrix} h_0 & h_{n-1} & \cdots & h_1 \\ h_1 & h_0 & \cdots & h_2 \\ \vdots & \vdots & \ddots & \vdots \\ h_{n-1} & h_{n-2} & \cdots & h_0 \end{pmatrix}$$

is a circulant matrix, and  $\mathbf{H} = [\mathbf{I}|\mathbf{rot}(\mathbf{h})]$  is the parity-check matrix of a systematic quasi-cyclic code of index 2. The secret key is composed of two vectors  $\mathbf{x}, \mathbf{y}$  that are sampled with a specified weight  $w$ .  $[\mathbf{x}|\mathbf{y}]$  can be viewed as a random codeword with a random error. Its syndrome is  $\mathbf{s} = [\mathbf{x}|\mathbf{y}]\mathbf{H}^T = \mathbf{x} + \mathbf{y} \times \mathbf{rot}(\mathbf{h})^T = \mathbf{x} + \mathbf{h} \cdot \mathbf{y}$ , where  $\mathbf{y} \times \mathbf{rot}(\mathbf{h})^T$  is defined as the general vector-matrix multiplication. The public key is composed of  $\mathbf{h}$  and the syndrome  $\mathbf{s}$ .

**Encryption.** Similar to key generation, three vectors  $\mathbf{r}_1, \mathbf{r}_2, \mathbf{e}$  are sampled with a specified weight  $w_r$ . Then the syndrome  $\mathbf{u}$  of  $[\mathbf{r}_1, \mathbf{r}_2]$  is computed. The message is encoded by concatenated Reed–Muller and Reed–Solomon code introduced previously, as well as added by  $\mathbf{s} \cdot \mathbf{r}_2 + \mathbf{e}$  to get  $\mathbf{v}$ . The final ciphertext comprises of  $\mathbf{u}$  and  $\mathbf{v}$ , i.e.,  $\mathbf{c} = [\mathbf{u}|\mathbf{v}]$ .

**Decryption.**  $\mathbf{v} - \mathbf{u} \cdot \mathbf{y}$  is decoded by concatenated Reed–Muller and Reed–Solomon code. The message can be correctly decoded whenever the Hamming weight of the given element is less than the minimum distance of the code. The probability that the message cannot be decoded from  $\mathbf{v} - \mathbf{u} \cdot \mathbf{y}$  is shown to be very small [AAB<sup>+</sup>20].

---

**Algorithm 1** HQC.PKE.KeyGen() and HQC.KEM.KeyGen()

```

1:  $\mathbf{h} \leftarrow \mathcal{R}$ 
2:  $(\mathbf{x}, \mathbf{y}) \xleftarrow{w} \mathcal{R}^2$ 
3:  $\mathbf{s} := \mathbf{x} + \mathbf{h} \cdot \mathbf{y}$ 
4: return  $(\mathbf{pk} := (\mathbf{h}, \mathbf{s}), \mathbf{sk} := (\mathbf{x}, \mathbf{y}))$ 

```

---



---

**Algorithm 2** HQC.PKE.Encrypt( $\mathbf{pk} = (\mathbf{h}, \mathbf{s}), \mathbf{m}, \theta$ )

```

1:  $(\mathbf{r}_1, \mathbf{r}_2, \mathbf{e}) \xleftarrow{w_r, \theta} \mathcal{R}^3$ 
2:  $\mathbf{u} := \mathbf{r}_1 + \mathbf{h} \cdot \mathbf{r}_2$ 
3:  $\mathbf{t} := \text{Encode}(\mathbf{m})$ 
4:  $\mathbf{v} := \mathbf{t} + \mathbf{s} \cdot \mathbf{r}_2 + \mathbf{e}$ 
5: return  $\mathbf{c} := (\mathbf{u}, \mathbf{v})$ 

```

---



---

**Algorithm 3** HQC.PKE.Decrypt( $\mathbf{sk} = (\mathbf{x}, \mathbf{y}), \mathbf{c} = (\mathbf{u}, \mathbf{v})$ )

```

1:  $\mathbf{m}' := \text{Decode}(\mathbf{v} - \mathbf{u} \cdot \mathbf{y})$ 
2: return  $\mathbf{m}'$ 

```

---

### 2.2.3 HQC KEM

**Key generation.** Key generation in KEM is the same as key generation in PKE.

**Encapsulation.** The message  $\mathbf{m}$  is sampled to generate the shared secret. The random seed  $\theta = \mathcal{G}(\mathbf{m})$ , which will be used in the encryption to control the randomness.  $\mathbf{m}$  is then encrypted to generate  $\mathbf{c}$ . Finally, the shared secret  $K = \mathcal{K}(\mathbf{m}, \mathbf{c})$ , and the ciphertext is  $[\mathbf{c}|\mathcal{H}(\mathbf{m})]$ .

**Decapsulation.**  $\mathbf{c}$  is used to retrieve the message  $\mathbf{m}'$ . The decryption process may not be correct and thus returns a wrong message. Therefore, the same process as encapsulation

**Table 1:** Parameter sets for HQC.  $n$  is the length of the vector (polynomial).  $n_1$  is the length of the Reed–Solomon code.  $n_2$  is the length of the Reed–Muller code.  $w$  is the weight of vectors  $\mathbf{x}, \mathbf{y}$ .  $w_r$  is the weight of vectors  $\mathbf{r}_1, \mathbf{r}_2, \mathbf{e}$ .  $[n, k, d]$  of Reed–Solomon and Reed–Muller codes are shown in the last two columns, and they are the length, the dimension, and the minimum distance of the code. In HQC, shortened Reed–Solomon code and duplicated Reed–Muller code are used. The multiplicity for duplicated Reed–Muller code is 3, 5, 5 for hqc-128, hqc-192, hqc-256.

Instance	$n$	$w$	$w_r$	security	$p_{\text{fail}}$	Reed–Solomon	Reed–Muller
hqc128	17,669	66	75	128	$< 2^{-128}$	[46, 16, 15]	[384, 8, 192]
hqc192	35,851	100	114	192	$< 2^{-192}$	[56, 24, 16]	[640, 8, 320]
hqc256	57,637	131	149	256	$< 2^{-256}$	[90, 32, 29]	[640, 8, 320]

needs to be done and the ciphertext needs to be checked with the received ciphertext. Finally, whether there are mistakes is returned.

---

**Algorithm 4** HQC.KEM.Encapsulate(pk = ( $\mathbf{h}, \mathbf{s}$ ))

---

```

1:  $\mathbf{m} \leftarrow \mathbb{F}_2^k$ 
2:  $\theta := \mathcal{G}(\mathbf{m})$ 
3:  $\mathbf{c} := (\mathbf{u}, \mathbf{v}) = \text{HQC.PKE.Encrypt}(\text{pk}, \mathbf{m}, \theta)$ 
4:  $K := \mathcal{K}(\mathbf{m}, \mathbf{c})$ 
5:  $\mathbf{d} := \mathcal{H}(\mathbf{m})$ 
6: return ( $K, (\mathbf{c}, \mathbf{d})$ )

```

---



---

**Algorithm 5** HQC.KEM.Decapsulate(sk = ( $\mathbf{x}, \mathbf{y}$ ),  $\mathbf{c}, \mathbf{d}$ )

---

```

1:  $\mathbf{m}' := \text{HQC.PKE.Decrypt}(\text{sk}, \mathbf{c})$ 
2:  $\theta' := \mathcal{G}(\mathbf{m}')$ 
3:  $\mathbf{c}' := (\mathbf{u}', \mathbf{v}') = \text{HQC.PKE.Encrypt}(\text{pk}, \mathbf{m}', \theta')$ 
4:  $\mathbf{d}' := \mathcal{H}(\mathbf{m}')$ 
5:  $K' := \mathcal{K}(\mathbf{m}', \mathbf{c})$ 
6: if  $\mathbf{c} \neq \mathbf{c}'$  or  $\mathbf{d} \neq \mathbf{d}'$  then
7:   return ( $K', 0$ )
8: else
9:   return ( $K', 1$ )
10: end if

```

---

### 3 Hardware Design of HQC

HQC Key Encapsulation Mechanism (HQC-KEM) consists of three main primitives: Key Generation, Encapsulation, and Decapsulation. The algorithms for each primitive were shown in Algorithm 1, Algorithm 4, and Algorithm 5, respectively. These primitives are built upon the HQC Public Key Encryption (HQC-PKE) primitives shown in Algorithm 1, Algorithm 2, and Algorithm 3. Which in turn are built upon other, more basic building blocks. In this work, we implement optimized and parameterizable hardware designs for all the primitives and the building blocks from scratch. In the following subsections we briefly discuss all the building blocks and provide comparisons with any existing designs. The main building blocks involved for each of the primitives are as follows:

- Key Generation: Fixed weight vector generator, PRNG based random vector generator, polynomial multiplication, modular addition, and SHAKE256
- Encapsulation: Encrypt, SHAKE256
- Decapsulation: Decrypt, Encrypt, SHAKE256

**Table 2:** SHAKE256 module area and timing information, data based on synthesis results for Artix 7 board with xc7a200t-3 FPGA chip. Formula for time-area product,  $T \times A$ , is  $(LUT * Time)/10^3$ .

Parallel Slices	Resources				F (MHz)	Cycles (cyc.)	Time (us)	T x A
	Logic (LUT) (DSP)		Memory (FF) (BR)					
1	1,437	0	498	0	163	5,010	30.74	44
2	1,558	0	466	0	167	2,306	13.81	21
4	1,625	0	370	0	157	1,086	6.92	11
8	1,958	0	280	0	158	542	3.43	6
16	2,819	0	236	0	164	270	1.65	4

### 3.1 Modules Common Across the Design

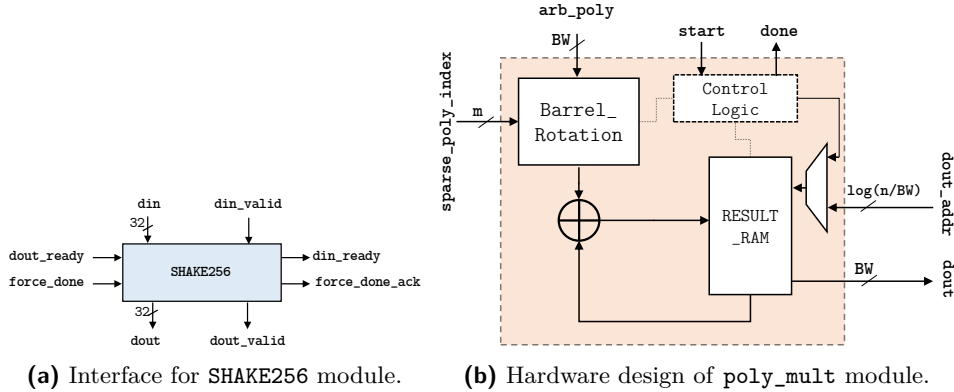
In this section we give a high-level overview of hardware designs of the building blocks that are used across the HQC-KEM and HQC-PKE.

#### 3.1.1 SHAKE256

HQC uses SHAKE256 for multiple purposes e.g., as a PRNG for fixed weight vector generation and random vector generation in Key Generation, as a PRNG for fixed weight vector generation in Encryption, and for hashing in encapsulation and decapsulation. We use the SHAKE256 module described in [CCD<sup>+</sup>22] (which was originally designed based on Keccak design from [WTJ<sup>+</sup>20]) to perform SHAKE256 operations. We further tailor the SHAKE256 hardware module as per the requirement for our hardware design:

- The existing SHAKE256 module [CCD<sup>+</sup>22] operates with command based interface where the number of input bytes to be processed and number of output bytes required are specified before starting the hash operation and there is no command to request for additional bytes. We modify the exiting design and add an additional command which provides the capability of requesting additional bytes. The purpose of adding this command is to support the fixed weight vector generation process described in Section 3.1.4.
- Since our modification of SHAKE256 holds the current state and does not automatically return to its new input loading state, we modify the operation of the existing forced exit signal to return the SHAKE256 module to default state. To support the dual clock domain design described in Section 3.7, we also add a forced exit acknowledgement port.

In addition to the aforementioned changes we further explored options for optimizing the maximum clock frequency by pipelining the critical path. We note that there are several such critical paths throughout the design and pipelining each path added severe overhead in terms of clock cycles with minimal improvement in the maximum clock frequency. Consequently the results presented in Table 2 are optimal time and area results for the given hardware architecture. We use a similar performance parameter `parallel_slices` as described in the original keccak design in [WTJ<sup>+</sup>20]. The SHAKE256 module has a fixed 32-bit data ports and data input and output is based on typical ready-valid protocol. The results targeting Xilinx Artix 7 xc7a200t FPGA are shown in Table 2. The clock cycle numbers provided in the Table 2 are for processing one block of input and generating one block of output (where each block size is 1088-bits). There are five different options to chose for the `parallel_slices` which provide different time-area trade-offs. We choose `parallel_slices = 16` as it provides the best time-area product. An interface diagram of the SHAKE256 module is shown in Figure 1a. For brevity, we represent all the ports interfacing with the SHAKE256 module with  $\Leftrightarrow$  in all further block diagrams in this paper.



(a) Interface for SHAKE256 module. (b) Hardware design of poly\_mult module.

**Figure 1:** Block diagram for interface of the SHAKE256 module and poly\_mult module.

### 3.1.2 Polynomial Multiplication

HQC uses polynomial multiplication operation in all the primitives of HQC-KEM. The polynomial multiplication operation is multiplication of two polynomials with  $n$  components in  $\mathbb{F}_2$ . After profiling all the polynomial multiplication operations from the HQC specification document and the reference design [AAB<sup>+</sup>20], we note that in all the polynomial multiplication operations, one of the inputs is a sparse fixed weight vector (with weight  $w$  or  $w_r$  in Table 1) of width  $n$ -bits. Consequently we design a sparse polynomial multiplication technique with an interleaved reduction  $X^n - 1$  (values of  $n$  can be found in Table 1).

The motivation behind our polynomial multiplication unit is as follows: we represent the non-sparse arbitrary polynomial as `arb_poly` and the sparse fixed-weight polynomial by `sparse_poly`. For `sparse_poly`, rather than storing the full polynomial we only store the indices for non-zero values. Then, the multiplication is performed by left shifting `arb_poly` with each index of `sparse_poly` and then performing reduction of the resultant vector in an interleaved fashion. Since the value of  $n$  is large in all parameter sets of HQC, we take a sequential approach for performing the left shift. We implement a sequential left shift module similar to one in [Gig04]. The shift module described [Gig04] uses a register based approach and is not scalable when the length of the input is as large as the  $n$  value for the HQC parameters (due to a larger resource utilization and complex routing). This issue is circumvented in our design by implementing a block RAM based sequential variable shift module with a dual port BRAM and small barrel rotation unit. The barrel rotation unit and the block RAM widths are used as performance parameter (BW - Block Width) for the shift module and in turn for the whole polynomial multiplication unit. A similar implementation of sequential variable shift module was previously described in [DdPM<sup>+</sup>21], however we could not readily use their implementation because the shift module is tightly embedded with the other modules for a different application and we re-implemented our version.

The hardware design of our polynomial multiplication module (`poly_mult`) is shown in Figure 1b. The `arb_poly` input to the `poly_mult` module is loaded sequentially and the width of each chunk of `arb_poly` is equal to BW (making total number of chunks in polynomial equal to  $\text{RAMDEPTH} = \text{ceil}(n/\text{BW})$ ). We store the least significant part of the polynomial at lowest address of the block RAM and most significant part at the highest address. Since the polynomial length in HQC parameters is equal to  $n$  and is not divisible by BW ( $n$  is a prime) we append the most significant part of the polynomial with zeros. For `sparse_poly`, one index is loaded at a time. While performing the shift operation we also perform the reduction ( $X^n - 1$ ) in an interleaved fashion. Since the



**Table 3:** Comparison of our area and timing information `poly_mult` module with the other sparse polynomial multiplication units targeting Artix 7 board with `xc7a200t` FPGA chip.

BW (bits)	Resources						Cycles (cyc.)	Time (us)	T x A
	Logic (LUT)	(DSP)	Memory (FF)	(BR)	F (MHz)				
Our <code>poly_mult</code> module, Polynomial Length* = 12,323, $W_{SPARSE}^* = 71$									
32	396	0	181	1	270	27,621	0.10	41	
64	599	0	205	2	277	13,918	0.05	30	
128	1,438	0	456	4	238	7,102	0.03	43	
General Sparse Multiplier, Polynomial Length* = 12,323, $W_{SPARSE}^* = 71$ [RBCGG21]									
32	319	0	127	2	234	27,691	0.12	38	
64	549	0	190	4	222	13,988	0.06	35	
128	1,136	0	381	8	185	7,172	0.04	44	
Sparse Multiplier, Polynomial Length* = 10,163, $W_{SPARSE}^* = 71$ [HWCW19]									
32	100 <sup>†</sup>	—	—	2	240	158,614	0.66	—	
64	157 <sup>†</sup>	—	—	3	220	90,880	0.41	—	
128	292 <sup>†</sup>	—	—	5	210	51,688	0.24	—	

<sup>†</sup> = Slices (no info on LUTs), <sup>+</sup> Length of the non-sparse arbitrary polynomial, \* = Weight of the sparse polynomial input

result of multiplying two  $n$ -bit polynomials could be a  $2n$ -bit polynomial and reduction of  $2n$ -bit polynomial to  $(X^n - 1)$  in  $\mathbb{F}_2$  is equivalent to slicing of the  $2n$ -bit polynomial into two parts of  $n$ -bit polynomials and then performing a bitwise XOR. As result, when the shift operation is performed on each chunk we also compute the address value (`ADDR_2N`) (signifying the degree of the resultant polynomial). If we notice that this degree of the resultant polynomial is greater than  $n$  we perform XOR of this chunk to the lower chunk by decoding the address based on the value of `ADDR_2N`. We perform similar operation over all the indices of the `sparse_poly` to achieve the final multiplied resultant value.

The clock cycles taken by our `poly_mult` module for one polynomial multiplication can be computed using the following formula where  $W_{SPARSE}$  is weight of the sparse polynomial,  $n$  is length of the polynomial, BW is the block width, 3 cycles represents the number of pipeline stages and 2 cycles are for the `start` and `done` synchronization with interfacing modules. The clock cycles taken for shift and interleaved reduction for one index is  $(3 + \text{ceil}(n/\text{BW}))$ . Our `poly_mult` module is constant time and we achieve that by fixing the  $W_{SPARSE}$  to a specific value ( $w$  and  $w_r$ ) based on the parameter set.

$$\text{latency}_{\text{poly\_mult}} = W_{SPARSE} \times (3 + \text{ceil}(n/\text{BW})) + 2$$

Table 3 shows the results for our `poly_mult` module compared with the related work. We note that our sparse polynomial multiplication module performs better in terms of time while utilizing half the Block RAM resources when compared to the existing designs. Table 4 shows results for our `poly_mult` module for the parameter sizes used for HQC hardware design.

### 3.1.3 Polynomial Addition/Subtraction

HQC uses polynomial addition/subtraction in all of its primitives. Since all addition and subtraction operations happen in  $\mathbb{F}_2$ , the addition and subtraction could be realized as the same operation. We design two variants of constant-time adders namely `xor_based_adder` and `location_based_adder` that could be attached with our polynomial multiplication module described in Section 3.1.2. We design our adder modules as an extension for polynomial multiplication because the addition/subtraction always appears with the polynomial multiplication as shown in Algorithm 1, Algorithm 2, and Algorithm 3. The adders operate on contents of block RAM since the polynomials are stored inside the block RAM. Both of the adder module designs do not use any additional block RAM resources,

**Table 4:** Time and area information of our `poly_mult` module for different HQC parameter sizes with different performance parameter (BW) sizes, data based on synthesis results for Artix 7 board with xc7a200t-3 FPGA chip.

Input Length <sup>†</sup> (bits)	$W_{sparse}^*$	Resources							
		Logic (LUT) (DSP)		Memory (FF) (BR)		F (MHz)	Cycles (cyc.)	Time (us)	T x A
Our <code>poly_mult</code> module (BW = 32)									
17,669 (hqc128)	66	412	0	189	1	287	36,698	0.13	53
35,851 (hqc192)	100	387	0	193	2	257	112,402	0.44	169
57,637 (hqc256)	131	397	0	199	2	267	236,457	0.89	352
Our <code>poly_mult</code> module (BW = 64)									
17,669 (hqc128)	66	620	0	245	2	270	18,482	0.07	42
35,851 (hqc192)	100	649	0	249	2	286	56,402	0.20	128
57,637 (hqc256)	131	644	0	223	2	283	118,426	0.42	269
Our <code>poly_mult</code> module (BW = 128)									
17,669 (hqc128)	66	1,439	0	496	4	238	9,374	0.04	57
35,851 (hqc192)	100	1,445	0	500	4	240	28,402	0.12	171
57,637 (hqc256)	131	1,448	0	474	4	245	59,476	0.24	352

<sup>†</sup> Length of the non-sparse polynomial, \* = Weight of the sparse polynomial input

**Table 5:** Polynomial addition modules (`xor_based_adder` and `loc_based_adder` with datapath width 128-bits) area and timing information, data based on synthesis results for Artix 7 board with xc7a200t-3 FPGA chip.

Input Length (bits)	Resources							
	Logic (LUT) (DSP)		Memory (FF) (BR)		F (MHz)	Cycles (cyc.)	Time (us)	T x A
<code>xor_based_adder</code> (BW = 128)								
17,669	143	0	159	0	330	142	0.43	0.06
35,851	142	0	161	0	318	284	0.89	0.12
57,637	142	0	161	0	311	455	1.46	0.20
<code>loc_based_adder</code>								
17,669	160	0	174	0	316	69	0.22	0.03
35,851	161	0	174	0	300	103	0.34	0.05
57,637	161	0	175	0	300	134	0.45	0.07

they load the polynomial multiplication output, perform the addition, and write the value back to the same block RAM inside the polynomial.

The `xor_based_adder` design performs addition in a regular  $\mathbb{F}_2$  fashion by performing bit-wise `exclusive-OR` operation. The module performs addition sequentially by generating one block RAM address per clock cycle to load inputs from two block RAMs and then performs addition and writes them back to one of the specified block RAMs at the same block RAM address.

The `location_based_adder` is an optimized adder designed to perform addition when one of the input is a sparse vector. This module is mainly designed to perform operations  $\mathbf{x} + \mathbf{h} \cdot \mathbf{y}$  from [Algorithm 1](#) and  $\mathbf{r}_1 + \mathbf{h} \cdot \mathbf{r}_2$  and  $\mathbf{s} \cdot \mathbf{r}_2 + \mathbf{e}$  from [Algorithm 2](#). In these operations the values of  $\mathbf{x}$ ,  $\mathbf{r}_1$ , and  $\mathbf{e}$  are sparse, fixed-weight vectors so the addition is optimized by only flipping the bits of the other input in the position of one. The `location_based_adder` module takes location of ones from the sparse vector as input and computes the address to load out the part of non-sparse polynomial from the block RAM and flips the bit on the appropriate location and writes it back to the same location. The process is repeated until all locations with ones are covered. Since there are a fixed, and known number of ones in the fixed-weight vector, there is a fixed number of operations and timing does not reveal any sensitive information. Results of our polynomial addition `location_based_adder` module for one performance parameter (width = 128) are shown in [Table 5](#).

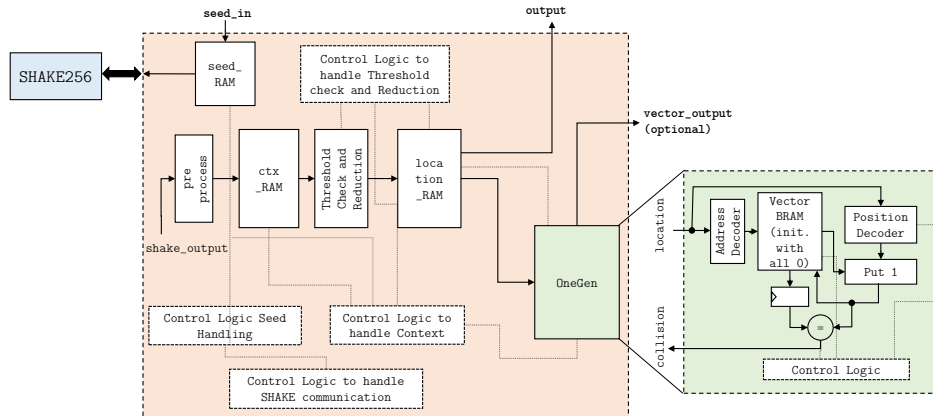


Figure 2: Hardware design of `fixed_weight_vector` module.

### 3.1.4 Fixed-Weight Vector Generator

The fixed-weight vector generator function generates a uniform random  $n$ -bit fixed-weight vector of a specified input weight ( $w$ ). The module assumes that there is a random number generator that can be used to generate uniformly random bits. The algorithm for fixed-weight generation as specified in [AAB<sup>+</sup>20] first generates  $24 \times w$  random bits. These random bits are then arranged into  $w$  24-bit integers. These 24-bit integers undergo a threshold check and are rejected if the integer value is beyond the threshold ( $949 \times 17,669$ ,  $467 \times 35,851$ ,  $291 \times 57,637$  for hqc-128, hqc-192 and hqc-256 respectively). After the threshold check, these integers are reduced modulo  $n$ . After the threshold check and reduction process if the weight is not equal to  $w$  then more random bits are drawn from RNG and the process is repeated until  $w$  integers are achieved. After the threshold check and reduction then a check for duplicates is performed over all the reduced integers. In case any duplicate is found, that integer is discarded and more random bits are requested drawn from the RNG which again undergo threshold check, reduction and duplicate check. This process is repeated until a uniform fixed weight vector is generated.

In our hardware design, we use a PRNG to generate the uniformly random bits required for the fixed weight vector generation from an input seed of length 320-bits. Our hardware design includes this PRNG in the form of SHAKE256. Our design assumes that the seed will be initialized by some other hardware module implementing a true random number generator.

The hardware design of `fixed_weight_vector` generation module is shown in Figure 2. We use SHAKE256 module described in Section 3.1.1 to expand 320-bit seed to a  $24 \times w$ -bit string. Since the SHAKE256 module has 32-bit interface the seed is loaded in 32-bit chunks and the seed is stored in `seed_RAM` as shown in the Figure 2. The 32-bit chunk from SHAKE256 is broken into 24-bit integer by `preprocess` unit and stored in the `ctx_RAM` then threshold check and reduction are performed. For the reduction, we use Barrett reduction. The Barrett reduction is optimized to reduce for specific value ( $n$ ) since that value is constant. After the reduction, the integer values are stored in the `locations_RAM`. Once the `locations_RAM` is filled the `OneGen` module is triggered. The `OneGen` module helps in detecting if there are any duplicates in the `locations_RAM`. Our `OneGen` module is inspired from duplicate checking logic described in [CCD<sup>+</sup>22]. While the `OneGen` module checks for duplicates, the SHAKE256 module generates the next  $24 \times w$ -bit string to tackle any potential duplicates and stores them in the `ctx_RAM`. This way we are able to mask any clock cycles taken for seed expansion.

The main pitfall where the fixed-weight vector generation process may show non constant-time behavior is the rejection sampling process (i.e., the threshold check and

**Table 6:** `fixed_weight_vector` module area and timing information, data based on synthesis results for Artix 7 board with `xc7a200t-3` FPGA chip and probability of failing constant-time behavior of our `fixed_weight_vector` generation module best and worst case values for the parameter, `ACCEPTABLE_REJECTIONS`.

Design	Resources					Cycles (cyc.)	Time (us)	T x A	Failure <sup>+</sup> Prob.
	Weight	Logic (LUT)	Memory (FF)	Memory (BR)	F (MHz)				
<b>non constant-time design (ACCEPTABLE_REJECTIONS = 0)</b>									
hqc128	75	240	111	2.0	226	709	3.14	0.75	$1.1 \times 2^{-11}$
hqc192	114	229	112	2.0	220	1,840	8.36	1.92	$1.1 \times 2^{-9}$
hqc256	149	234	117	2.0	228	2,106	9.24	2.16	$1.1 \times 2^{-12}$
<b>constant-time design (ACCEPTABLE_REJECTIONS = <math>w_r</math>)</b>									
hqc128	75	316	124	2.0	223	3,649	16.36	5.17	$2.8 \times 2^{-199}$
hqc192	114	295	125	2.0	236	4,200	17.80	5.25	$1.1 \times 2^{-280}$
hqc256	149	314	192	2.5	242	5,935	24.52	7.70	$4.9 \times 2^{-355}$

<sup>+</sup> = Probability of our design failing to behave constant-time.

duplicate detection as discussed earlier). A timing attack on existing software reference implementation of HQC [AAB<sup>+</sup>20] was performed in [GHJ<sup>+</sup>22]. The authors use the information of rejection sampling routine (that is part of fixed-weight generation) being invoked during the deterministic re-encryption process in decapsulation and show that this leaks secret-dependent timing information. The timing of the rejection sampling routine depends upon the given seed. This seed is derived for the encrypt function in encapsulation and decapsulation procedures using the message. The decapsulation operation is dependent on the decoded message and this dependency allows to construct a plaintext distinguisher (described in detail in [GHJ<sup>+</sup>22]) which is then used to mount the timing attack.

Our fixed-weight generation module can be parametrized to create design with arbitrarily small probability of timing attack being possible, we note that probability of absolute zero cannot be practically reached. In our hardware module, we make the constant time behavior parameterizable (parameter name is `ACCEPTABLE_REJECTIONS`). We can specify how many indices could be rejected and still the design will behave constant time (at the cost of extra area for more storage and extra cycles). The extra area is because we generate additional (based on parameter value) uniformly random bits in advance and store them in the `ctx_RAM` (shown in Figure 2). The extra clock cycles are needed because even after we found the required number of indices under the threshold value, we still go over all the `ctx_RAM` locations and, for the duplicate detection logic inside `OneGen` module (shown in Figure 2), the control logic is programmed to take the same cycles in both cases of duplicate being detected or not. The parameter can be set based on user’s target failure probability. If the actual failures are within the failure probability set by the selected parameter value, then the timing side channel given in [GHJ<sup>+</sup>22] is not possible.

The right most column in Table 6 shows probability of non constant-time design (`ACCEPTABLE_REJECTIONS = 0`) versus constant-time design (`ACCEPTABLE_REJECTIONS =  $w_r$` ) failing to behave in constant-time manner. The choice of the demonstrated parameter is made based on the values of  $w_r$  given in Table 1. We choose  $w_r$  as the parameter value for demonstration because when the rejection sampling procedure rejects the first index the `seedexpander` function in software reference implementation [AAB<sup>+</sup>20] generates  $w_r$  new set of indices and first index from the new set is used to replace the old rejected index and for the second rejection next index from the new set is chosen to replace the old rejected index, etc. This process is repeated until specified weight for the vector is achieved. To compute the failure probability (given in Table 6) for each parameter set, we take in to account both threshold check failure and duplicate detection probabilities for the respective parameter sets.

Table 6 shows the results of the fixed weight hardware module targeting Artix 7 `xc7a200t` FPGA. The area excludes `SHAKE256` module as the `SHAKE256` is shared among

**Table 7:** encode module area and timing information, data based on synthesis results for Artix 7 board with xc7a200t-3 FPGA chip.

Design	Resources							
	Logic		Memory		F	Cycles	Time	T x A
	(LUT)	(DSP)	(FF)	(BR)	(MHz)	(cyc.)	(us)	
hqc128	858	0	922	2	270.34	97	0.36	307.86
hqc192	1,011	0	1,088	2	298.32	131	0.44	443.96
hqc256	1,503	0	1,689	2	293.51	189	0.64	967.83

all primitives. The reported frequency in Table 6 is the `fixed_weight` core frequency. We discuss later our dual clock domain design (in Section 3.7) that allows the modules to run at their core frequency while SHAKE256 runs on a slower clock.

## 3.2 Encode and Decode Modules

The encode and decode modules are building blocks of the encrypt and decrypt modules, respectively. We describe the encode and decode modules here, before describing the bigger encrypt and decrypt modules in Section 3.3.

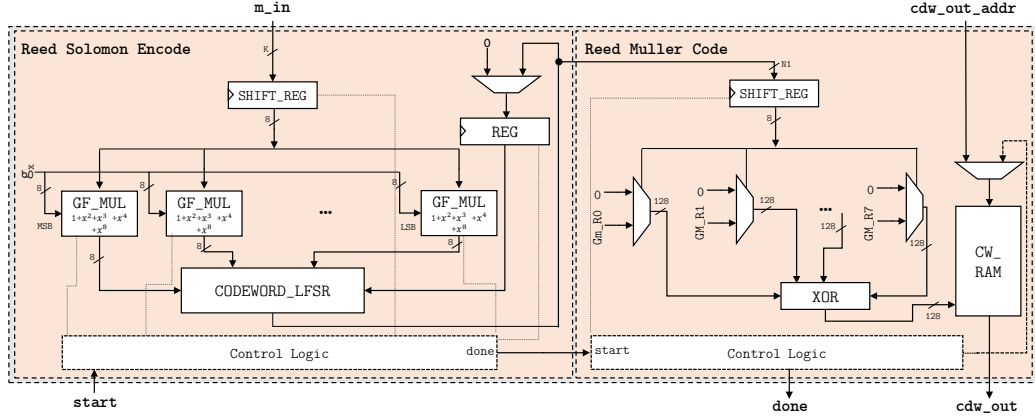
### 3.2.1 Encode Module

As specified in Section 2.2.1, HQC Encode uses concatenation of two codes namely Reed–Muller and Reed–Solomon codes. The hardware design of our `encode` module is shown in Figure 3. The Encode function takes  $K$ -bit input and first encodes it with the Reed–Solomon code. The Reed–Solomon encoding process involves systematic encoding using a linear feedback shift register (LFSR) with a feedback connection based on the generator polynomial (shown in Section 2.2.1). The Reed–Solomon code generates a  $n_1$ -bit output (as given in [AAB<sup>+</sup>20] the value for  $n_1$  is 368, 448, and 720 for `hqc128`, `hqc192`, and `hqc256` respectively). For the Galois field multiplication unit (for the field  $\mathbb{F}_2[x]/(x^8 + x^4 + x^3 + x^2 + 1)$ ) we design an LFSR-based optimized multiplication unit similar to the one described in [SR17]. The number of Galois field multiplication units we run in parallel is equal to the degree of the generator polynomial. The outputs from Galois field multipliers are fed in to a LFSR after each cycle. At the end of encoding process the module generates a  $n_1$ -bit output.

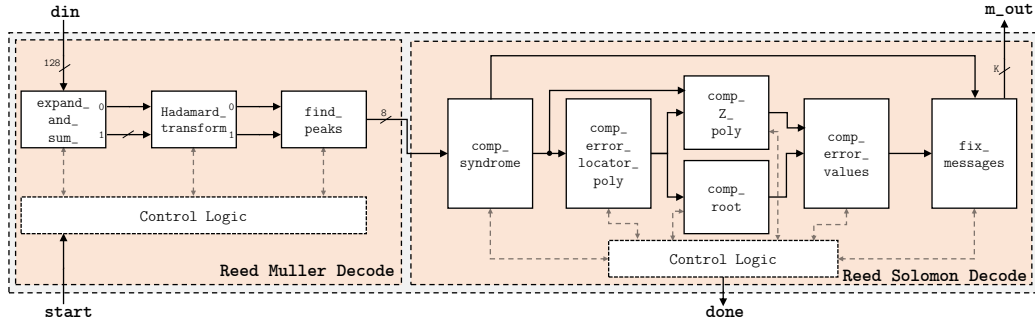
The  $n_1$ -bit output from Reed–Solomon code is then encoded by Reed–Muller code. The Reed–Muller encoding is achieved by performing vector-matrix multiplication where each byte from input is the vector and the matrix is generator matrix ( $\mathbf{G}$ ) given in Section 2.2.1. In our design we store the generator matrix rows (each row is of length 128-bits) in ROM and we select the matrix rows based on each input byte. We store the output after multiplying input byte into a block RAM in chunks of 128-bits. Based on the security parameter set the code word output from Reed–Muller code has a multiplicity value (i.e., number of times a code word or in our case number of times each block RAM location is repeated). As per the specification [AAB<sup>+</sup>20], `hqc128` has multiplicity value of 3 and `hqc192` and `hqc256` have multiplicity value of 5. To optimize the storage, we only store one copy of code word, and while accessing the code word we compute the block RAM address in a way that the multiplicity is achieved. The time and area results for our `encode` module targeting Artix 7 board with xc7a200t-3 FPGA are shown in Table 7.

### 3.2.2 Decode Module

As introduced in Section 2.2.1, the ciphertext is first decoded with duplicated Reed-Muller code and then with shortened Reed-Solomon code. To decode duplicated Reed-Muller code, the `transformation` module expands and adds multiple code words into expanded



**Figure 3:** Hardware design of **encode** module (formed by concatenating two encode functionalities, Reed-Solomon on the the left side and Reed-Muller on the right side).



**Figure 4:** Hardware design of **decode** module (formed by concatenating two decode functionalities, Reed-Muller on the the left side and Reed-Solomon on the right side).

code word, and then the Hadamard transformation is applied to the expanded code word. Finally, Find\_Peak module finds the location of the highest absolute value of the Hadamard\_Transformation output. Figure 5 describes detailed hardware design of Reed-Muller Decoder. `expand_and_sum` module collects data inputs into  $m \times 128$ -bit shift register, then add and shift the last 2-bit lsb of each shift register to produce a pair of data outputs. The data pair is then processed in `hadamard_transformation` module which consist of 7 layers of similar blocks of radix-2 butterfly structure. With the outputs from `hadamard_transformation` coming in pairs, finding peak can be done in parallel inside the `Find_Peak` module and compare the peaks of each to be the final peak. The whole processes then repeated  $n_1$  times to produce  $n_1$  data output to Reed-Solomon Decoder.

To decode Reed-Solomon code, we need to sequentially compute syndromes  $S_i$ , coefficients  $\sigma_i$  of error location polynomial  $\sigma(x)$ , roots of error location polynomial  $(\alpha^i)^{-1}$ , pre-defined helper polynomial  $Z((\alpha^i)^{-1})$ , errors  $e_i$ , and finally correct the output of decode of Reed-Muller code based on the errors.

### 3.3 Encrypt and Decrypt Modules

The encrypt and decrypt modules are building blocks of the encapsulation and decapsulation modules, respectively. We describe the encrypt and decrypt modules here, before describing the bigger encapsulation and decapsulation modules later.

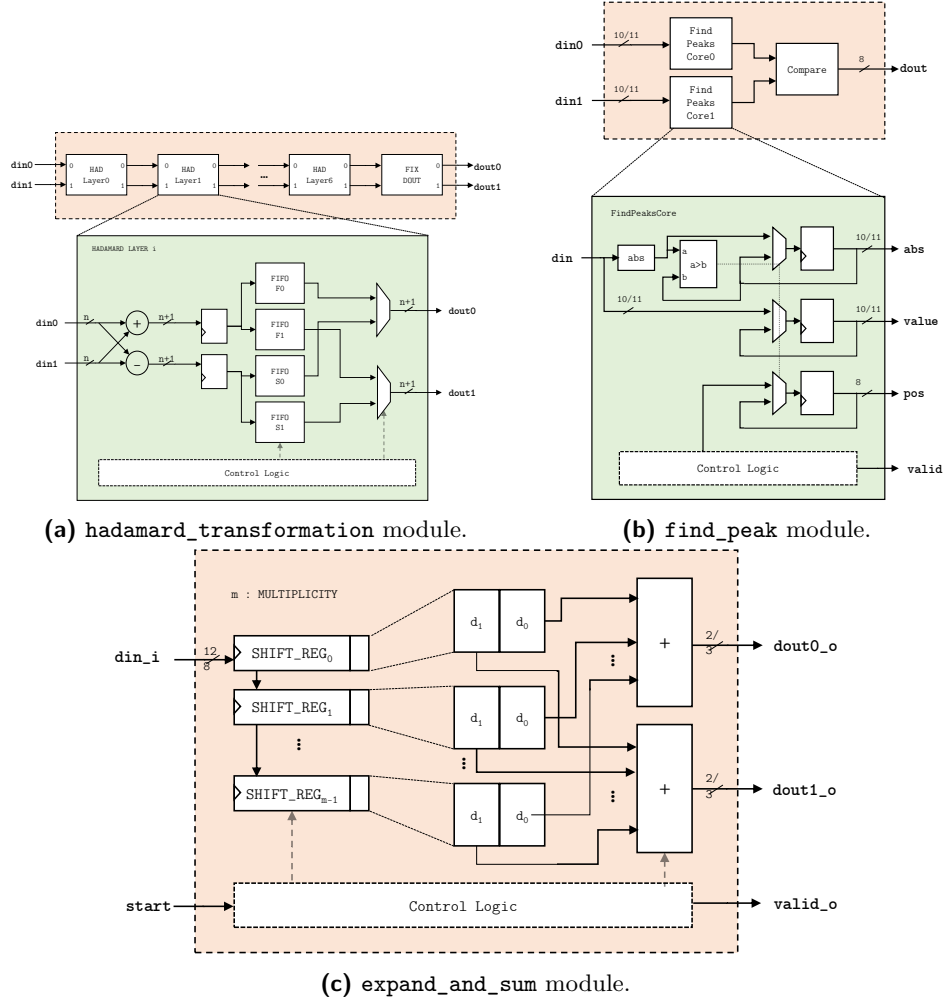


Figure 5: Hardware design of Reed-Muller Decoder.

### 3.3.1 Encrypt Module

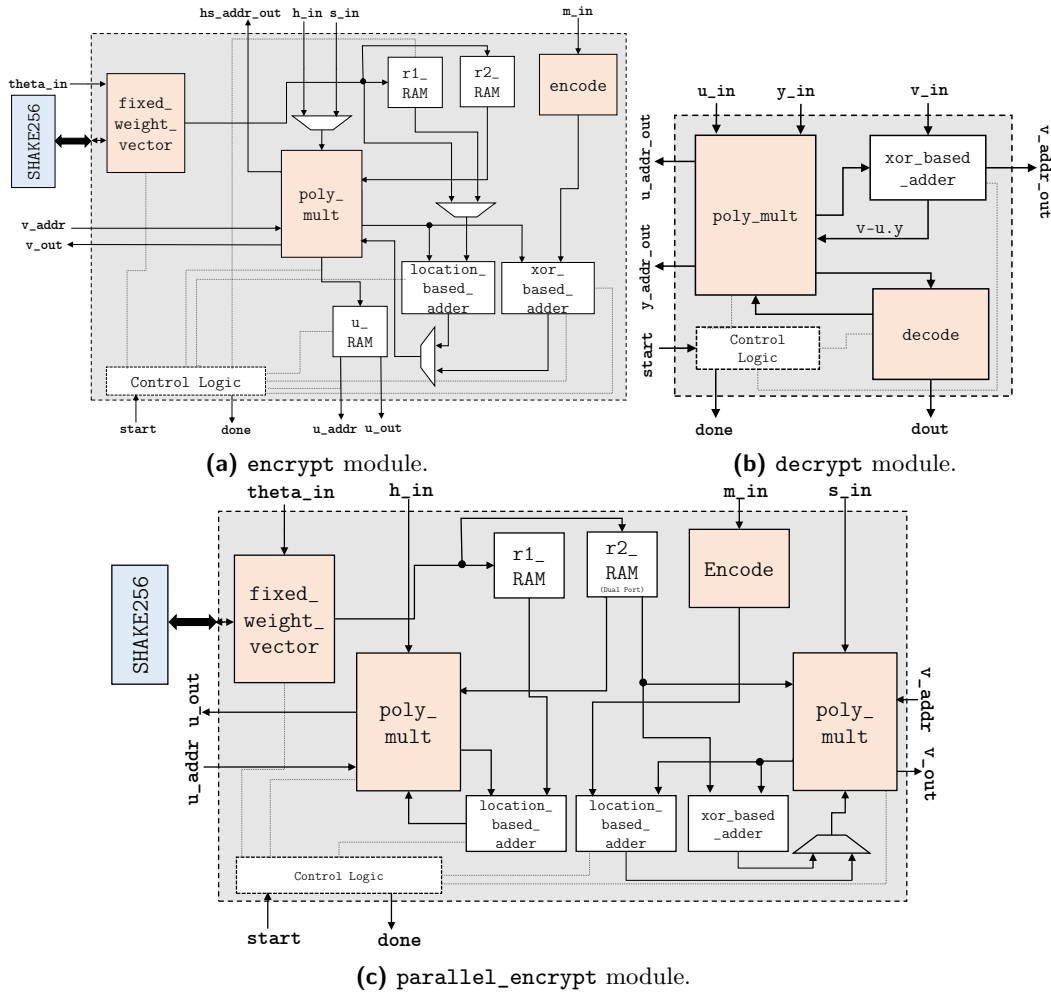
The `encrypt` module (shown in Algorithm 2) takes public key  $(\mathbf{h}, \mathbf{s})$ , message  $\mathbf{m}$ , and seed  $(\theta)$  and generates a ciphertext  $(\mathbf{u}, \mathbf{v})$  as the output. The hardware design for the `encrypt` module is shown in Figure 6a. We use the constant-time version of our `fixed_weight_vector` module with `ACCEPTABLE_REJECTIONS = w_r` described in Section 3.1.4 to generate  $\mathbf{r}_1$ ,  $\mathbf{r}_2$ , and  $\mathbf{e}$  fixed-weight vectors of weight  $w_r$  by expanding `theta_in` and in parallel we run `encode` module (described in Section 3.2.1). After the generation of  $\mathbf{r}_2$  we start the polynomial multiplication of  $\mathbf{h} \cdot \mathbf{r}_2$  in parallel to the  $\mathbf{e}$  generation. For polynomial multiplication we use the `poly_mult` module with `BW = 128` described in Section 3.1.2. The addition of  $\mathbf{r}_1$  in  $\mathbf{u}$  computation and  $\mathbf{e}$  in  $\mathbf{v}$  computation is performed by our `location_based_adder` and addition with  $\mathbf{t}$  is performed by `xor_based_adder` (described in Section 3.1.3).

From Algorithm 2, we observe that both  $\mathbf{h} \cdot \mathbf{r}_2$  and  $\mathbf{s} \cdot \mathbf{r}_2$  multiplications can be performed in parallel, consequently we design a `parallel_encrypt` module targeting higher performance where we use two polynomial multiplications in parallel (shown in Figure 6c). We provide a choice of using either `encrypt` or `parallel_encrypt` module as a parameter. Table 9 shows results for both the `encrypt` hardware implementations targeting Xilinx



**Table 8:** decode module area and timing information, data based on synthesis results for Artix 7 board with xc7a200t-3 FPGA chip.

Design	Resources							T x A
	Logic (LUT)	DSP	Memory (FF)	(BR)	F (MHz)	Cycles (cyc.)	Time (ms)	
hqc128	2,817	0	3,779	2.5	205	4,611	0.02	63
hqc192	3,257	0	4,727	2.5	212	5,485	0.03	84
hqc256	3,679	0	5,574	2.5	206	9,199	0.04	164

**Figure 6:** Hardware design of encrypt, parallel\_encrypt, and decrypt modules.

Artix 7 xc7a200t FPGA. We note that the major contributor to the overall time in encrypt operation is due to polynomial multiplication and using two `poly_mult` modules in parallel reduces the overall time by 40-60% across different parameter sets. The area results do not include the SHAKE256 module as the SHAKE256 is shared among all primitives. Figure 4 shows the hardware block design on our module and Table 8 shows the time and area results for our module.



**Table 9:** encrypt module area and timing information, data based on synthesis results for Artix 7 board with xc7a200t-3 FPGA chip.

Design	Resources <sup>†</sup>						T x A
	Logic (LUT)	DSP	Memory (FF) (BR)		F (MHz)	Cycles (cyc.)	
encrypt module – uses one poly_mult module with with BW = 128							
hqc128	2,108	0	1,611	10	233	26,418	0.11 239
hqc192	3,936	0	1,882	13	196	72,870	0.37 1,463
hqc256	2,855	0	2,363	10.5	232	146,363	0.63 1,801
parallel_encrypt module – two poly_mult modules with BW = 128 running in parallel							
hqc128	5,055	0	2,041	12.5	213	15,403	0.07 365
hqc192	5,511	0	2,292	15	193	39,838	0.21 1,137
hqc256	5,652	0	2,779	13	210	77,736	0.37 2,092

<sup>†</sup> = Given resources does not include the area for SHAKE256 module.

**Table 10:** decrypt module area and timing information, data based on synthesis results for Artix 7 board with xc7a200t-3 FPGA chip.

Design	Resources <sup>†</sup>						T x A
	Logic (LUT)	DSP	Memory (FF) (BR)		F (MHz)	Cycles (cyc.)	
hqc128	6,352	0	5,730	10.5	194	14,198	0.07 464
hqc192	7,038	0	6,787	10.5	187	34,313	0.18 1,291
hqc256	8,544	0	8,740	13	186	69,356	0.37 3,185

<sup>†</sup> = Given resources does not include the area for SHAKE256 module.

### 3.3.2 Decrypt Module

The decrypt module (shown in Algorithm 3) takes secret key  $(\mathbf{x}, \mathbf{y})$ , ciphertext  $(\mathbf{u}, \mathbf{v})$ , and generates the message  $(\mathbf{m}')$ . Figure 6b shows our hardware design for decrypt module. The module accepts part of the secret key  $(\mathbf{y})$  as locations with ones (since it is a sparse fixed weight vector). We use our poly\_mult module with BW = 128 described in Section 3.1.2 to compute  $\mathbf{u} \cdot \mathbf{y}$  and use xor\_based\_adder module (described in Section 3.1.3) to compute  $\mathbf{v} - \mathbf{u} \cdot \mathbf{y}$ . We then use the decode module (described in Section 3.2.2) to decode  $\mathbf{v} - \mathbf{u} \cdot \mathbf{y}$  and retrieve the message. Table 10 shows our hardware implementation results for decrypt module targeting Xilinx Artix 7 xc7a200t FPGA.

## 3.4 Key Generation

We now begin to describe the top-level modules, starting with the key generation, following in later sections with encapsulation and decapsulation. The discussion here focuses on single clock domain design. The novel dual clock domain design which allows the top-level modules run at different frequencies from the SHAKE256 module that they depend on is described in Section 3.7.

The key generation (shown in Algorithm 1) takes secret key seed and public key seed as an input and generates secret key  $(\mathbf{x}, \mathbf{y})$  and public key  $(\mathbf{h}, \mathbf{s})$  respectively as output. Figure 7 shows the hardware design of our keygen module. Our keygen module assumes that the public key seed and the secret key seed are generated by some other hardware module implementing a true random number generator. We use the constant-time version fixed\_weight\_vector module with ACCEPTABLE\_REJECTIONS =  $w$  described in Section 3.1.4 to generate  $(\mathbf{x}, \mathbf{y})$  from the secret key seed.  $\mathbf{x}$  and  $\mathbf{y}$  are fixed weight vectors of weight  $w$  and length  $n$ -bits. To optimize the storage, rather than storing full  $n$ -bit sparse vector we only output locations of ones. There is also an optional provision to output the full vector as described in Section 3.1.4. The vector\_set\_random uses the SHAKE256 module to expand the public key seed and generates  $\mathbf{h}$ . We then use poly\_mult module with BW = 128 (described in Section 3.1.2) to compute  $(\mathbf{h} \cdot \mathbf{y})$  and finally

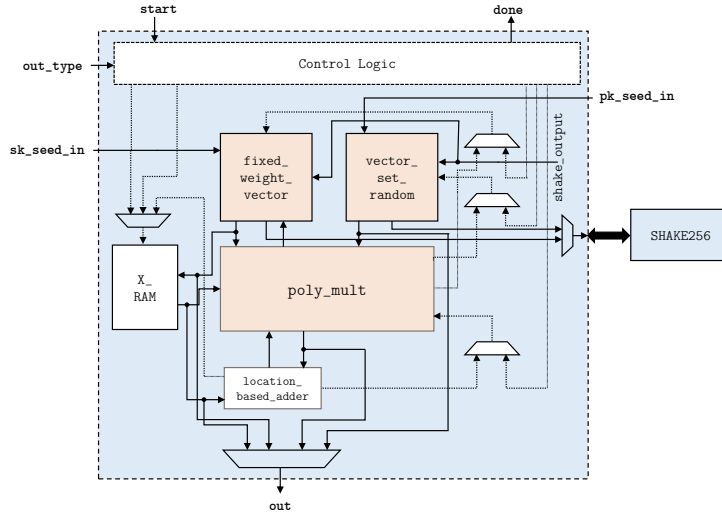


Figure 7: Hardware design of the keygen module.

Table 11: keygen module area and timing information, data based on synthesis results for Artix 7 board with xc7a200t-3 FPGA chip. Data for single clock design, using slow clock due to hash module critical path.

Design	Resources <sup>†</sup>				F (MHz)	Cycles (cyc.)	Time (ms)	T x A
	Logic (LUT)	DSP	Memory (FF)	Memory (BR)				
hqc128	2,561	0	1,116	5.5	164	15,752	0.09	245
hqc192	1,630	0	866	9.5	164	39,580	0.24	393
hqc256	1,725	0	891	9.5	164	76,111	0.46	800
hqc128-perf HLS*[AAB <sup>+</sup> 20]	12,000	0	9,000	3	150	40,000	0.27	3,240
hqc128-comp. HLS*[AAB <sup>+</sup> 20]	4,700	0	2,700	3	129	630,000	4.80	22,560
hqc128-optimized. HLS*[DNN <sup>+</sup> 22]	11,484	0	8,798	6	150	40,427	0.27	3,095
hqc128-pure HLS*[DNN <sup>+</sup> 22]	24,746	0	21,746	7	153	40,427	0.27	6,539

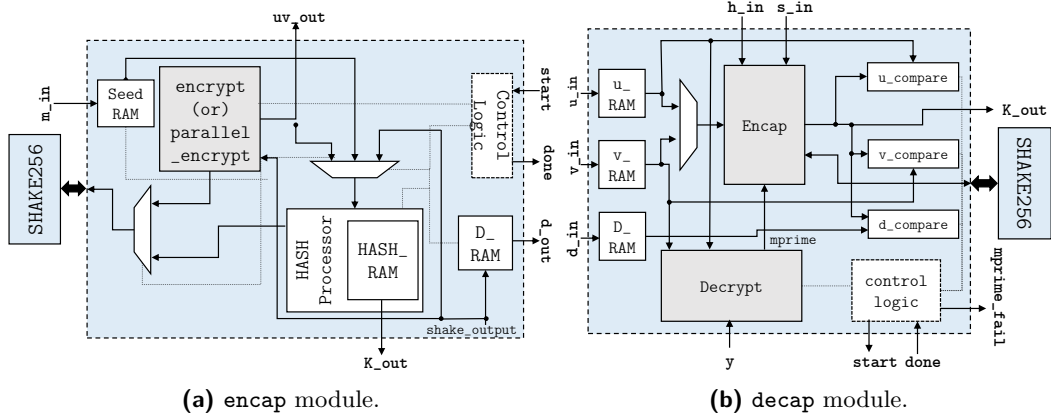
<sup>†</sup> = Given resources does not include the area for SHAKE256 module, \* = Target FGPA is Artix-7 xc7a100t-1

use `location_based_adder` module (described in Section 3.1.3) to compute  $\mathbf{s}$ . We note that in the Figure 7 only a block RAM for  $\mathbf{x}$  storage (`X_RAM`) is visible because the  $\mathbf{y}$ ,  $\mathbf{h}$ ,  $\mathbf{s}$  are stored in the block RAMs which are inside `fixed_weight_vector`, `poly_mult`, and `location_based_adder` modules respectively.

Table 11 shows the results for the `keygen` module. We note that the maximum clock frequency of `keygen` module alone (without `SHAKE256`) module is in range 250-259 MHz based on the parameter set selected, but since the critical path lies inside the `SHAKE256` module we report `SHAKE256`'s frequency in the Table 11. The area results do not include the `SHAKE256` module because it is shared among all other primitives. Our dual clock domain design can be applied to the key generation to run the key generation and `SHAKE256` at two different frequencies (discussed in Section 3.7).

### 3.5 Encapsulation Module

The encapsulate operation (shown in Algorithm 4) takes public key  $(\mathbf{h}, \mathbf{s})$  and message  $\mathbf{m}$  as an input and generates shared secret  $(K)$  and ciphertext  $(\mathbf{c} = (\mathbf{u}, \mathbf{v}))$  and  $\mathbf{d}$ . The hardware design of the `encap` module is shown in Figure 8a. Our `encap` module assumes that  $\mathbf{m}$  is generated by some other hardware module implementing a true random number generator and provided as an input to our module. Since the `SHAKE256` module is extensively



**Figure 8:** Hardware design of `encap` and `decap` modules.

**Table 12:** `encap` module area and timing information, data based on synthesis results for Artix 7 board with `xc7a200t-3` FPGA chip. Design with single clock, using the slow clock due to hash module critical path.

Design	Resources <sup>†</sup>				Cycles (cyc.)	Time (ms)	T x A
	Logic (LUT)	DSP	Memory (FF)	Memory (BR)			
our <code>encap</code> module with <code>encrypt</code>							
<code>hqc128</code>	2,628	0	1,965	13	164	34,521	0.21 553
<code>hqc192</code>	4,318	0	2,328	18	164	89,065	0.54 2,345
<code>hqc256</code>	3,295	0	2,850	15.5	164	172,126	1.05 3,458
our <code>encap</code> module with <code>parallel_encrypt</code>							
<code>hqc128</code>	5,545	0	2,393	15.5	164	23,506	0.14 794
<code>hqc192</code>	6,678	0	2,737	20	164	56,033	0.34 2,281
<code>hqc256</code>	6,632	0	3,265	18	164	103,499	0.63 4,185
<code>hqc128 perf HLS*[AAB+20]</code>	16,000	0	13,000	5.0	151	89,000	0.59 9,440
<code>hqc128 comp. HLS*[AAB+20]</code>	6,400	0	4,100	5.0	127	1,500,000	12.00 76,800
<code>hqc128 optimized HLS*[DNN+22]</code>	16,487	0	13,390	10.0	152	89,110	0.59 9,665
<code>hqc128 pure HLS*[DNN+22]</code>	29,496	0	26,333	11.0	148	89,131	0.59 17,764

<sup>†</sup> = Given resources does not include the area for `SHAKE256` module, \* = Target FPGA is Artix-7 `xc7a100t-1`

used in encapsulate operation we design a `HASH_processor` module which handles all the communication with the `SHAKE256` module. `HASH_processor` modules reduces the multiplexing logic of inputs to the `SHAKE256` module significantly.

The `Hash_processor` modules helps in expanding  $\mathbf{m}$  to generate  $\theta$ . We then use our `encrypt` module (described in Section 3.3.1) to encrypt  $\mathbf{m}$  using  $\theta$  and the public key as inputs and generates ciphertext. After the generation of  $\mathbf{r}_1$ ,  $\mathbf{r}_2$ , and  $\mathbf{e}$  inside the `encrypt` module (described in Section 3.3.1) we then run `HASH_processor` module in parallel to `encrypt` module to generate  $\mathbf{d}$ . After the encryption of  $\mathbf{m}$  we then use the `HASH_processor` to compute  $\mathcal{K}(\mathbf{m}, \mathbf{c})$  to generate the shared secret  $K$ . Our design is constant-time since all the underlying modules are constant-time and the control logic from the `encap` module does not depend on any secret input.

Table 12 shows the results for the `encap` module with our `encrypt` and `parallel_encrypt`. The maximum clock frequency of `encap` module alone (without `SHAKE256` module) is in range 196-232 MHz for the design with `encrypt` and 192-213 MHz for the design with `parallel_encrypt` based on the parameter set selected but since the critical path lies inside the `SHAKE256` module we report that frequency in the Table 12. The area results do not include the `SHAKE256` module because it is shared among all other primitives.



**Table 14:** Designs with dual clocks, using slow clock for hash module, and fast clock for remainder of the design. `keygen`, `encap`, and `decap` modules area and timing information, data based on synthesis results for Artix 7 board with `xc7a200t-3` FPGA chip.

Design	Resources <sup>†</sup>								Time (ms)	T x A
	Logic (LUT) (DSP)		Memory (FF) (BR)		F <sub>C</sub> (MHz)	F <sub>S</sub> (MHz)	Cyc <sub>C</sub> (cyc.)	Cyc <sub>S</sub> (cyc.)		
	<b>keygen</b>									
hqc128	2,611	0	905	10.5	250	164	11,821	7,958	0.10	250
hqc192	2,568	0	917	11.5	259	164	32,405	14,445	0.21	547
hqc256	2,579	0	939	12	250	164	65,082	22,153	0.40	1,020
	<b>encap with encrypt module</b>									
hqc128	2,498	0	1,975	13	222	164	25,484	17,935	0.22	560
hqc192	4,199	0	2,338	18	196	164	72,263	33,452	0.57	2,404
hqc256	3,152	0	2,860	15.5	232	164	146,444	51,209	0.94	2,972
	<b>encap with parallel_encrypt module</b>									
hqc128	5,532	0	2,403	15.5	213	164	14,470	17,935	0.18	981
hqc192	6,526	0	2,747	20	192	164	39,231	33,452	0.41	2,662
hqc256	6,090	0	3,276	18	206	164	77,818	51,209	0.69	4,207
	<b>decap with encrypt module</b>									
hqc128	8,142	0	6,455	20	194	164	40,254	17,935	0.32	2,581
hqc192	10,373	0	7,837	25	178	164	107,721	33,452	0.81	8,393
hqc256	9,598	0	9,243	22.5	186	164	217,628	51,209	1.48	14,218
	<b>decap with parallel_encrypt module</b>									
hqc128	12,623	0	6,892	22.5	194	164	29,240	17,935	0.26	3,284
hqc192	12,463	0	8,246	27	178	164	74,689	33,452	0.62	7,772
hqc256	13,004	0	9,652	25	186	164	149,002	51,209	1.11	14,469

F<sub>C</sub> = Core Frequency, F<sub>S</sub> = SHAKE256 Frequency, Cyc<sub>C</sub> = Core Cycles, Cyc<sub>S</sub> = SHAKE256 Cycles, † = Given resources does not include the area for SHAKE256 module

### 3.7 Dual Clock Design

By profiling our `keygen`, `encap`, and `decap` modules we observe that the maximum number of clock cycles in all the operations are taken by our `poly_mult` module (described in Section 3.1.2) and we note that (as specified in Section 3.1.4) the maximum clock frequency of `keygen`, `encap`, `decap` modules is limited by the SHAKE256 module since the critical path lies inside the round function of SHAKE256. To optimize the time taken by the time consuming modules such as `poly_mult` whose frequency is higher than that of SHAKE256 module, we implement a `core_wrapper` module (shown in Figure 9) that has a capability to support two asynchronous clocks.

The `core_wrapper` block represented in Figure 9 can be any of the our modules that needs to interface with the SHAKE256 module, i.e. key generation, encapsulation or decapsulation module. In the `core_wrapper` design we use two FIFOs (generated using the Xilinx IP generator) one in each direction (`Core_to_SHAKE256_FIFO`, `SHAKE256_to_Core_FIFO`). We compute the depth of the FIFO considering the worst case scenario and `CoreClock` to be of higher frequency than `SHAKE256Clock`. For `Core_to_SHAKE_FIFO` we note the depth be 36 (with width of each location to be 32). Out of 36 locations first two block are the SHAKE256 module commands describing input and output width and rest of  $34 \times 32 = 1088$ -bits represents the block size of the SHAKE256. We select FIFO width to be 32 since the SHAKE256 module has a 32-bit interface (as described in Section 3.1.1). For `SHAKE_to_CORE_FIFO` we compute the depth of the FIFO to be 1. For the `force_done` and `force_done_ack` signals we use a `Dual_Flop_Synchronizer`.

Table 14 shows the results for our dual clock hardware designs of `keygen`, `encap`, `decap` modules. We note that there is some overhead in terms of additional clock cycles when moving the data through FIFOs but the overhead is overcome through the higher clock frequency of the core. The area results do not include the SHAKE256 module as the SHAKE256 is shared among all primitives. Overall we note that the dual clock design is helpful when either the difference in core cycles and SHAKE cycles is higher or the difference between the SHAKE256 module and `core` module is higher or both (e.g., `keygen`

**Table 15:** Comparison of the time and area for our HQC hardware design with the related work.

Design	Resources										
	Logic		Memory		F	Encap		Decap		KeyGen	
	(LUT)	(DSP)	(FF)	(BR)	(MHz)	(Mcy.)	(ms)	(Mcy.)	(ms)	(Mcy.)	(ms)
<b>Security Level 1 — Classical 128-bit Security</b>											
<b>HQC – Our Work, HDL design, Artix 7 (xc7a200t)</b>											
<i>Balanced – SC</i>	16,389	0	9,762	38.5	164	0.03	0.21	0.05	0.30	0.02	0.10
<i>Balanced – DC</i>	16,070	0	9,571	43.5	194	0.04	0.24	0.06	0.32	0.02	0.11
<i>HighSpeed – SC</i>	23,680	0	10,624	43.5	164	0.02	0.14	0.04	0.23	0.02	0.10
<i>HighSpeed – DC</i>	23,585	0	10,436	48.5	194	0.03	0.18	0.05	0.26	0.02	0.11
<b>HQC – [AAB<sup>+</sup>20], HLS design, Artix 7 (xc7a100t)</b>											
<i>LightWeight</i>	8,900	0	6,400	14.0	132	1.50	12.00	2.10	16.00	0.63	4.80
<i>HighSpeed</i>	20,000	0	16,000	12.5	148	0.09	0.60	0.19	1.20	0.04	0.30
<b>HQC – [DNN<sup>+</sup>22], HLS design, Artix 7 (xc7a100t)</b>											
<i>optimized</i>	32,423	0	10,084	25	150	0.09	0.59	0.19	1.29	0.04	0.27
<i>pure</i>	32,398	0	10,068	26	148	0.09	0.60	0.19	1.30	0.04	0.27

*SC* = SingleClockDomain, *DC* = DualClockDomain, *FF* = flip-flop, *F* =  $F_{\max}$ , *BR* = BRAM,

module results from Table 14). If not, the extra clock cycles taken for synchronization and the extra area over head would deteriorate the overall performance of the design (e.g., *encap* and *decap* modules' results from Table 14).

## 4 Related Work

This section presents related work, focusing on full hardware designs of the four fourth-round public-key encryption and key-establishment algorithms in NIST's standardization process: BIKE, Classic McEliece, HQC, and SIKE. We also include the CRYSTALS-Kyber which is a public-key encryption and key-establishment algorithm selected for standardization at the end of the prior third-round. Due to limited space, related work on software implementations is omitted.

### 4.1 Comparison to Existing HQC Hardware Designs

For most other designs there is usually a light-weight and high-speed version. For our design we also present multiple variants: *Balanced* and *HighSpeed* with single clock domain (abbreviated *SC*) and *Balanced* and *HighSpeed* with dual clock domain (abbreviated *DC*). The main difference between our *Balanced* and *HighSpeed* designs is, our *Balanced* designs uses the regular `encrypt` module (shown in Figure 6a) and our *HighSpeed* design uses the `parallel_encrypt` module (shown in Figure 6c) for performing the encryption and re-encryption operations in encapsulation and decapsulation operations. For our *SC* designs, the resources are the sum of all the resources used by the key generation, encapsulation, decapsulation and shared hash module. The *SC* designs' frequency is limited by the hash module frequency. For our *DC* designs, we also sum the resources. The resource usage increases due to the asynchronous FIFOs used to bridge the two clock domains and associated control logic. The cycles increase due to the extra cycles waiting for FIFOs to be filled with data before being read. However, we note that if most of the cycles are spent in the faster clock domain, the overall times are reduced.

As a comparison, hardware design for HQC has been previously reported in [AAB<sup>+</sup>20]. The design was generated using high-level synthesis (HLS) as opposed to hand-written HDL code. The code can be generated to obtain performance numbers: 0.3ms for key generation, 0.6ms for encapsulation, and 1.2ms for decapsulation, the times correspond to the high-speed implementation of the lowest security level. Our design is faster for all three operations. Authors also provide light-weight version for the lowest security level, but did not provide hardware designs for other levels. Our implementation covers all

**Table 16:** Comparison of the time and area of hardware designs of other (NIST PQC competition) round 4 KEM candidates.

Design	Resources										
	Logic		Memory		F	Encap		Decap		KeyGen	
	(LUT)	(DSP)	(FF)	(BR)	(MHz)	(Mcycc.)	(ms)	(Mcycc.)	(ms)	(Mcycc.)	(ms)
<b>Security Level 1 — Classical 128-bit Security</b>											
<b>BIKE – [RBMG22], HDL design, Artix 7 (xc7a35t)</b>											
<i>LightWeight</i>	12,868	7	5,354	17.0	121	0.20	1.2	1.62	13.3	2.67	21.9
<i>HighSpeed</i>	52,967	13	7,035	49.0	96	0.01	0.1	0.19	1.9	0.26	2.7
<b>BIKE – [RBCGG21], HDL design, Artix 7 (xc7a200t)</b>											
<i>LightWeight</i>	12,319	7	3,896	9.0	121	0.05	0.4	0.84	6.89	0.46	3.8
<i>TradeOff</i>	19,607	9	5,008	17.0	100	0.03	0.3	0.42	4.2	0.18	1.9
<i>HighSpeed</i>	25,549	13	5,462	34.0	113	0.01	0.1	0.21	1.9	0.19	1.7
<b>Classic McEliece – [CCD<sup>+</sup>22], HDL design, Artix 7 (xc7a200t)</b>											
<i>LightWeight</i>	23,890	5	45,658	138.5	112	0.13	1.1	0.17	1.5	8.88	79.2
<i>HighSpeed</i>	40,018	4	61,881	177.5	113	0.03	0.3	0.10	0.9	0.97	8.6
<b>SIKE – [MLRB20], HDL design, Artix 7 (xc7a100t)</b>											
<i>LightWeight</i>	11,943	57	7,202	21	145	—	25.6	—	27.2	—	15.1
<i>HighSpeed</i>	22,673	162	11,661	37	109	—	15.3	—	16.3	—	9.1
<b>Kyber – [JGCS21], HDL design, (xc7a35t-2)</b>											
<i>CB</i>	5,269	2	2,422	6	—	0.67	2.67	0.73	2.93	0.69	2.75
<i>RB</i>	7,151	2	2,422	5	—	0.03	0.10	0.03	0.12	0.04	0.15
<b>Kyber – [DMG21], HDL design, (xc7a200t)</b>											
<i>HighSpeed</i>	9,457	4	8,543	4.5	220	0.003	0.01	0.004	0.02	0.002	0.01
<b>Kyber – [XL21], HDL design, (xc7a12t-1)</b>											
<i>Balanced</i>	7,412	2	4,644	3	161	0.005	0.23	0.006	0.04	0.003	0.02

*CB* = CoProcessorBased, *RB* = RoundBased, *FF* = flip-flop, *F* =  $F_{\max}$ , *BR* = BRAM

three security levels. The authors provide code to generate VHDL implementation, for an Artix-7, from the HLS-compatible sources.<sup>2</sup> A different HLS based design has been presented in [DNN<sup>+</sup>22], which achieves similar results to [AAB<sup>+</sup>20], and is also slower than our design.

The area and timing of the different HQC hardware designs is listed in Table 15. The table compares our HQC design to other existing HQC HLS designs from the literature. Our data is from synthesis reports, while data for the other algorithms is from the cited papers.

## 4.2 Comparison to Hardware Designs for Other Round 4 Algorithms

We also provide Table 16 where we tabulate latest hardware implementations of all other post-quantum cryptographic algorithm hardware implementations from the fourth round of NIST’s standardization process, plus the to-be standardized Kyber algorithm. We focus on comparison of the hardware designs for lowest level of security, Level 1, as all publications give clear time and area numbers. Majority of related work provides hardware designs for more than the lowest security level, but the timing and area numbers are not clearly broken down in the respective publications, so we focus only on comparing among the lowest security level designs.

Among the other existing designs, a hardware design for BIKE has been presented in [RBMG22]. The work investigated different strategies to efficiently implement the BIKE algorithm on FPGAs. The authors improved already existing polynomial multipliers, proposed efficient designs to realize polynomial inversions, and implement the Black-Gray-Flip (BGF) decoder. The authors provided VHDL designs for key generation, encapsulation, and decapsulation. For the fastest design, the authors showed 2.7ms for the key generation, 0.1ms for the encapsulation, and 1.9ms for the decapsulation, the times correspond to the

<sup>2</sup><https://pqc-hqc.org/implementation.html>



high-speed implementation for the lowest security level. The authors also provide data for light-weight implementation for the lowest security level. Their paper further discusses Level 3 parameters for BIKE, but does not give final hardware data for the that security level. The authors provide free, non-commercial license for the hardware code.<sup>3</sup>

Another BIKE hardware implementation in [RBCGG21] provides similar results but at much smaller area for their high-speed version. Two key arithmetic components from BIKE, polynomial multiplication and inversion were improved by implementing a sparse polynomial multiplier and extended euclidean algorithm based inversion unit due to which substantial amount of improvement was seen in all the primitives. The authors provided verilog designs for key generation, encapsulation, and decapsulation and they are available free under non-commercial license.<sup>4</sup>

Apart from earlier mentioned BIKE hardware implementations, [GGM<sup>+</sup>22] presents a practical approach of client (decapsulation and keygen) - server(Encapsulation) based model for generating the shared secret using BIKE. The presented design outperforms [RBMG22] in terms of time for all primitives but has significantly larger area footprint. [MGCZ22] presents a comparison between pure software implementation, pure HLS design, and HLS based HW/SW codesign for BIKE. However the performance of any of these design do not outperform the performance results tabulated in Table 16.

Classic McEliece has been most recently implemented in [CCD<sup>+</sup>22]. This is the first complete implementation of Classic McEliece KEM. The design provided Verilog code for encapsulation and decapsulation modules as well as key generation module with seed expansion. The authors presented three new algorithms that can be used for systemization of the public key matrix during key generation. The authors showed that the complete Classic McEliece design can perform key generation in 8.6ms, encapsulation in 0.3ms, and decapsulation in 0.9ms, the times correspond to the high-speed implementation for the lowest security level. The authors also provide hardware implementation for other security levels, and light-weight and high-speed versions for all the levels. The authors provide open-source code for the hardware.<sup>5</sup> Apart from the earlier implementation, [ZZC<sup>+</sup>22] presented a high-throughput and compact key generation module. The authors presented improvements in Gaussian elimination, sorting unit and other hardware optimizations such as algorithm level pipelining. Overall, 11% reduction in clock cycles, 31% reduction in BRAM utilization was observed with 11% increase in the logic utilization.

Hardware implementation of SIKE has been provided in [MLRB20]. The authors created VHDL implementation of SIKE as a hardware co-processor. Their design can realize any of the SIKE security levels. For the high-speed design for the lowest security level, authors report the time for encapsulation, decapsulation, and keygen as 15.3ms, 16.3ms, and 9.1ms respectively. The authors make the code available under Creative Commons public domain license.<sup>6</sup>

Different hardware implementations of CRYSTALS-Kyber are available in [JGCS21, DMG21, XL21]. The authors presented designs configurable for different performance, area requirements, and parameter sets. The high-speed design provided in [DMG21] outperforms all other algorithms in terms of time for key generation, encapsulation, and decapsulation. For the lowest security level, the authors reported 0.02ms for key generation, 0.03ms for encapsulation, and 0.04ms for decapsulation. The authors did not provide access to the code for their hardware design.

<sup>3</sup><https://github.com/Chair-for-Security-Engineering/BIKE>

<sup>4</sup><https://github.com/Chair-for-Security-Engineering/RacingBIKE>

<sup>5</sup><https://caslab.csl.yale.edu/code/pqc-classic-mceliece/>

<sup>6</sup><https://github.com/pmassolino/hw-sike>



## 5 Conclusion

This work presented hardware design for constant-time implementation of the HQC code-based key encapsulation mechanism. This work presented first, hand-optimized design of HQC key generation, encapsulation, and decapsulation written in Verilog targeting implementation on FPGAs. The three modules further share a common SHAKE256 hash module to reduce area overhead. The architecture of the hardware modules included novel, dual clock domain design, allowing the common SHAKE module to run at slower clock speed compared to the rest of the design, while other faster modules run at their optimal clock rate. The design currently outperforms the other hardware designs for HQC, and is competitive with other fourth-round Post-Quantum Cryptography standardization process. As this work showed, code-based designs can be competitive with other schemes when optimized hardware is developed.

## References

- [AAB<sup>+</sup>20] Carlos Aguilar Melchor, Nicolas Aragon, Slim Bettaieb, Loïc Bidoux, Olivier Blazy, Jean-Christophe Deneuville, Philippe Gaborit, Edoardo Persichetti, Gilles Zémor, and Jurjen Bos. HQC. Technical report, National Institute of Standards and Technology, 2020. available at [https://pqc-hqc.org/doc/hqc-specification\\_2021-06-06.pdf](https://pqc-hqc.org/doc/hqc-specification_2021-06-06.pdf).
- [CCD<sup>+</sup>22] Po-Jen Chen, Tung Chou, Sanjay Deshpande, Norman Lahr, Ruben Niederhagen, Jakub Szefer, and Wen Wang. Complete and improved FPGA implementation of Classic McEliece. *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2022(3):71–113, 2022.
- [DdPM<sup>+</sup>21] Sanjay Deshpande, Santos Merino del Pozo, Victor Mateu, Marc Manzano, Najwa Aaraj, and Jakub Szefer. Modular inverse for integers using fast constant time gcd algorithm and its applications. In *Proceedings of the International Conference on Field Programmable Logic and Applications*, FPL, August 2021.
- [DMG21] Viet Ba Dang, Kamyar Mohajerani, and Kris Gaj. High-speed hardware architectures and fpga benchmarking of crystals-kyber, ntru, and saber. *Cryptology ePrint Archive*, Paper 2021/1508, 2021. <https://eprint.iacr.org/2021/1508>.
- [DNN<sup>+</sup>22] Sanjay Deshpande, Mamuri Nawan, Kashif Nawaz, Jakub Szefer, and Chuanqi Xu. Don't Wait for SHAKE256: A Fast HQC Hardware Implementation. 2022.
- [GGM<sup>+</sup>22] Andrea Galimberti, Davide Galli, Gabriele Montanaro, William Fornaciari, and Davide Zoni. On the use of hardware accelerators in qc-mdpc code-based cryptography. In *Proceedings of the 19th ACM International Conference on Computing Frontiers*, CF '22, page 193–194, New York, NY, USA, 2022. Association for Computing Machinery.
- [GHJ<sup>+</sup>22] Qian Guo, Clemens Hlauschek, Thomas Johansson, Norman Lahr, Alexander Nilsson, and Robin Leander Schröder. Don't reject this: Key-recovery timing attacks due to rejection-sampling in hqc and bike. *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2022, Issue 3:223–263, 2022.
- [Gig04] Paul Gigliotti. Implementing barrel shifters using multipliers. Technical Report XAPP195, Xilinx, 2004.

- [HHK17] Dennis Hofheinz, Kathrin Hövelmanns, and Eike Kiltz. A modular analysis of the Fujisaki–Okamoto transformation. In Yael Kalai and Leonid Reyzin, editors, *Theory of Cryptography*, pages 341–371, Cham, 2017. Springer International Publishing.
- [HWCW19] Jingwei Hu, Wen Wang, Ray C.C. Cheung, and Huaxiong Wang. Optimized polynomial multiplier over commutative rings on fpgas: A case study on bike. In *2019 International Conference on Field-Programmable Technology (ICFPT)*, pages 231–234, 2019.
- [JGCS21] Arpan Jati, Naina Gupta, Anupam Chattopadhyay, and Somitra Kumar Sanadhya. A configurable CRYSTALS-Kyber hardware implementation with side-channel protection. *Cryptology ePrint Archive*, 2021.
- [MGCZ22] Gabriele Montanaro, Andrea Galimberti, Ernesto Colizzi, and Davide Zoni. Hardware-software co-design of bike with hls-generated accelerators, 2022.
- [MLRB20] Pedro Maat C. Massolino, Patrick Longa, Joost Renes, and Lejla Batina. A compact and scalable hardware/software co-design of SIKE. *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2020(2):245–271, Mar. 2020.
- [RBCGG21] Jan Richter-Brockmann, Ming-Shing Chen, Santosh Ghosh, and Tim Güneysu. Racing bike: Improved polynomial multiplication and inversion in hardware. *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2022(1):557–588, Nov. 2021.
- [RBMG22] Jan Richter-Brockmann, Johannes Mono, and Tim Guneyusu. Folding bike: Scalable hardware implementation for reconfigurable devices. *IEEE Transactions on Computers*, 71(5):1204–1215, 2022.
- [SR17] Cecilia Sandoval-Ruiz. Vhdl optimized model of a multiplier in finite fields. *Ingenieria y Universidad*, 21(2):195–212, Jun. 2017.
- [WTJ<sup>+</sup>20] Wen Wang, Shanquan Tian, Bernhard Jungk, Nina Bindel, Patrick Longa, and Jakub Szefer. Parameterized hardware accelerators for lattice-based cryptography and their application to the hw/sw co-design of qTESLA. *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2020(3):269–306, Jun. 2020.
- [XL21] Yufei Xing and Shuguo Li. A compact hardware implementation of cca-secure key exchange mechanism crystals-kyber on fpga. *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2021(2):328–356, Feb. 2021.
- [ZZC<sup>+</sup>22] Yihong Zhu, Wenping Zhu, Chen Chen, Min Zhu, Zhengdong Li, Shaojun Wei, and Leibo Liu. Compact  $gf(2)$  systemizer and optimized constant-time hardware sorters for key generation in classic mceliece. *Cryptology ePrint Archive*, Paper 2022/1277, 2022. <https://eprint.iacr.org/2022/1277>.