

Hamming Quasi-Cyclic (HQC)

HQC is an IND-CCA2 KEM running for standardization to NIST's competition in the category "post-quantum public key encryption scheme". Different sets of parameters are proposed for security strength categories 1, 3, and 5.

Principal Submitters (by alphabetical order):

- Carlos AGUILAR MELCHOR
- Nicolas ARAGON
- Slim BETTAIEB
- Loïc BIDOUX
- Olivier BLAZY
- Jean-Christophe DENEUVILLE
- Philippe GABORIT
- Edoardo PERSICHETTI
- Gilles ZÉMOR

Inventors: Same as submitters

Developers: Same as submitters

Owners: Same as submitters

Main contact

✉ Philippe GABORIT
@ philippe.gaborit@unilim.fr
☎ +33-626-907-245
≡ University of Limoges
✉ 123 avenue Albert Thomas
87 060 Limoges Cedex
France

Backup point of contact

✉ Jean-Christophe DENEUVILLE
@ jch.deneuville@gmail.com
☎ +33-631-142-705
≡ INSA-CVL Bourges &
University of Limoges
✉ 4 rue Jean le Bail
87 000 Limoges
France

Signatures

Digital copies of the signed statements are provided in Appendix A. The original paper versions will be given to Dustin MOODY directly at the First PQC Standardization Conference.

Contents

1	Specifications	3
1.1	Preliminaries	3
1.1.1	General definitions	3
1.1.2	Difficult problems for cryptography	5
1.2	Encryption and security	7
1.3	Presentation of the scheme	9
1.3.1	Public key encryption version (HQC.PKE)	9
1.3.2	KEM/DEM version (HQC.KEM)	10
1.3.3	A hybrid encryption scheme (HQC.HE)	11
1.4	Analysis of the error vector distribution for Hamming distance	11
1.5	Decoding codes with low rates and good decoding properties	12
1.5.1	Tensor product codes	13
1.5.2	BCH codes	14
1.5.3	Decoding BCH codes	15
1.5.4	Decryption Failure Probability	17
1.6	Parameters	18
2	Performance Analysis	20
2.1	Reference Implementation	22
2.2	Optimized Implementation	22
3	Known Answer Test Values	23
4	Security	23
5	Known Attacks	27
6	Advantages and Limitations	28
6.1	Advantages	28
6.2	Limitations	28
	References	29
A	Signed statements by the submitters	32

1 Specifications

In this section, we introduce HQC, an efficient encryption scheme based on coding theory. HQC stands for Hamming Quasi-Cyclic. This proposal is currently under revision for publication in IEEE Transactions on Information Theory. Many notations, definitions and properties are very similar to [12]. We nevertheless include them in this proposal for completeness.

HQC is a code-based public key cryptosystem with several desirable properties:

- It is proved IND-CPA assuming the hardness of (a decisional version of) the Syndrome Decoding on structured codes. By construction, HQC perfectly fits the recent KEM-DEM transformation of [21], and allows to get an hybrid encryption scheme with strong security guarantees (IND-CCA2) and good efficiency,
- In contrast with most code-based cryptosystems, the assumption that the family of codes being used is indistinguishable among random codes is no longer required, and
- It features a decryption failure probability analysis.

Organization of the Specifications. This section is organized as follows: we provide the required background in Sec. 1.1, we make some recalls on encryption and security in Sec. 1.2 then present our proposal in Sec. 1.3. An analysis of the decryption failure rate is proposed in Sec. 1.4. Details about codes being used are provided in Sec. 1.5, together with a specific analysis for these codes. Finally, concrete sets of parameters are provided in Sec. 1.6.

1.1 Preliminaries

1.1.1 General definitions

Throughout this document, \mathbb{Z} denotes the ring of integers and \mathbb{F}_2 the binary finite field. Additionally, we denote by $\omega(\cdot)$ the Hamming weight of a vector *i.e.* the number of its non-zero coordinates, and by $\mathcal{S}_w^n(\mathbb{F}_2)$ the set of words in \mathbb{F}_2^n of weight w . Formally:

$$\mathcal{S}_w^n(\mathbb{F}_2) = \{\mathbf{v} \in \mathbb{F}_2^n, \text{ such that } \omega(\mathbf{v}) = w\}.$$

\mathcal{V} denotes a vector space of dimension n over \mathbb{F} for some positive $n \in \mathbb{Z}$. Elements of \mathcal{V} can be interchangeably considered as row vectors or polynomials in $\mathcal{R} = \mathbb{F}[X]/(X^n - 1)$. Vectors/Polynomials (resp. matrices) will be represented by lower-case (resp. upper-case) bold letters. A prime integer n is said primitive if the polynomial $X^n - 1/(X - 1)$ is irreducible in \mathcal{R} .

For $\mathbf{u}, \mathbf{v} \in \mathcal{V}$, we define their product similarly as in \mathcal{R} , *i.e.* $\mathbf{uv} = \mathbf{w} \in \mathcal{V}$ with

$$w_k = \sum_{i+j \equiv k \pmod n} x_i y_j, \text{ for } k \in \{0, 1, \dots, n-1\}. \quad (1)$$

Our new protocol takes great advantage of the cyclic structure of matrices. In the same fashion as [1], $\mathbf{rot}(\mathbf{h})$ for $\mathbf{h} \in \mathcal{V}$ denotes the circulant matrix whose i^{th} column is the vector corresponding to $\mathbf{h}X^i$. This is captured by the following definition.

Definition 1.1.1 (Circulant Matrix). *Let $\mathbf{v} = (v_1, \dots, v_n) \in \mathbb{F}_2^n$. The circulant matrix induced by \mathbf{v} is defined and denoted as follows:*

$$\mathbf{rot}(\mathbf{v}) = \begin{pmatrix} v_0 & v_{n-1} & \dots & v_1 \\ v_1 & v_0 & \dots & v_2 \\ \vdots & \vdots & \ddots & \vdots \\ v_{n-1} & v_{n-2} & \dots & v_0 \end{pmatrix} \in \mathbb{F}^{n \times n} \quad (2)$$

As a consequence, it is easy to see that the product of any two elements $\mathbf{x}, \mathbf{y} \in \mathcal{R}$ can be expressed as a usual vector-matrix (or matrix-vector) product using the $\mathbf{rot}(\cdot)$ operator as

$$\mathbf{u} \cdot \mathbf{v} = \mathbf{u} \times \mathbf{rot}(\mathbf{v})^\top = (\mathbf{rot}(\mathbf{u}) \times \mathbf{v}^\top)^\top = \mathbf{v} \times \mathbf{rot}(\mathbf{u})^\top = \mathbf{v} \cdot \mathbf{u}. \quad (3)$$

Coding Theory. We now recall some basic definitions and properties about coding theory that will be useful to our construction. We mainly focus on general definitions, and refer the reader to Sec. 1.3 the description of the scheme, and also to [22] for a complete survey on code-based cryptography.

Definition 1.1.2 (Linear Code). *A Linear Code \mathcal{C} of length n and dimension k (denoted $[n, k]$) is a subspace of \mathcal{R} of dimension k . Elements of \mathcal{C} are referred to as codewords.*

Definition 1.1.3 (Generator Matrix). *We say that $\mathbf{G} \in \mathbb{F}^{k \times n}$ is a Generator Matrix for the $[n, k]$ code \mathcal{C} if*

$$\mathcal{C} = \{\mathbf{m}\mathbf{G}, \text{ for } \mathbf{m} \in \mathbb{F}^k\}. \quad (4)$$

Definition 1.1.4 (Parity-Check Matrix). *Given an $[n, k]$ code \mathcal{C} , we say that $\mathbf{H} \in \mathbb{F}^{(n-k) \times n}$ is a Parity-Check Matrix for \mathcal{C} if \mathbf{H} is a generator matrix of the dual code \mathcal{C}^\perp , or more formally, if*

$$\mathcal{C}^\perp = \{\mathbf{v} \in \mathbb{F}^n \text{ such that } \mathbf{H}\mathbf{v}^\top = \mathbf{0}\}. \quad (5)$$

Definition 1.1.5 (Syndrome). *Let $\mathbf{H} \in \mathbb{F}_2^{(n-k) \times n}$ be a parity-check matrix of some $[n, k]$ code \mathcal{C} , and $\mathbf{v} \in \mathbb{F}_2^n$ be a word. Then the syndrome of \mathbf{v} is $\mathbf{H}\mathbf{v}^\top$, and we have $\mathbf{v} \in \mathcal{C} \Leftrightarrow \mathbf{H}\mathbf{v}^\top = \mathbf{0}$.*

Definition 1.1.6 (Minimum Distance). *Let \mathcal{C} be an $[n, k]$ linear code over \mathcal{R} and let ω be a norm on \mathcal{R} . The Minimum Distance of \mathcal{C} is*

$$d = \min_{\mathbf{u}, \mathbf{v} \in \mathcal{C}, \mathbf{u} \neq \mathbf{v}} \omega(\mathbf{u} - \mathbf{v}). \quad (6)$$

A code with minimum distance d is capable of decoding arbitrary patterns of up to $\delta = \lfloor \frac{d-1}{2} \rfloor$ errors. Code parameters are denoted $[n, k, d]$.

Code-based cryptography usually suffers from huge keys. In order to keep our cryptosystem efficient, we will use the strategy of Gaborit [17] for shortening keys. This results in Quasi-Cyclic Codes, as defined below.

Definition 1.1.7 (Quasi-Cyclic Codes [28]). View a vector $\mathbf{c} = (\mathbf{c}_1, \dots, \mathbf{c}_s)$ of \mathbb{F}_2^{sn} as s successive blocks (n -tuples). An $[sn, k, d]$ linear code \mathcal{C} is Quasi-Cyclic (QC) of index s if, for any $\mathbf{c} = (\mathbf{c}_1, \dots, \mathbf{c}_s) \in \mathcal{C}$, the vector obtained after applying a simultaneous circular shift to every block $\mathbf{c}_1, \dots, \mathbf{c}_s$ is also a codeword.

More formally, by considering each block \mathbf{c}_i as a polynomial in $\mathcal{R} = \mathbb{F}[X]/(X^n - 1)$, the code \mathcal{C} is QC of index s if for any $\mathbf{c} = (\mathbf{c}_1, \dots, \mathbf{c}_s) \in \mathcal{C}$ it holds that $(X \cdot \mathbf{c}_1, \dots, X \cdot \mathbf{c}_s) \in \mathcal{C}$.

Definition 1.1.8 (Systematic Quasi-Cyclic Codes). A systematic Quasi-Cyclic $[sn, n]$ code of index s and rate $1/s$ is a quasi-cyclic code with an $(s-1)n \times sn$ parity-check matrix of the form:

$$\mathbf{H} = \begin{bmatrix} \mathbf{I}_n & 0 & \cdots & 0 & \mathbf{A}_1 \\ 0 & \mathbf{I}_n & & & \mathbf{A}_2 \\ & & \ddots & & \vdots \\ 0 & & \cdots & \mathbf{I}_n & \mathbf{A}_{s-1} \end{bmatrix} \quad (7)$$

where $\mathbf{A}_1, \dots, \mathbf{A}_{s-1}$ are circulant $n \times n$ matrices.

Remark 1.1. The definition of systematic quasi-cyclic codes of index s can of course be generalized to all rates ℓ/s , $\ell = 1 \dots s-1$, but we shall only use systematic QC-codes of rates $1/2$ and $1/3$ and wish to lighten notation with the above definition. In the sequel, referring to a systematic QC-code will imply by default that it is of rate $1/s$. Note that arbitrary QC-codes are not necessarily equivalent to a systematic QC-code.

1.1.2 Difficult problems for cryptography

In this section we describe difficult problems which can be used for cryptography and discuss their complexity.

All problems are variants of the *decoding problem*, which consists of looking for the closest codeword to a given vector: when dealing with linear codes, it is readily seen that the decoding problem stays the same when one is given the *syndrome* of the received vector rather than the received vector. We therefore speak of *Syndrome Decoding* (SD).

Definition 1.1.9 (SD Distribution). For positive integers, n , k , and w , the $\text{SD}(n, k, w)$ Distribution chooses $\mathbf{H} \xleftarrow{\$} \mathbb{F}^{(n-k) \times n}$ and $\mathbf{x} \xleftarrow{\$} \mathbb{F}^n$ such that $\omega(\mathbf{x}) = w$, and outputs $(\mathbf{H}, \sigma(\mathbf{x}) = \mathbf{H}\mathbf{x}^\top)$.

Definition 1.1.10 (Search SD Problem). Let ω be a norm over \mathcal{R} . On input $(\mathbf{H}, \mathbf{y}^\top) \in \mathbb{F}^{(n-k) \times n} \times \mathbb{F}^{(n-k)}$ from the SD distribution, the Syndrome Decoding Problem $\text{SD}(n, k, w)$ asks to find $\mathbf{x} \in \mathbb{F}^n$ such that $\mathbf{H}\mathbf{x}^\top = \mathbf{y}^\top$ and $\omega(\mathbf{x}) = w$.

For the Hamming distance the SD problem has been proven to be NP-complete in [4]. This problem can also be seen as the Learning Parity with Noise (LPN) problem with a fixed number of samples [2]. For cryptography we also need a decision version of the problem, which is given in the following definition.

Definition 1.1.11 (Decision SD Problem). *On input $(\mathbf{H}, \mathbf{y}^\top) \xleftarrow{\$} \mathbb{F}^{(n-k) \times n} \times \mathbb{F}^{(n-k)}$, the Decision SD Problem $\text{DSD}(n, k, w)$ asks to decide with non-negligible advantage whether $(\mathbf{H}, \mathbf{y}^\top)$ came from the $\text{SD}(n, k, w)$ distribution or the uniform distribution over $\mathbb{F}^{(n-k) \times n} \times \mathbb{F}^{(n-k)}$.*

As mentioned above, this problem is the problem of decoding random linear codes from random errors. The random errors are often taken as independent Bernoulli variables acting independently on vector coordinates, rather than uniformly chosen from the set of errors of a given weight, but this hardly makes any difference and one model rather than the other is a question of convenience. The DSD problem has been shown to be polynomially equivalent to its search version in [2].

Finally, as our cryptosystem will use QC-codes, we explicitly define the problem on which our cryptosystem will rely. The following definitions describe the DSD problem in the QC configuration, and are just a combination of Def. 1.1.7 and 1.1.11. Quasi-Cyclic codes are very useful in cryptography since their compact description allows to decrease considerably the size of the keys. In particular the case $s = 2$ corresponds to double circulant codes with generator matrices of the form $(\mathbf{I}_n \mid \mathbf{A})$ for \mathbf{A} a circulant matrix. Such double circulant codes have been used for almost 10 years in cryptography (cf [18]) and more recently in [28]. Quasi-cyclic codes of index 3 are also considered in [28].

Definition 1.1.12 (s -QCSD Distribution). *For positive integers n , w and s , the s -QCSD(n, w) Distribution chooses uniformly at random a parity matrix $\mathbf{H} \xleftarrow{\$} \mathbb{F}^{(sn-n) \times sn}$ of a systematic QC code \mathcal{C} of index s and rate $1/s$ (see Def. 1.1.8) together with a vector $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_s) \xleftarrow{\$} \mathbb{F}^{sn}$ such that $\omega(\mathbf{x}_i) = w$, $i = 1..s$, and outputs $(\mathbf{H}, \mathbf{H}\mathbf{x}^\top)$.*

Definition 1.1.13 ((Search) s -QCSD Problem). *For positive integers n , w , s , a random parity check matrix \mathbf{H} of a systematic QC code \mathcal{C} of index s and $\mathbf{y} \xleftarrow{\$} \mathbb{F}^{sn-n}$, the Search s -Quasi-Cyclic SD Problem s -QCSD(n, w) asks to find $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_s) \in \mathbb{F}^{sn}$ such that $\omega(\mathbf{x}_i) = w$, $i = 1..s$, and $\mathbf{y} = \mathbf{x}\mathbf{H}^\top$.*

It would be somewhat more natural to choose the parity-check matrix \mathbf{H} to be made up of independent uniformly random circulant submatrices, rather than with the special form required by (7). We choose this distribution so as to make the security reduction to follow less technical. It is readily seen that, for fixed s , when choosing quasi-cyclic codes with this more general distribution, one obtains with non-negligible probability, a quasi-cyclic code that admits a parity-check matrix of the form (7). Therefore requiring quasi-cyclic codes to be systematic does not hurt the generality of the decoding problem for quasi-cyclic codes. A similar remark holds for the slightly special form of weight distribution of the vector \mathbf{x} .

Assumption 1. *Although there is no general complexity result for quasi-cyclic codes, decoding these codes is considered hard by the community. There exist general attacks which uses the cyclic structure of the code [31] but these attacks have only a very limited impact on the practical complexity of the problem. The conclusion is that in practice, the best attacks are the same as those for non-circulant codes up to a small factor.*

The problem has a decisional form:

Definition 1.1.14 (Decision s -QCSD Problem). *For positive integers n, w, s , a random parity check matrix \mathbf{H} of a systematic QC code \mathcal{C} and $\mathbf{y} \xleftarrow{\$} \mathbb{F}^{sn}$, the Decision s -Quasi-Cyclic SD Problem s -DQCSD(n, w) asks to decide with non-negligible advantage whether $(\mathbf{H}, \mathbf{y}^\top)$ came from the s -QCSD(n, w) distribution or the uniform distribution over $\mathbb{F}^{(sn-n) \times sn} \times \mathbb{F}^{(sn-n)}$.*

As for the ring-LPN problem, there is no known reduction from the search version of s -QCSD problem to its decision version. The proof of [2] cannot be directly adapted in the quasi-cyclic case, however the best known attacks on the decision version of the problem s -QCSD remain the direct attacks on the search version of the problem s -QCSD.

1.2 Encryption and security

Encryption Scheme. An encryption scheme is a tuple of four polynomial time algorithms (Setup, KeyGen, Encrypt, Decrypt):

- **Setup**(1^λ), where λ is the security parameter, generates the global parameters **param** of the scheme;
- **KeyGen**(**param**) outputs a pair of keys, a (public) encryption key **pk** and a (private) decryption key **sk**;
- **Encrypt**(**pk**, **m**, θ) outputs a ciphertext **c**, on the message **m**, under the encryption key **pk**, with the randomness θ . We also use **Encrypt**(**pk**, **m**) for the sake of clarity;
- **Decrypt**(**sk**, **c**) outputs the plaintext **m**, encrypted in the ciphertext **c** or \perp .

Such an encryption scheme has to satisfy both *Correctness* and *Indistinguishability under Chosen Plaintext Attack* (IND-CPA) security properties.

Correctness: For every λ , every **param** \leftarrow **Setup**(1^λ), every pair of keys (**pk**, **sk**) generated by **KeyGen**, every message **m**, we should have $P[\text{Decrypt}(\text{sk}, \text{Encrypt}(\text{pk}, \mathbf{m}, \theta)) = \mathbf{m}] = 1 - \text{negl}(\lambda)$ for $\text{negl}(\cdot)$ a negligible function, where the probability is taken over varying randomness θ .

IND-CPA [19]: This notion formalized by the game depicted in Fig. 1, states that an adversary should not be able to efficiently guess which plaintext has been encrypted even if he knows it is one among two plaintexts of his choice.

In the following, we denote by $|\mathcal{A}|$ the running time of an adversary \mathcal{A} . The global advantage for polynomial time adversaries running in time less than t is:

$$\text{Adv}_{\mathcal{E}}^{\text{ind}}(\lambda, t) = \max_{|\mathcal{A}| \leq t} \text{Adv}_{\mathcal{E}, \mathcal{A}}^{\text{ind}}(\lambda), \quad (8)$$

where $\text{Adv}_{\mathcal{E}, \mathcal{A}}^{\text{ind}}(\lambda)$ is the advantage the adversary \mathcal{A} has in winning game $\text{Exp}_{\mathcal{E}, \mathcal{A}}^{\text{ind}-b}(\lambda)$:

Exp _{\mathcal{E}, \mathcal{A}} ^{ind- b} (λ)

1. $\text{param} \leftarrow \text{Setup}(1^\lambda)$
2. $(\text{pk}, \text{sk}) \leftarrow \text{KeyGen}(\text{param})$
3. $(\mathbf{m}_0, \mathbf{m}_1) \leftarrow \mathcal{A}(\text{FIND} : \text{pk})$
4. $\mathbf{c}^* \leftarrow \text{Encrypt}(\text{pk}, \mathbf{m}_b, \theta)$
5. $b' \leftarrow \mathcal{A}(\text{GUESS} : \mathbf{c}^*)$
6. RETURN b'

Figure 1: Game for the IND-CPA security of an asymmetric encryption scheme.

$$\text{Adv}_{\mathcal{E}, \mathcal{A}}^{\text{ind}}(\lambda) = \left| \Pr[\mathbf{Exp}_{\mathcal{E}, \mathcal{A}}^{\text{ind}-1}(\lambda) = 1] - \Pr[\mathbf{Exp}_{\mathcal{E}, \mathcal{A}}^{\text{ind}-0}(\lambda) = 1] \right|. \quad (9)$$

IND-CPA and IND-CCA2: Note that the standard security requirement for a public key cryptosystem is IND-CCA2, *indistinguishability against adaptive chosen-ciphertext attacks*, and not just IND-CPA. The main difference is that for IND-CCA2 indistinguishability must hold even if the attacker is given a *decryption oracle* first when running the FIND algorithm and also when running the GUESS algorithm (but cannot query the oracle on the challenge ciphertext \mathbf{c}^*). We do not present the associated formal game and definition as an existing (and inexpensive) transformation can be used [21] for our scheme to pass from IND-CPA to IND-CCA2. Various generic techniques transforming a IND-CPA scheme into an IND-CCA2 scheme are known [15, 16, 29, 11] but cannot be applied to our scheme due to potential decryption errors.

In [21] Hofheinz et al. present a generic transformation that takes into account decryption errors and can be applied directly to our scheme. Roughly, their construction provides a way to convert a guarantee against passive adversaries into indistinguishability against active ones by turning a public key cryptosystem into a KEM-DEM. The tightness (the quality factor) of the reduction depends on the ciphertext distribution. Regarding our scheme, random words only have a negligible (in the security parameter) probability of being valid ciphertexts. In other words, the γ -spreadness factor of [21] is small enough so that there is no loss between the IND-CPA security of our public key cryptosystem and the IND-CCA2 security of the KEM-DEM version presented in Fig. 3.

The security reduction is tight in the random oracle model and does not require any supplemental property from our scheme as we have the IND-CPA property (instead of just a weaker property called *One-Wayness*). Let us denote by $\text{Encrypt}(\text{pk}, \mathbf{m}, \theta)$ the encryption function defined in Fig. 2 that uses randomness θ to generate uniformly random values \mathbf{r}_1 , \mathbf{r}_2 , and \mathbf{e} . The idea of [21] transformation is to de-randomize the encryption function $\text{Encrypt}(\text{pk}, \mathbf{m}, \theta)$ by using a hash function \mathcal{G} and do a deterministic encryption of \mathbf{m} by calling $c = \text{Encrypt}(\text{pk}, \mathbf{m}, \mathcal{G}(\mathbf{m}))$. The ciphertext is sent together with a hash $K = \mathcal{H}(\mathbf{c}, \mathbf{m})$ that ties the ciphertext to the plaintext. The receiver then decrypts \mathbf{c} into \mathbf{m} , checks the hash value, and uses again the deterministic encryption to check that \mathbf{c} is indeed *the* ciphertext associated to \mathbf{m} .

As the reduction is tight we do not need to change our parameters when we pass from IND-CPA to IND-CCA2. From a computational point of view, the overhead for the sender is two hash calls and for the receiver it is two hash calls and an encrypt call. From a communication point of view the overhead is the bitsize of a hash (or two if the reduction must hold in the Quantum Random Oracle Model, see [21] for more details).

1.3 Presentation of the scheme

In this section, we describe our proposal: HQC. We begin with the PKE version, then describe the transformation of [21] to obtain a KEM-DEM that achieves IND-CCA2. Parameter sets can be found in Sec. 1.6.

1.3.1 Public key encryption version (HQC.PKE)

Presentation of the scheme. HQC uses two types of codes: a decodable $[n, k]$ code \mathcal{C} , generated by $\mathbf{G} \in \mathbb{F}^{k \times n}$ and which can correct at least δ errors via an efficient algorithm $\mathcal{C}.\text{Decode}(\cdot)$; and a random double-circulant $[2n, n]$ code, of parity-check matrix $(\mathbf{1}, \mathbf{h})$. The four polynomial-time algorithms constituting our scheme are depicted in Fig. 2.

- **Setup**(1^λ): generates and outputs the global parameters $\text{param} = (n, k, \delta, w, w_{\mathbf{r}}, w_{\mathbf{e}})$.
- **KeyGen**(param): samples $\mathbf{h} \xleftarrow{\$} \mathcal{R}$, the generator matrix $\mathbf{G} \in \mathbb{F}^{k \times n}$ of \mathcal{C} , $\mathbf{s} \mathbf{k} = (\mathbf{x}, \mathbf{y}) \xleftarrow{\$} \mathcal{R}^2$ such that $\omega(\mathbf{x}) = \omega(\mathbf{y}) = w$, sets $\mathbf{p} \mathbf{k} = (\mathbf{h}, \mathbf{s} = \mathbf{x} + \mathbf{h} \cdot \mathbf{y})$, and returns $(\mathbf{p} \mathbf{k}, \mathbf{s} \mathbf{k})$.
- **Encrypt**($\mathbf{p} \mathbf{k}, \mathbf{m}$): generates $\mathbf{e} \xleftarrow{\$} \mathcal{R}$, $\mathbf{r} = (\mathbf{r}_1, \mathbf{r}_2) \xleftarrow{\$} \mathcal{R}^2$ such that $\omega(\mathbf{e}) = w_{\mathbf{e}}$ and $\omega(\mathbf{r}_1) = \omega(\mathbf{r}_2) = w_{\mathbf{r}}$, sets $\mathbf{u} = \mathbf{r}_1 + \mathbf{h} \cdot \mathbf{r}_2$ and $\mathbf{v} = \mathbf{m} \mathbf{G} + \mathbf{s} \cdot \mathbf{r}_2 + \mathbf{e}$, returns $\mathbf{c} = (\mathbf{u}, \mathbf{v})$.
- **Decrypt**($\mathbf{s} \mathbf{k}, \mathbf{c}$): returns $\mathcal{C}.\text{Decode}(\mathbf{v} - \mathbf{u} \cdot \mathbf{y})$.

Figure 2: Description of our proposal HQC.PKE.

Notice that the generator matrix \mathbf{G} of the code \mathcal{C} is publicly known, so the security of the scheme and the ability to decrypt do not rely on the knowledge of the error correcting code \mathcal{C} being used.

Correctness. The correctness of our new encryption scheme clearly relies on the decoding capability of the code \mathcal{C} . Specifically, assuming $\mathcal{C}.\text{Decode}$ correctly decodes $\mathbf{v} - \mathbf{u} \cdot \mathbf{y}$, we have:

$$\text{Decrypt}(\mathbf{s} \mathbf{k}, \text{Encrypt}(\mathbf{p} \mathbf{k}, \mathbf{m})) = \mathbf{m}. \quad (10)$$

And $\mathcal{C}.\text{Decode}$ correctly decodes $\mathbf{v} - \mathbf{u} \cdot \mathbf{y}$ whenever

$$\omega(\mathbf{s} \cdot \mathbf{r}_2 - \mathbf{u} \cdot \mathbf{y} + \mathbf{e}) \leq \delta \quad (11)$$

$$\omega((\mathbf{x} + \mathbf{h} \cdot \mathbf{y}) \cdot \mathbf{r}_2 - (\mathbf{r}_1 + \mathbf{h} \cdot \mathbf{r}_2) \cdot \mathbf{y} + \mathbf{e}) \leq \delta \quad (12)$$

$$\omega(\mathbf{x} \cdot \mathbf{r}_2 - \mathbf{r}_1 \cdot \mathbf{y} + \mathbf{e}) \leq \delta \quad (13)$$

In order to provide an upper bound on the decryption failure probability, an analysis of the distribution of the error vector $\mathbf{e}' = \mathbf{x} \cdot \mathbf{r}_2 - \mathbf{r}_1 \cdot \mathbf{y} + \mathbf{e}$ is provided in Sec. 1.4.

1.3.2 KEM/DEM version (HQC.KEM)

Let \mathcal{E} be an instance of the HQC cryptosystem as described above. Let \mathcal{G} , \mathcal{H} , and \mathcal{K} be hash functions, typically SHA512 as advised by NIST¹. The KEM-DEM version of the HQC cryptosystem is defined as follows:

- **Setup**(1^λ): as before, except that k will be the length of the symmetric key being exchanged, typically $k = 256$.
- **KeyGen**(param): exactly as before.
- **Encapsulate**(pk): generate $\mathbf{m} \xleftarrow{\$} \mathbb{F}^k$ (this will serve as a seed to derive the shared key). Derive the randomness $\theta \leftarrow \mathcal{G}(\mathbf{m})$. Generate the ciphertext $c \leftarrow (\mathbf{u}, \mathbf{v}) = \mathcal{E}.\text{Encrypt}(\text{pk}, \mathbf{m}, \theta)$, and derive the symmetric key $K \leftarrow \mathcal{K}(\mathbf{m}, c)$. Let $\mathbf{d} \leftarrow \mathcal{H}(\mathbf{m})$, and send (c, \mathbf{d}) .
- **Decapsulate**(sk, c, d): Decrypt $\mathbf{m}' \leftarrow \mathcal{E}.\text{Decrypt}(\text{sk}, c)$, compute $\theta' \leftarrow \mathcal{G}(\mathbf{m}')$, and (re-)encrypt \mathbf{m}' to get $c' \leftarrow \mathcal{E}.\text{Encrypt}(\text{pk}, \mathbf{m}', \theta')$. If $c \neq c'$ or $\mathbf{d} \neq \mathcal{H}(\mathbf{m}')$ then abort. Otherwise, derive the shared key $K \leftarrow \mathcal{K}(\mathbf{m}, c)$.

Figure 3: Description of our proposal HQC.KEM.

According to [21], the KEM-DEM version of HQC is IND-CCA2. More details regarding the tightness of the reduction are provided at the end of Sec. 1.6.

Security concerns and implementation details. Notice that while NIST only recommends SHA512 as a hash function (or TupleHash256 for hardware efficiency purposes), the transformation of [21] would be dangerous – at least in our setting – if one sets $\mathcal{G} = \mathcal{H}$. Indeed, publishing the randomness $\theta = \mathcal{G}(\mathbf{m}) = \mathcal{H}(\mathbf{m}) = \mathbf{d}$ used to generate \mathbf{r}_1 , \mathbf{r}_2 , and \mathbf{e} , would allow one to retrieve \mathbf{s} , the secret key of \mathcal{E} .

We therefore suggest to use a pseudo-random function for \mathcal{G} , such as an AES-based seed expander, and SHA512 for \mathcal{H} .

¹See Dustin Moody’s mail entitled “new FAQ question” on PQC-forum (20/07/2017 – 12:58 CET)

1.3.3 A hybrid encryption scheme (HQC.HE)

While NIST claimed that they will be using generic transformations to convert any IND-CCA2 KEM into an IND-CCA2 PKE, no detail on these conversions have been provided. We therefore refer to HQC.HE to designate the PKE scheme resulting from applying a generic conversion to HQC.KEM.

1.4 Analysis of the error vector distribution for Hamming distance

The aim of this section is to determine the probability that the condition in Eq. (13) holds. In order to do so, we study the error distribution of the error vector $\mathbf{e}' = \mathbf{x} \cdot \mathbf{r}_2 - \mathbf{r}_1 \cdot \mathbf{y} + \mathbf{e}$.

The vectors $\mathbf{x}, \mathbf{y}, \mathbf{r}_1, \mathbf{r}_2, \mathbf{e}$ have been taken to be uniformly and independently chosen among vectors of weight w , $w_{\mathbf{r}}$ or $w_{\mathbf{e}}$. This distribution can be closely approximated by a binomial distribution \mathcal{B} , where a vector consists of n Bernoulli variables of parameter $p = w/n$ (or $p_{\mathbf{r}} = w_{\mathbf{r}}/n$ and $p_{\mathbf{e}} = w_{\mathbf{e}}/n$ respectively). In other words, $\mathcal{S}_w^n(\mathbb{F}_2)$ is close to $\mathcal{B}(n, w/n)$, similarly for $w_{\mathbf{r}}$ and $w_{\mathbf{e}}$. To simplify the analysis we shall assume this model rather than the constant weight uniform model. Both models are very close, and our cryptographic protocols work just as well in both settings.

We first evaluate the distributions of the products $\mathbf{x} \cdot \mathbf{r}_2$ and $\mathbf{r}_1 \cdot \mathbf{y}$.

Proposition 1.4.1. *Let $\mathbf{x} = (X_1, \dots, X_n)$ (resp. $\mathbf{r} = (R_1, \dots, R_n)$) be a random vector where the X_i (resp. R_i) are independent Bernoulli variables of parameter p (resp. $p_{\mathbf{r}}$), $P(X_i = 1) = p$ and $P(R_i = 1) = p_{\mathbf{r}}$. Assuming \mathbf{x} and \mathbf{r} are independent, and denoting $\mathbf{z} = \mathbf{x} \cdot \mathbf{r} = (Z_1, \dots, Z_n)$ as defined in Eq. (1), we have:*

$$\begin{cases} \Pr[Z_k = 1] = \frac{1}{2} - \frac{1}{2}(1 - 2pp_{\mathbf{r}})^n, \\ \Pr[Z_k = 0] = \frac{1}{2} + \frac{1}{2}(1 - 2pp_{\mathbf{r}})^n. \end{cases} \quad (14)$$

Proof. We have

$$Z_k = \sum_{i+j=k+1 \bmod n} X_i R_j \bmod 2. \quad (15)$$

Every term $X_i R_j$ is the product of two independent Bernoulli variables of parameter respectively p and $p_{\mathbf{r}}$, and is therefore a Bernoulli variable of parameter $p \times p_{\mathbf{r}}$. The variable Z_k is the sum modulo 2 of n such products, which are all independent since every variable X_i is involved exactly once in (15), for $1 \leq i \leq n$, and similarly every variable R_j is involved once in (15). Therefore Z_k is the sum modulo 2 of n independent Bernoulli variables of parameter $p \times p_{\mathbf{r}}$, and we have

$$\Pr[Z_k = 1] = \sum_{0 \leq i \leq n, i \text{ odd}} \binom{n}{i} (pp_{\mathbf{r}})^i (1 - pp_{\mathbf{r}})^{n-i}$$

which, using the equations:

$$\sum_{\substack{0 \leq i \leq n, \\ i \text{ odd}}} \binom{n}{i} a^i b^{n-i} = \frac{(a+b)^n - (a-b)^n}{2}, \text{ and } \sum_{\substack{0 \leq i \leq n, \\ i \text{ even}}} \binom{n}{i} a^i b^{n-i} = \frac{(a+b)^n + (a-b)^n}{2} \quad (16)$$

with $a = pp_{\mathbf{r}}$ and $b = 1 - pp_{\mathbf{r}}$, simplifies into the claimed result. \square

Let us denote by $\tilde{p} = \tilde{p}(n, w) = \Pr[Z_k = 1]$ from Eq. (14). Let \mathbf{x}, \mathbf{y} (resp. $\mathbf{r}_1, \mathbf{r}_2$) be independent random vectors whose coordinates are independently Bernoulli distributed with parameter p (resp. $p_{\mathbf{r}}$). Then the k -th coordinates of $\mathbf{x} \cdot \mathbf{r}_2$ and of $\mathbf{r}_1 \cdot \mathbf{y}$ are independent and Bernoulli distributed with parameter \tilde{p} . Therefore their modulo 2 sum $\mathbf{t} = \mathbf{x} \cdot \mathbf{r}_2 - \mathbf{r}_1 \cdot \mathbf{y}$ is Bernoulli distributed with

$$\begin{cases} \Pr[t_k = 1] = 2\tilde{p}(1 - \tilde{p}), \\ \Pr[t_k = 0] = (1 - \tilde{p})^2 + \tilde{p}^2. \end{cases} \quad (17)$$

Finally, by adding the term \mathbf{e} to \mathbf{t} , we obtain the distribution of the coordinates of the error vector $\mathbf{e}' = \mathbf{x} \cdot \mathbf{r}_2 - \mathbf{r}_1 \cdot \mathbf{y} + \mathbf{e}$. Since the coordinates of \mathbf{e} are Bernoulli of parameter $p_{\mathbf{e}}$ and those of \mathbf{t} are Bernoulli distributed as (17) and independent from \mathbf{e} , we obtain :

Proposition 1.4.2. *Let $\mathbf{x}, \mathbf{y} \sim \mathcal{B}(n, \frac{w}{n})$, $\mathbf{r}_1, \mathbf{r}_2 \sim \mathcal{B}(n, \frac{w_{\mathbf{r}}}{n})$ and $\mathbf{e} \sim \mathcal{B}(n, \frac{w_{\mathbf{e}}}{n})$, and let $\mathbf{e}' = \mathbf{x} \cdot \mathbf{r}_2 - \mathbf{r}_1 \cdot \mathbf{y} + \mathbf{e}$. Then*

$$\begin{cases} \Pr[e'_k = 1] = 2\tilde{p}(1 - \tilde{p})(1 - \frac{w_{\mathbf{e}}}{n}) + ((1 - \tilde{p})^2 + \tilde{p}^2) \frac{w_{\mathbf{e}}}{n}, \\ \Pr[e'_k = 0] = ((1 - \tilde{p})^2 + \tilde{p}^2) (1 - \frac{w_{\mathbf{e}}}{n}) + 2\tilde{p}(1 - \tilde{p}) \frac{w_{\mathbf{e}}}{n}. \end{cases} \quad (18)$$

Proposition 1.4.2 gives us the probability that a coordinate of the error vector \mathbf{e}' is 1. In our simulations to follow, which occur in the regime $p = \alpha\sqrt{n}$ with constant α , we make the simplifying assumption that the coordinates of \mathbf{e}' are independent, meaning that the weight of \mathbf{e}' follows a binomial distribution of parameter p^* , where p^* is defined as in Eq. (18): $p^* = 2\tilde{p}(1 - \tilde{p})(1 - \frac{w_{\mathbf{e}}}{n}) + ((1 - \tilde{p})^2 + \tilde{p}^2) \frac{w_{\mathbf{e}}}{n}$. This approximation will give us, for $0 \leq d \leq \min(2 \times w \times w_{\mathbf{r}} + w_{\mathbf{e}}, n)$,

$$\Pr[\omega(\mathbf{e}') = d] = \binom{n}{d} (p^*)^d (1 - p^*)^{(n-d)}. \quad (19)$$

In practice, the results obtained by simulation on the decryption failure are very coherent with this assumption.

1.5 Decoding codes with low rates and good decoding properties

The previous section allowed us to determine the distribution of the error vector \mathbf{e} in the configuration where a simple linear code is used. Now the decryption part corresponds to decoding the error described in the previous section. Any decodable code can be used at this point, depending on the considered application: clearly small dimension codes will allow better decoding, but at the cost of a lower encryption rate. The particular case that we consider corresponds typically to the case of key exchange or authentication, where only a small amount of data needs to be encrypted (typically 80, 128 or 256 bits, a symmetric secret key size). We therefore need codes with low rates which are able to correct many

errors. Again, a tradeoff is necessary between efficiently decodable codes but with a high decoding cost and less efficiently decodable codes but with a smaller decoding cost.

An example of such a family of codes with good decoding properties, meaning a simple decoding algorithm which can be analyzed, is given by Tensor Product Codes, which are used for biometry [7], where the same type of issue appears. More specifically, we will consider a special simple case of Tensor Product Codes (BCH codes and repetition codes), for which a precise analysis of the decryption failure can be obtained in the Hamming distance case.

1.5.1 Tensor product codes

Definition 1.5.1 (Tensor Product Code). *Let \mathcal{C}_1 (resp. \mathcal{C}_2) be a $[n_1, k_1, d_1]$ (resp. $[n_2, k_2, d_2]$) linear code over \mathbb{F} . The Tensor Product Code of \mathcal{C}_1 and \mathcal{C}_2 denoted $\mathcal{C}_1 \otimes \mathcal{C}_2$ is defined as the set of all $n_2 \times n_1$ matrices whose rows are codewords of \mathcal{C}_1 and whose columns are codewords of \mathcal{C}_2 .*

More formally, if \mathcal{C}_1 (resp. \mathcal{C}_2) is generated by \mathbf{G}_1 (resp. \mathbf{G}_2), then

$$\mathcal{C}_1 \otimes \mathcal{C}_2 = \{ \mathbf{G}_2^\top \mathbf{X} \mathbf{G}_1 \text{ for } \mathbf{X} \in \mathbb{F}^{k_2 \times k_1} \} \quad (20)$$

Remark 1.2. *Using the notation of the above definition, the tensor product of two linear codes is a $[n_1 n_2, k_1 k_2, d_1 d_2]$ linear code.*

Specifying the tensor product code. Even if tensor product codes seem well-suited for our purpose, an analysis similar to the one in Sec. 1.4 becomes much more complicated. Therefore, in order to provide strong guarantees on the decryption failure probability for our cryptosystem, we chose to restrict ourselves to a tensor product code $\mathcal{C} = \mathcal{C}_1 \otimes \mathcal{C}_2$, where \mathcal{C}_1 is a BCH(n_1, k_1, δ_1) code of length n_1 , dimension k_1 , and correcting capability δ_1 (i.e. it can correct up to δ_1 errors), and \mathcal{C}_2 is the repetition code of length n_2 and dimension 1, denoted $\mathbb{1}_{n_2}$. (Notice that $\mathbb{1}_{n_2}$ can decode up to $\delta_2 = \lfloor \frac{n_2-1}{2} \rfloor$.) Subsequently, the analysis becomes possible and remains accurate but the negative counterpart is that there probably are some other tensor product codes achieving better efficiency (or smaller key sizes).

In the Hamming metric version of the cryptosystem we propose, a message $\mathbf{m} \in \mathbb{F}^{k_1}$ is first encoded into $\mathbf{m}_1 \in \mathbb{F}^{n_1}$ with a BCH($n_1, k_1 = k, \delta_1$) code, then each coordinate $\mathbf{m}_{1,i}$ of \mathbf{m}_1 is re-encoded into $\tilde{\mathbf{m}}_{1,i} \in \mathbb{F}^{n_2}$ with a repetition code $\mathbb{1}_{n_2}$. We denote $n = n_1 n_2$ the length of the tensor product code² (its dimension is $k = k_1 \times 1$), and by $\tilde{\mathbf{m}}$ the resulting encoded vector, i.e. $\tilde{\mathbf{m}} = (\tilde{\mathbf{m}}_{1,1}, \dots, \tilde{\mathbf{m}}_{1,n_1}) \in \mathbb{F}^{n_1 n_2}$.

The efficient algorithm used for the repetition code is the majority decoding, i.e. more formally:

$$\mathbb{1}_{n_2}.\text{Decode}(\tilde{\mathbf{m}}_{1,j}) = \begin{cases} 1 & \text{if } \sum_{i=0}^{n_2-1} \tilde{\mathbf{m}}_{1,j,i} \geq \lceil \frac{n_2+1}{2} \rceil, \\ 0 & \text{otherwise.} \end{cases} \quad (21)$$

²In practice, the length is the smallest primitive prime greater than n to avoid algebraic attacks.

1.5.2 BCH codes

For any positive integers $m \geq 3$ and $t \leq 2^{m-1}$, there exists a binary BCH code with the following parameters [24]:

- Block length $n = 2^m - 1$
- Number of parity-check digits $n - k \leq m\delta$, with δ , the correcting capacity of the code and k the number of information bits
- Minimum distance $d_{min} \geq 2\delta + 1$

We denote this code by $\text{BCH}[n, k, \delta]$. Let α be the primitive element in \mathbb{F}_{2^m} , the generator polynomial $g(x)$ of the $\text{BCH}[n, k, \delta]$ code is given by:

$$g(x) = \text{LCM} \{ \phi_1(x), \phi_2(x), \dots, \phi_{2\delta}(x) \}$$

with $\phi_i(x)$ being the minimal polynomial of α^i (refer to [24] for more details on generator polynomial).

Depending on the parameters of the HQC scheme, we construct shortened BCH codes such that $k = 256$ from the two following BCH codes BCH-1 and BCH-2 (codes from [24]):

code	n	k	δ
BCH-1	1023	513	57
BCH-2	1023	483	60

We obtain the following shortened codes

code	n	k	δ
BCH-S1	766	256	57
BCH-S2	796	256	60

The shortened codes are obtained by subtracting 257 (and 227) from BCH-1 (from BCH-2):

- $\text{BCH-S1}[766 = 1023 - 257, 256 = 513 - 257, 57]$
- $\text{BCH-S2}[796 = 1023 - 227, 256 = 483 - 227, 60]$

Notice that shortening the BCH code does not affect the correcting capacity.

In our case, we will be working in \mathbb{F}_{2^m} , for that we use the primitive polynomial of degree $1 + X^3 + X^{10}$ to build this field (polynomial from [24]). We precomputed the generator polynomials for the two codes that we will be using in our implementation (BCH-S1 and BCH-S2) and we included their Hexadecimal formats in the file `parameters.h`.

1.5.3 Decoding BCH codes

We give a brief reminder on decoding BCH codes following [24]. Consider the BCH code defined by $[n, k, \delta]$, with $n = 2^m - 1$ ($m \geq 0$ of positive integer) and suppose that a code word $v(x) = v_0 + v_1x + \dots + v_{n-1}x^{n-1}$ is transmitted and that during transmission, error occurred in the following received vector:

$$r(x) = r_0 + r_1x + r_2x^2 + \dots + r_{n-1}x^{n-1}$$

We have that the location of errors are given by the error polynomial $e(x) = e_0 + e_1x + e_2x^2 + \dots + e_{n-1}x^{n-1}$, if $e_i = 1$, then there is an error occurred at that location. Then we can write

$$r(x) = v(x) + e(x)$$

We define the set of syndromes $S_1, S_2, \dots, S_{2\delta}$ as $S_i = r(\alpha^i)$, with α being the primitive element in \mathbb{F}_{2^m} . We have that $r(\alpha^i) = e(\alpha^i)$, since $v(\alpha^i) = 0$ (v is a code word). Suppose that $e(x)$ has t errors at locations j_1, \dots, j_t , then

$$e(x) = x^{j_1} + x^{j_2} + \dots + x^{j_t},$$

we obtain the following set of equations, where $\alpha^{j_1}, \alpha^{j_2}, \dots, \alpha^{j_t}$ are unknown:

$$\begin{aligned} S_1 &= \alpha^{j_1} + \alpha^{j_2} + \dots + \alpha^{j_t} \\ S_2 &= (\alpha^{j_1})^2 + (\alpha^{j_2})^2 + \dots + (\alpha^{j_t})^2 \\ S_3 &= (\alpha^{j_1})^3 + (\alpha^{j_2})^3 + \dots + (\alpha^{j_t})^3 \\ &\vdots \\ S_{2\delta} &= (\alpha^{j_1})^{2\delta} + (\alpha^{j_2})^{2\delta} + \dots + (\alpha^{j_t})^{2\delta} \end{aligned}$$

The goal of a BCH decoding algorithm is to solve this system of equations. We define the error location numbers by $\beta_i = \alpha^{j_i}$, which indicate the location of the errors. The equations above, can be expressed as follows:

$$\begin{aligned} S_1 &= \beta_1 + \beta_2 + \dots + \beta_t \\ S_2 &= \beta_1^2 + \beta_2^2 + \dots + \beta_t^2 \\ S_3 &= \beta_1^3 + \beta_2^3 + \dots + \beta_t^3 \\ &\vdots \\ S_{2\delta} &= \beta_1^{2\delta} + \beta_2^{2\delta} + \dots + \beta_t^{2\delta} \end{aligned}$$

we define the error location polynomial as:

$$\begin{aligned} \sigma(x) &= (1 + \beta_1x)(1 + \beta_2x) \dots (1 + \beta_tx) \\ &= 1 + \sigma_1x + \sigma_2x^2 + \dots + \sigma_tx^t \end{aligned}$$

We can see that, the roots of $\sigma(x)$ are $\beta_1^{-1}, \beta_2^{-1}, \dots, \beta_t^{-1}$ which are the inverses of the error location numbers. By inverting those roots we can construct the error polynomial $e(x)$.

We can summarize the decoding procedure of a BCH $[n, k, \delta]$ code by the following steps:

1. The first step is the computation of $2 \times \delta$ syndromes using the received polynomial
2. The second step is the computation of the error-location polynomial $\sigma(x)$ from the $2 \times \delta$ syndromes computed in the first step (in our implementation we will use the Simplified Berlekamp's Algorithm [23])
3. The third step is to find the error-location numbers by calculating the roots of the polynomial $\sigma(x)$ and returning their inverse (in our implementation we will be using the Chien search algorithm [10])
4. The fourth step is the correction of errors in the received polynomial

Remark 1.3. *As mentioned before, in our implementation, we deal with shortened BCH code. We notice that we will be using the same decoding procedure described above.*

Step 1. Syndrome computations. The following function computes the syndromes.

```
void syndrome_gen(syndrome_set* synd_set, gf_tables* tables, vector_u32* v); //
    bch.h
```

The syndromes are computed by evaluating the received polynomial stored in the vector v at the $2 \times \text{PARAM DELTA}$ consecutive roots of the generator polynomial α^i for $i = 1, 2, \dots, 2 \times \text{PARAM DELTA}$. Let us denote by $r(x)$ the polynomial in the vector v , thus the syndromes are

$$r(\alpha), r(\alpha^2), \dots, r(\alpha^{2 \times \text{PARAM DELTA}})$$

and they are stored as \mathbb{F}_{2^m} elements in the structure `synd set` which is the output of the function.

Step 2. Computing the Error-Location Polynomial. The following function computes the error location polynomial $\sigma(x)$ as defined above and store it in the vector `sigma`

```
void get_error_location_poly(sigma_poly* sigma, gf_tables* tables, syndrome_set*
    synd_set); // bch.h
```

This function implements the simplified Berlekamp's algorithm for finding the error location polynomial for binary **BCH** codes given by Joiner and Komo in [23].

Step 3. Finding the Error-Location Numbers. The following function computes the roots of the error location polynomial and finds their inverses which are the error location numbers.

```
void chien_search(uint16_t* error_pos, uint16_t* size, gf_tables* tables,
    sigma_poly* sigma); // bch.h
```

To find the roots of the polynomial $\sigma(x)$ stored in the structure `sigma`, we have to evaluate $\sigma(x)$ in all the elements of the Galois Field: let α be the generator of the field then we have to check for $j = 1, 2, \dots$ if $\sigma(\alpha^j) = 0$. Then if α^k is a root we store α^{-k} in the output array of the function. The Chien procedure permits to compute $\sigma(\alpha^{k+1})$ from $\sigma(\alpha^k)$, in fact :

- Suppose that σ is of degree t . If we have evaluated α^k , we obtain

$$\sigma(\alpha^k) = 1 + \sigma_1 \alpha^k + \sigma_2 \alpha^{2k} + \dots + \sigma_t \alpha^{tk}$$

- Then, we can obtain $\sigma(\alpha^{k+1})$ in $O(t)$ operation. In fact the i -th term in $\sigma(\alpha^{k+1})$ can be obtained from the i -th term of $\sigma(\alpha^k)$ by multiplying that term by α^i .

Suppose that we are using BCH $[n, k, \delta]$ one of the shortened BCH codes described bellow. Then, we have that the inverses of the roots of the elements α^i with $i \in \{1, \dots, 2^{10} - 1 - n\}$ will not be a valid error positions. In fact the location number obtained will be greater than n . For that it is useless to evaluate the error location polynomial $\sigma(x)$ in the element α^i for $i \in \{1, \dots, 2^{10} - 1 - n\}$. Therefore, in our implementation we starts the evaluation at α^i with $i = 2^{10} - n$.

Step 4. Error correction. To correct the errors in the received polynomial: we have to build the error polynomial $e(x)$ using the error location numbers obtained by the Chien search algorithm, then we add the error polynomial to the received polynomial. The following function builds $e(x)$ and store the result in the vector `e`

```
void error_poly_gen(vector_u32* e, uint16_t* error_pos, uint16_t size); // bch.h
```

1.5.4 Decryption Failure Probability

With a tensor product code $\mathcal{C} = \text{BCH}(n_1, k_1, \delta) \otimes \mathbb{1}_{n_2}$ as defined above, a decryption failure occurs whenever the decoding algorithm of the BCH code does not succeed in correcting errors that would have arisen after wrong decodings by the repetition code. Therefore, the analysis of the decryption failure probability is again split into three steps: evaluating the probability that the repetition code does not decode correctly, the conditional probability of a wrong decoding for the BCH code given an error weight and finally, the decryption failure probability using the law of total probability.

Step 1. We now focus on the probability that an error occurs while decoding the repetition code. As shown in Sec. 1.4, the probability for a coordinate of $\mathbf{e}' = \mathbf{x} \cdot \mathbf{r}_2 - \mathbf{r}_1 \cdot \mathbf{y} + \mathbf{e}$ to be 1 is

p^* (see Eq. (18)). As mentioned above, $\mathbb{1}_{n_2}$ can decode up to $\delta_2 = \lfloor \frac{n_2-1}{2} \rfloor$ errors. Therefore, assuming that the error vector \mathbf{e}' has weight γ (which occurs with the probability given in Eq. (19)), the probability of getting a decoding error on a single block of the repetition code $\mathbb{1}_{n_2}$ is hence given by:

$$\bar{p}_\gamma = \bar{p}_\gamma(n_1, n_2) = \sum_{i=\lfloor \frac{n_2-1}{2} \rfloor + 1}^{n_2} \binom{n_2}{i} \left(\frac{\gamma}{n_1 n_2} \right)^i \left(1 - \frac{\gamma}{n_1 n_2} \right)^{n_2-i}. \quad (22)$$

Step 2. We now focus on the $\text{BCH}(n_1, k_1, \delta_1)$ code, and recall that it can correct up to δ_1 errors. Now the probability \mathcal{P} that the $\text{BCH}(n_1, k_1, \delta_1)$ code fails to decode correctly the encoded message \mathbf{m}_1 back to \mathbf{m} is given by the probability that an error occurred on at least $\delta_1 + 1$ blocks of the repetition code. Therefore, we have

$$\mathcal{P} = \mathcal{P}(\delta_1, n_1, n_2, \gamma) = \sum_{i=\delta_1+1}^{n_1} \binom{n_1}{i} (\bar{p}_\gamma)^i (1 - \bar{p}_\gamma)^{n_1-i}. \quad (23)$$

Step 3. Finally, using the law of total probability, we have that the decryption failure probability is given by the sum *over all the possible weights* of the probability that the error has this specific weight times the probability of a decoding error for this weight. This is captured in the following theorem, whose proof is a straightforward consequence of the formulae of Sec. 1.4 and 1.5.1.

Theorem 1.4. *Let $\mathcal{C} = \text{BCH}(n_1, k_1, \delta) \otimes \mathbb{1}_{n_2}$, $(\text{pk}, \text{sk}) \leftarrow \text{KeyGen}$, and $\mathbf{m} \xleftarrow{\$} \mathbb{F}_2^{k_1}$, then with the notations above, the decryption failure probability is*

$$p_{\text{fail}} = \Pr[\text{Decrypt}(\text{sk}, \text{Encrypt}(\text{pk}, \mathbf{m})) \neq \mathbf{m}] \quad (24)$$

$$= \sum_{\gamma=0}^{\min(2 \times w \times w_{\mathbf{r}} + w_{\mathbf{e}}, n_1 n_2)} \Pr[\omega(\mathbf{e}') = \gamma] \mathcal{P}(\delta_1, n_1, n_2, \gamma) \quad (25)$$

Eq. (25) gives a theoretical approximation of the decryption failure rate. The parameters presented in Tab. 1 were obtained using this formula. Experimental evidences supporting the validity of the assumptions made to obtain this formula are provided in Fig. 4.

1.6 Parameters

In this section, we specify which codes are used for our HQC and give concrete sets of parameters. As mentioned in the previous section, we use a tensor product code (Def. 1.5.1) $\mathcal{C} = \text{BCH}(n_1, k, \delta) \otimes \mathbb{1}_{n_2}$. A message $\mathbf{m} \in \mathbb{F}^k$ is encoded into $\mathbf{m}_1 \in \mathbb{F}^{n_1}$ with the BCH code, then each coordinate $\mathbf{m}_{1,i}$ of \mathbf{m}_1 is encoded into $\tilde{\mathbf{m}}_{1,i} \in \mathbb{F}^{n_2}$ with $\mathbb{1}_{n_2}$. To match the

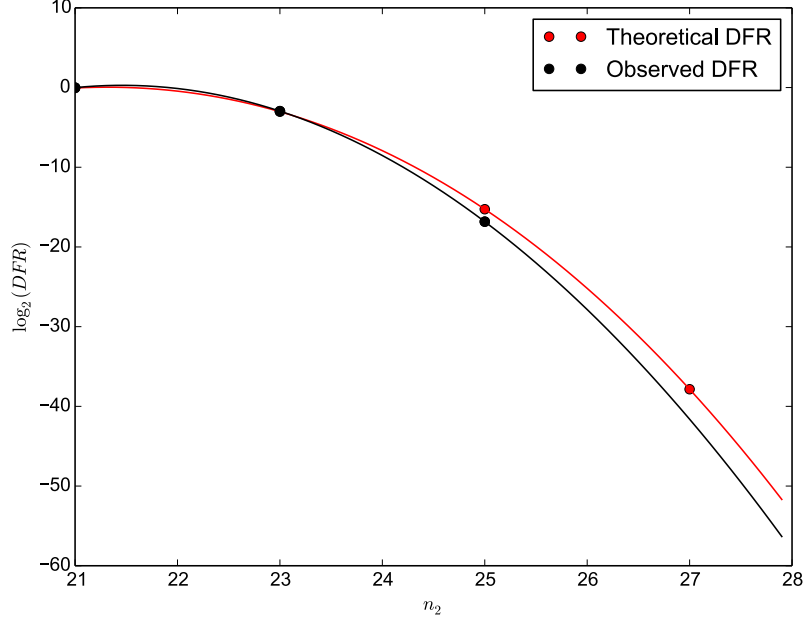


Figure 4: Logarithm of theoretical and observed decryption failure rates (DFR). The red curve corresponding to theoretical DFR was obtained using Eq. (25) while the black curve corresponding to experimental DFR was obtained by running 10^5 encryption/decryption over 10^3 codes with $n_1 = 766$, $k_1 = 256$, $\delta_1 = 57$, $w = 67$, $w_r = 77$. The parameters have been selected to make the theoretical DFR sufficiently high to compare it to experiments. Finally, the curves have been interpolated to the second order on the logarithm of the probability.

description of our cryptosystem in Sec. 1.3, we have $\mathbf{mG} = \tilde{\mathbf{m}} = (\tilde{\mathbf{m}}_{1,1}, \dots, \tilde{\mathbf{m}}_{1,n_1}) \in \mathbb{F}^{n_1 n_2}$. To obtain the ciphertext, $\mathbf{r} = (\mathbf{r}_1, \mathbf{r}_2) \xleftarrow{\$} \mathcal{R}^2$ and $\mathbf{e} \xleftarrow{\$} \mathcal{R}$ are generated and the encryption of \mathbf{m} is $\mathbf{c} = (\mathbf{u} = \mathbf{r}_1 + \mathbf{h} \cdot \mathbf{r}_2, \mathbf{v} = \mathbf{mG} + \mathbf{s} \cdot \mathbf{r}_2 + \mathbf{e})$.

We propose several sets of parameters, targeting different levels of security. According to NIST, it may be assumed that an attacker can only make 2^{64} queries to the decryption oracle. In this sense, we propose several decryption failure rates ranging from 2^{-64} to $2^{-\lambda}$ where λ is the security parameter. The proposed sets of parameters cover security categories 1, 3, and 5 (for respectively 128, 192, and 256 bits of security). For each parameter set, the parameters are chosen so that the minimal workfactor of the best known attack exceeds the security parameter. For classical attacks, best known attacks include the works from [8, 6, 14, 3] and for quantum attacks, the work of [5]. We consider $w = \mathcal{O}(\sqrt{n})$ and follow the complexity described in [9] (see Sec. 5 for more details).

In Tab. 1, n_1 denotes the length of the BCH code, n_2 the length of the repetition code 1 so that the length of the tensor product code \mathcal{C} is $n \approx n_1 n_2$ (actually the smallest primitive

prime greater than $n_1 n_2$). k is the dimension of the BCH code and hence also the dimension of \mathcal{C} . δ is the decoding capability of the BCH code, *i.e.* the maximum number of errors that the BCH can decode. w is the weight of the n -dimensional vectors \mathbf{x} , \mathbf{y} , $w_{\mathbf{r}}$ the weight of \mathbf{r}_1 , and \mathbf{r}_2 and similarly $w_{\mathbf{e}} = \omega(\mathbf{e})$ for our cryptosystem.

Instance	n_1	n_2	n	k	δ	w	$w_{\mathbf{r}} = w_{\mathbf{e}}$	security	p_{fail}
Basic-I	766	29	22,229	256	57	67	77	128	$< 2^{-64}$
Basic-II	766	31	23,747	256	57	67	77	128	$< 2^{-96}$
Basic-III	796	31	24,677	256	60	67	77	128	$< 2^{-128}$
Advanced-I	796	51	40,597	256	60	101	117	192	$< 2^{-64}$
Advanced-II	766	57	43,669	256	57	101	117	192	$< 2^{-128}$
Advanced-III	766	61	46,747	256	57	101	117	192	$< 2^{-192}$
Paranoiac-I	766	77	59,011	256	57	133	153	256	$< 2^{-64}$
Paranoiac-II	766	83	63,587	256	57	133	153	256	$< 2^{-128}$
Paranoiac-III	796	85	67,699	256	60	133	153	256	$< 2^{-192}$
Paranoiac-IV	796	89	70,853	256	60	133	153	256	$< 2^{-256}$

Table 1: Parameter sets for our cryptosystem in Hamming metric. The tensor product code used is $\mathcal{C} = \text{BCH}(n_1, k_1, \delta_1) \otimes \mathbf{1}_{n_2}$ (see Sec. 1.5.1). The considered BCH codes are initially of length 1023, then shortened to support 256 bits dimension (see Sec. 1.5.2). For the resulting public key, secret key and ciphertext sizes, please see Tab. 2 below. One may use seeds to shorten keys thus obtaining sizes presented in Tab. 3. The aforementioned sizes are the ones used in our reference implementation except that we also concatenate the public key within the secret key in order to respect the NIST API.

Computational costs of the system. For encryption the main cost is a product of a cyclic matrix of size n with a vector of weight $\mathcal{O}(\sqrt{n})$. Using the Fourier transform the asymptotical cost is in $\mathcal{O}(n \log(n))$ but for our range of parameters, taking into account the weight $\mathcal{O}(\sqrt{n})$ allows to obtain a cost in $\mathcal{O}(n^{\frac{3}{2}})$ which is better in practice than what is obtained with Fourier transform. For decryption, there is always the cost of a matrix times a small vector in $\mathcal{O}(n^{\frac{3}{2}})$, plus the cost of decoding. For our proposition the decoding consists in a repetition code of length n_2 and the decoding of BCH code of length n_1 ($766 \leq n_1 \leq 796$), the cost of the repetition code decoding is hence linear, when the cost of the BCH is quadratic in the length n_1 of the BCH code. Overall the main cost remains the computation of the matrix-vector product in $\mathcal{O}(n^{\frac{3}{2}})$.

2 Performance Analysis

In this section, we provide concrete performance measures of our implementation. For each parameter set, results have been obtained by running 100,000 random instances and computing their average execution time. The benchmarks have been performed on a machine

Instance	pk size	sk size	ct size	ss size
Basic-I	5,558	252	5,622	64
Basic-II	5,938	252	6,001	64
Basic-III	6,170	252	6,234	64
Advanced-I	10,150	404	10,214	64
Advanced-II	10,918	404	10,982	64
Advanced-III	11,688	404	11,752	64
Paranoiac-I	14,754	532	14,818	64
Paranoiac-II	15,898	532	15,962	64
Paranoiac-III	16,926	566	16,990	64
Paranoiac-IV	17,714	566	17,778	64

Table 2: Resulting theoretical sizes in bytes for HQC. The public key **pk** is composed of (**h**, **s**) and has size $2n$ bits. The secret key **sk** is composed of (**x**, **y**) and has size $2w\lceil\log_2(n)\rceil$ bits. The ciphertext **ct** is composed of (**u**, **v**, **d**) and has size $2n + 512$ bits. The shared secret **ss** is composed of K and has size 512 bits (SHA512 output size).

Instance	pk size	sk size	ct size	ss size
Basic-I	2,819	40	5,622	64
Basic-II	3,009	40	6,002	64
Basic-III	3,125	40	6,234	64
Advanced-I	5,115	40	10,214	64
Advanced-II	5,499	40	10,982	64
Advanced-III	5,884	40	11,752	64
Paranoiac-I	7,417	40	14,818	64
Paranoiac-II	7,989	40	15,962	64
Paranoiac-III	8,503	40	16,990	64
Paranoiac-IV	8,897	40	17,778	64

Table 3: Resulting sizes in bytes for HQC using NIST seed expander initialized with 40 bytes long seeds. The public key **pk** is composed of (**seed1**, **s**) and has size $320 + n$ (in bits). The secret key **sk** is composed of (**seed2**) and has size 320 (in bits). The ciphertext **ct** is composed of (**u**, **v**, **d**) and has size $2n + 512$ (in bits). The shared secret **ss** is composed of K and has size 512 bits (SHA512 output size).

running Ubuntu 16.04 LTS. The latter has 32GB of memory and an Intel® Core™ i7-4770 CPU @ 3.4GHz for which the Hyper-Threading, Turbo Boost and SpeedStep features were disabled. The scheme have been compiled with gcc (version 7.2.0) using the compilation flags `-O3 -std=c99 -pedantic`. The following third party library have been used: `openssl` (version 1.1.0f).

2.1 Reference Implementation

The performances of our reference implementation on the aforementioned benchmark platform are described in Tab. 4 (timings in ms) and Tab. 5 (millions of CPU cycles required).

Instance	KeyGen	Encaps	Decaps
Basic-I	0.17	0.36	0.57
Basic-II	0.18	0.38	0.61
Basic-III	0.19	0.40	0.63
Advanced-I	0.37	0.77	1.13
Advanced-II	0.40	0.83	1.21
Advanced-III	0.43	0.89	1.28
Paranoiac-I	0.65	1.38	1.96
Paranoiac-II	0.76	1.60	2.22
Paranoiac-III	0.78	1.65	2.35
Paranoiac-IV	0.82	1.76	2.50

Table 4: Timings (in ms) of the reference implementation for different instances of HQC.

Instance	KeyGen	Encaps	Decaps
Basic-I	0.57	1.22	1.95
Basic-II	0.61	1.28	2.07
Basic-III	0.63	1.35	2.15
Advanced-I	1.26	2.61	3.82
Advanced-II	1.37	2.81	4.11
Advanced-III	1.47	3.02	4.35
Paranoiac-I	2.21	4.67	6.67
Paranoiac-II	2.52	5.37	7.51
Paranoiac-III	2.66	5.62	8.03
Paranoiac-IV	2.81	5.95	8.46

Table 5: Millions of cycles for the reference implementation for different instances of HQC.

2.2 Optimized Implementation

No optimized implementation has been provided. As a consequence, the folders `Optimized_Implementation/` and `Reference_Implementation/` are identical. Additional implementation (optimized variant using vectorization, constant-time implementation...) might be provided later.

3 Known Answer Test Values

Known Answer Test (KAT) values have been generated using the script provided by the NIST. They are available in the folder `KAT/Reference_Implementation/`. As mentioned in Sec. 2.2, since the reference and optimized implementations are identical, `KAT/Optimized_Implementation/` is just a copy of `KAT/Reference_Implementation/`.

In addition, we provide, for each parameter set, an example with *intermediate values* in the folder `KAT/Reference_Implementation/`.

Notice that one can generate the aforementioned test files using respectively the `kat` and `verbose` modes of our implementation. The procedure to follow in order to do so is detailed in the technical documentation.

4 Security

In this section we prove the security of our encryption scheme viewed as a PKE scheme (IND-CPA). The security of the KEM/DEM version is provided by the transformation described in [21], and the tightness of the reduction provided by this transformation has been discussed at the end of Sec. 1.2.

Theorem 4.1. *The scheme presented above is IND-CPA under the 2-DQCSD and 3-DQCSD assumptions.*

Proof. To prove the security of the scheme, we are going to build a sequence of games transitioning from an adversary receiving an encryption of message \mathbf{m}_0 to an adversary receiving an encryption of a message \mathbf{m}_1 and show that if the adversary manages to distinguish one from the other, then we can build a simulator breaking the DQCSD assumption, for QC codes of index 2 or 3 (codes with parameters $[2n, n]$ or $[3n, n]$), and running in approximately the same time.

Game G_1 : This is the real game, which we can state algorithmically as follows:

- Game $_{\mathcal{E}, \mathcal{A}}^1(\lambda)$**
1. $\text{param} \leftarrow \text{Setup}(1^\lambda)$
 2. $(\text{pk}, \text{sk}) \leftarrow \text{KeyGen}(\text{param})$ with $\text{pk} = (\mathbf{h}, \mathbf{s} = \text{sk} \cdot \mathbf{h}^\top)$
 3. $(\mathbf{m}_0, \mathbf{m}_1) \leftarrow \mathcal{A}(\text{FIND} : \text{pk})$
 4. $\mathbf{c}^* \leftarrow \text{Encrypt}(\text{pk}, \mathbf{m}_0)$
 5. $\mathbf{b}' \leftarrow \mathcal{A}(\text{GUESS} : \mathbf{c}^*)$
 6. RETURN \mathbf{b}'

Game G_2 : In this game we start by forgetting the decryption key sk , and taking \mathbf{s} at random, and then proceed honestly:

Game $_{\mathcal{E},\mathcal{A}}^2(\lambda)$

1. $\text{param} \leftarrow \text{Setup}(1^\lambda)$
- 2a. $(\text{pk}, \text{sk}) \leftarrow \text{KeyGen}(\text{param})$ with $\text{pk} = (\mathbf{h}, \mathbf{s} = \text{sk} \cdot \mathbf{h}^\top)$
- 2b. $\mathbf{s} \xleftarrow{\$} \mathcal{R}$
- 2c. $(\text{pk}, \text{sk}) \leftarrow ((\mathbf{h}, \mathbf{s}), \mathbf{0})$
3. $(\mathbf{m}_0, \mathbf{m}_1) \leftarrow \mathcal{A}(\text{FIND} : \text{pk})$
4. $\mathbf{c}^* \leftarrow \text{Encrypt}(\text{pk}, \mathbf{m}_0)$
5. $\mathbf{b}' \leftarrow \mathcal{A}(\text{GUESS} : \mathbf{c}^*)$
6. RETURN \mathbf{b}'

The adversary has access to pk and \mathbf{c}^* . As he has access to pk and the **Encrypt** function, anything that is computed from pk and \mathbf{c}^* can also be computed from just pk . Moreover, the distribution of \mathbf{c}^* is independent of the game we are in, and therefore we can suppose the only input of the adversary is pk . Suppose he has an algorithm \mathcal{D}_λ , taking pk as input, that distinguishes with advantage ϵ **Game** \mathbf{G}_1 and **Game** \mathbf{G}_2 , for some security parameter λ . Then he can also build an algorithm $\mathcal{D}'_{\mathcal{E},\mathcal{D}_\lambda}$ which solves the 2-DQCSD(n, w) problem for parameters (n, w) resulting from **Setup**(λ), with the same advantage ϵ , when given as input a challenge $(\mathbf{H}, \mathbf{y}^\top) \in \mathbb{F}^{n \times 2n} \times \mathbb{F}^n$.

$\mathcal{D}'_{\mathcal{E},\mathcal{D}_\lambda}((\mathbf{H}, \mathbf{y}^\top))$

1. Set $\text{param} \leftarrow \text{Setup}(\lambda)$ and get \mathbf{G} from **KeyGen**(param)
2. $\text{pk} \leftarrow (\mathbf{h}, \mathbf{y})$
2. $b' \leftarrow \mathcal{D}_\lambda(\text{pk})$
4. If $b' == 0$ output QCS
5. If $b' == 1$ output UNIFORM

Note that if we define pk as (\mathbf{h}, \mathbf{y}) with \mathbf{G} generated by **KeyGen**(n, k, δ, w) and $(\mathbf{H}, \mathbf{y}^\top)$ from a 2-QCSD(n, w) distribution pk follows exactly the same distribution as in **Game** \mathbf{G}_1 . On the other hand if $(\mathbf{H}, \mathbf{y}^\top)$ comes from a uniform distribution, pk follows exactly the same distribution as in **Game** \mathbf{G}_2 . Thus we have

$$\Pr [\mathcal{D}'_{\mathcal{E},\mathcal{D}_\lambda}((\mathbf{h}, \mathbf{y}^\top)) = \text{QCS} | (\mathbf{h}, \mathbf{y}^\top) \leftarrow 2\text{-QCSD}(n, w)] = \Pr [\mathcal{D}_\lambda(\text{pk}) = 0 | \text{pk from } \mathbf{Game}_{\mathcal{E},\mathcal{A}}^0(\lambda)]$$

and

$$\Pr [\mathcal{D}'_{\mathcal{E},\mathcal{D}_\lambda}((\mathbf{h}, \mathbf{y}^\top)) = \text{UNIFORM} | (\mathbf{h}, \mathbf{y}^\top) \leftarrow 2\text{-QCSD}(n, w)] = \Pr [\mathcal{D}_\lambda(\text{pk}) = 1 | \text{pk from } \mathbf{Game}_{\mathcal{E},\mathcal{A}}^0(\lambda)]$$

And similarly when $(\mathbf{h}, \mathbf{y}^\top)$ is uniform the probabilities of $\mathcal{D}'_{\mathcal{E},\mathcal{D}_\lambda}$ outputs match those of \mathcal{D}_λ when pk is from $\mathbf{Game}_{\mathcal{E},\mathcal{A}}^1(\lambda)$. The advantage of $\mathcal{D}'_{\mathcal{E},\mathcal{D}_\lambda}$ is therefore equal to the advantage of \mathcal{D}_λ .

Game \mathbf{G}_3 : Now that we no longer know the decryption key, we can start generating random ciphertexts. So instead of picking correctly weighted $\mathbf{r}_1, \mathbf{r}_2, \mathbf{e}$, the simulator now picks random vectors in the full space.

Game $_{\mathcal{E},\mathcal{A}}^3(\lambda)$

1. $\text{param} \leftarrow \text{Setup}(1^\lambda)$
- 2a. $(\text{pk}, \text{sk}) \leftarrow \text{KeyGen}(\text{param})$ with $\text{pk} = (\mathbf{h}, \mathbf{s} = \text{sk} \cdot \mathbf{h}^\top)$
- 2b. $\mathbf{s} \xleftarrow{\$} \mathcal{R}$
- 2c. $(\text{pk}, \text{sk}) \leftarrow ((\mathbf{h}, \mathbf{s}), \mathbf{0})$
3. $(\mathbf{m}_0, \mathbf{m}_1) \leftarrow \mathcal{A}(\text{FIND} : \text{pk})$
- 4a. Generate $\mathbf{e} \xleftarrow{\$} \mathcal{R}$, $\mathbf{r} = (\mathbf{r}_1, \mathbf{r}_2) \xleftarrow{\$} \mathcal{R}^2$ uniformly at random
- 4b. $\mathbf{u}^\top \leftarrow \mathbf{H}\mathbf{r}^\top$ and $\mathbf{v} \leftarrow \mathbf{m}_0\mathbf{G} + \mathbf{s} \cdot \mathbf{r}_2 + \mathbf{e}$
- 4c. $\mathbf{c}^* \leftarrow (\mathbf{u}, \mathbf{v})$
5. $\mathbf{b}' \leftarrow \mathcal{A}(\text{GUESS} : \mathbf{c}^*)$
6. RETURN \mathbf{b}'

As we have

$$(\mathbf{u}, \mathbf{v} - \mathbf{m}_0\mathbf{G})^\top = \begin{pmatrix} \mathbf{I}_n & \mathbf{0} & \text{rot}(\mathbf{h}) \\ \mathbf{0} & \mathbf{I}_n & \text{rot}(\mathbf{s}) \end{pmatrix} \cdot (\mathbf{r}_1, \mathbf{e}, \mathbf{r}_2)^\top,$$

the difference between Game \mathbf{G}_2 and Game \mathbf{G}_3 is that in the former

$$\left(\begin{pmatrix} \mathbf{I}_n & \mathbf{0} & \text{rot}(\mathbf{h}) \\ \mathbf{0} & \mathbf{I}_n & \text{rot}(\mathbf{s}) \end{pmatrix}, (\mathbf{u}, \mathbf{v} - \mathbf{m}_0\mathbf{G})^\top \right)$$

follows the 3-QCSD distribution (for a $2n \times 3n$ QC matrix of index 3), and in the latter it follows a uniform distribution (as \mathbf{r}_1 and \mathbf{e} are uniformly distributed and independently chosen One-Time Pads).

Note that an adversary is not able to obtain \mathbf{c}^* from pk any more, as depending on which game we are \mathbf{c}^* is generated differently. The input of a game distinguisher will therefore be $(\text{pk}, \mathbf{c}^*)$. As it must interact with the challenger as usually we suppose it has two access modes **FIND** and **GUESS** to process first pk and later \mathbf{c}^* .

Suppose the adversary is able to distinguish Game \mathbf{G}_2 and Game \mathbf{G}_3 , with a distinguisher \mathcal{D}_λ , which takes as input $(\text{pk}, \mathbf{c}^*)$ and outputs a guess $b' \in \{1, 2\}$ of the game we are in.

Again, we can build a distinguisher $\mathcal{D}'_{\mathcal{E}, \mathcal{D}_\lambda}$ that will break the 3-DQCSD(n, w) assumption for parameters (n, w) from $\text{Setup}(1^\lambda)$ with the same advantage as the game distinguisher, when given an input $(\mathbf{H}, \mathbf{y}^\top) \in \mathbb{F}^{2n \times 3n} \times \mathbb{F}^{2n}$. In the 3-DQCSD(n, w) problem, matrix \mathbf{H} is assumed to be of the form

$$\begin{pmatrix} \mathbf{I}_n & \mathbf{0} & \text{rot}(\mathbf{a}) \\ \mathbf{0} & \mathbf{I}_n & \text{rot}(\mathbf{b}) \end{pmatrix}.$$

In order to use explicitly \mathbf{a} and \mathbf{b} we note the matrix $\mathbf{H}_{\mathbf{a}, \mathbf{b}}$ instead of just \mathbf{H} . We will also note $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2)$.

- $$\mathcal{D}'_{\mathcal{E}, \mathcal{D}_\lambda}((\mathbf{H}_{\mathbf{a}, \mathbf{b}}, (\mathbf{y}_1, \mathbf{y}_2)^\top))$$
1. $\text{param} \leftarrow \text{Setup}(1^\lambda)$
 - 2a. $(\text{pk}, \text{sk}) \leftarrow \text{KeyGen}(\text{param})$ with $\text{pk} = (\mathbf{h}, \mathbf{s})$
 - 2b. $(\text{pk}, \text{sk}) \leftarrow ((\mathbf{G}, (\mathbf{I}_n \text{rot}(\mathbf{a})), \mathbf{b}), \mathbf{0})$
 3. $(\mathbf{m}_0, \mathbf{m}_1) \leftarrow \mathcal{D}_\lambda(\text{FIND} : \text{pk})$
 4. $\mathbf{u} \leftarrow \mathbf{y}_1, \mathbf{v} \leftarrow \mathbf{m}_0 \mathbf{G} + \mathbf{y}_2$ and $\mathbf{c}^* \leftarrow (\mathbf{u}, \mathbf{v})$
 5. $\mathbf{b}' \leftarrow \mathcal{D}_\lambda(\text{GUESS} : \mathbf{c}^*)$
 4. If $\mathbf{b}' == \mathbf{1}$ output QCS
 5. If $\mathbf{b}' == \mathbf{2}$ output UNIFORM

The distribution of pk is unchanged with respect to the games as the first matrix is from KeyGen , the second matrix follows the same distribution as in KeyGen , and the vectors \mathbf{b} and \mathbf{s} are both uniformly chosen. If $(\mathbf{H}_{\mathbf{a}, \mathbf{b}}, (\mathbf{y}_1, \mathbf{y}_2)^\top)$ follows the 3-QCSD(n, w) distribution, then

$$(\mathbf{y}_1, \mathbf{y}_2)^\top = \begin{pmatrix} \mathbf{I}_n & \mathbf{0} & \text{rot}(\mathbf{a}) \\ \mathbf{0} & \mathbf{I}_n & \text{rot}(\mathbf{b}) \end{pmatrix} \cdot (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)^\top$$

with $\omega(\mathbf{x}_1) = \omega(\mathbf{x}_2) = \omega(\mathbf{x}_3) = w$. Thus, \mathbf{c}^* follows the same distribution as in Game \mathbf{G}_2 . If $(\mathbf{H}_{\mathbf{a}, \mathbf{b}}, (\mathbf{y}_1, \mathbf{y}_2)^\top)$ follows an uniform distribution, then \mathbf{c}^* follows the same distribution as in Game \mathbf{G}_3 . We obtain therefore the same equalities for the output probabilities of $\mathcal{D}'_{\mathcal{E}, \mathcal{D}_\lambda}$ and \mathcal{D}_λ as with the previous games and therefore the advantages of both distinguishers are equal.

Game \mathbf{G}_4 : We now encrypt the other plaintext. We chose $\mathbf{r}'_1, \mathbf{r}'_2, \mathbf{e}'$ uniformly and set $\mathbf{u}^\top = \mathbf{h} \mathbf{r}'^\top$ and $\mathbf{v} = \mathbf{m}_1 \mathbf{G} + \mathbf{s} \cdot \mathbf{r}'_2 + \mathbf{e}'$. This is the last game we describe explicitly, since, even if it is a mirror of Game \mathbf{G}_3 , it involves a new proof.

- $$\text{Game}_{\mathcal{E}, \mathcal{A}}^4(\lambda)$$
1. $\text{param} \leftarrow \text{Setup}(1^\lambda)$
 - 2a. $(\text{pk}, \text{sk}) \leftarrow \text{KeyGen}(\text{param})$ with $\text{pk} = (\mathbf{h}, \mathbf{s})$
 - 2b. $\mathbf{s} \xleftarrow{\$} \mathcal{R}$
 - 2c. $(\text{pk}, \text{sk}) \leftarrow ((\mathbf{h}, \mathbf{s}), \mathbf{0})$
 3. $(\mathbf{m}_0, \mathbf{m}_1) \leftarrow \mathcal{A}(\text{FIND} : \text{pk})$
 - 4a. Generate $\mathbf{e}' \xleftarrow{\$} \mathcal{R}, \mathbf{r} = (\mathbf{r}'_1, \mathbf{r}'_2) \xleftarrow{\$} \mathcal{R}^2$ uniformly at random
 - 4b. $\mathbf{u}^\top \leftarrow \mathbf{Q} \mathbf{r}'^\top$ and $\mathbf{v} \leftarrow \mathbf{m}_1 \mathbf{G} + \mathbf{s} \cdot \mathbf{r}'_2 + \mathbf{e}'$
 - 4c. $\mathbf{c}^* \leftarrow (\mathbf{u}, \mathbf{v})$
 5. $\mathbf{b}' \leftarrow \mathcal{A}(\text{GUESS} : \mathbf{c}^*)$
 6. RETURN \mathbf{b}'

The outputs from Game \mathbf{G}_3 and Game \mathbf{G}_4 follow the exact same distribution, and therefore the two games are indistinguishable from an information-theoretic point of view. Indeed, for each tuple (\mathbf{r}, \mathbf{e}) of Game \mathbf{G}_3 , resulting in a given (\mathbf{u}, \mathbf{v}) , there is a one to one mapping to a couple $(\mathbf{r}', \mathbf{e}')$ resulting in Game \mathbf{G}_4 in the *same* (\mathbf{u}, \mathbf{v}) ,

namely $\mathbf{r}' = \mathbf{r}$ and $\mathbf{e}' = \mathbf{m}_0\mathbf{G} + \mathbf{m}_1\mathbf{G}$. This implies that choosing uniformly (\mathbf{r}, \mathbf{e}) in Game \mathbf{G}_3 and choosing uniformly $(\mathbf{r}', \mathbf{e}')$ in Game \mathbf{G}_4 leads to the same output distribution for (\mathbf{u}, \mathbf{v}) .

Game \mathbf{G}_5 : In this game, we now pick $\mathbf{r}'_1, \mathbf{r}'_2, \mathbf{e}'$ with the correct weight.

Game \mathbf{G}_6 : We now conclude by switching the public key to an honestly generated one.

We do not explicit these last two games as Game \mathbf{G}_4 and Game \mathbf{G}_5 are the equivalents of Game \mathbf{G}_3 and Game \mathbf{G}_2 except that \mathbf{m}_1 is used instead of \mathbf{m}_0 . A distinguisher between these two games breaks therefore the 3-DQCSD assumption too. Similarly Game \mathbf{G}_5 and Game \mathbf{G}_6 are the equivalents of Game \mathbf{G}_2 and Game \mathbf{G}_1 and a distinguisher between these two games breaks the 2-DQCSD assumption.

We managed to build a sequence of games allowing a simulator to transform a ciphertext of a message \mathbf{m}_0 to a ciphertext of a message \mathbf{m}_1 . Hence, the advantage of an adversary against the IND-CPA experiment is bounded as:

$$\text{Adv}_{\mathcal{E}, \mathcal{A}}^{\text{ind}}(\lambda) \leq 2 \left(\text{Adv}^{2\text{-DQCSD}}(\lambda) + \text{Adv}^{3\text{-DQCSD}}(\lambda) \right). \quad (26)$$

□

5 Known Attacks

The practical complexity of the SD problem for the Hamming metric has been widely studied for more than 50 years. Most efficient attacks are based on Information Set Decoding, a technique first introduced by Prange in 1962 [30] and improved later by Stern [32], then Dumer [13]. Recent works [26, 3, 27] suggest a complexity of order $2^{cw(1+\text{negl}(1))}$, for some constant c . A particular work focusing on the regime $w = \text{negl}(n)$ confirms this formula, with a close dependence between c and the rate k/n of the code being used [9].

Specific structural attacks. Quasi-cyclic codes have a special structure which may potentially open the door to specific structural attacks. A first generic attack is the DOOM attack [31] which because of cyclicity implies a gain of $\mathcal{O}(\sqrt{n})$ (when the gain is in $\mathcal{O}(n)$ for MDPC codes, since the code is generated by a small weight vector basis). It is also possible to consider attacks on the form of the polynomial generating the cyclic structure. Such attacks have been studied in [20, 25, 31], and are especially efficient when the polynomial $x^n - 1$ has many low degree factors. These attacks become inefficient as soon as $x^n - 1$ has only two irreducible factors of the form $(x - 1)$ and $x^{n-1} + x^{n-2} + \dots + x + 1$, which is the case when n is prime and q generates the multiplicative group $(\mathbb{Z}/n\mathbb{Z})^*$. Such numbers are known up to very large values. We consider such n for our parameters.

Parameters and tightness of the reduction. We proposed different sets of parameters in Sec. 1.6: basic, advanced, and paranoid which respectively provide 128 (category 1), 192 (category 3), and 256 (category 5) bits of classical (*i.e.* pre-quantum) security. The

quantum-safe security is obtained by dividing the security bits by two (taking the square root of the complexity) [5]. For each security level, we provide different decryption failure rates to better adapt to the adversary computing power. Notice that even if the adversary has access to a quantum computer, this *does not* change the decryption failure rate.³ Best known attacks include the works from [8, 6, 14, 26, 3, 27] and for quantum attacks, the work of [5]. In the setting $w = \mathcal{O}(\sqrt{n})$, best known attacks have a complexity in $2^{-t \ln(1-R)(1+o(1))}$ where $t = \mathcal{O}(w)$ and R is the rate of the code [9]. In our configuration, we have $t = 2w$ and $R = 1/2$ for the reduction to the 2-DQCSD problem, and $t = 3w_{\mathbf{r}}$ and $R = 1/3$ for the 3-DQCSD problem. By taking into account the DOOM attack [31], and also the fact that we consider balanced vectors (\mathbf{x}, \mathbf{y}) and $(\mathbf{r}_1, \mathbf{e}, \mathbf{r}_2)$ for the attack (which costs only a very small factor, since random words have a good probability to be balanced on each block), we need to divide this complexity by approximately \sqrt{n} (up to polylog factor). The term $o(1)$ is respectively $\log \left(\binom{n}{w}^2 / \binom{2n}{2w} \right)$ and $\log \left(\binom{n}{w_{\mathbf{r}}}^3 / \binom{3n}{3w_{\mathbf{r}}} \right)$ for the 2-DQCSD and 3-DQCSD problems. Overall our security reduction is tight corresponding to generic instances of the classical 2-DQCSD and 3-DQCSD problems according to the best attacks of [9].

6 Advantages and Limitations

6.1 Advantages

The main advantages of HQC over existing code-based cryptosystems are:

- its IND-CPA reduction to a well-understood problem on coding theory: the Quasi-Cyclic Syndrome Decoding problem,
- its immunity against attacks aiming at recovering the hidden structure of the code being used,
- close estimations of its decryption failure rate.

The last item allows to achieve a tight reduction for the IND-CCA2 security of the KEM-DEM version through the recent transformation of [21].

6.2 Limitations

We have proposed an instantiation of HQC using BCH codes tensored with repetition codes. As seen above, this construction presents the major advantage of making possible and easy to conduct a study of the error vector distribution, yielding a good estimation of the decryption failure rate. HQC might be more efficient with other families of codes, but another analysis would have to be done.

³We do not consider the very strong adversarial model where the adversary is given access to a *quantum* decryption oracle.

A first limitation to our cryptosystem (at least for the PKE version) is the low encryption rate. It is possible to encrypt 256 bits of plaintext as required by NIST, but increasing this rate also increases the parameters.

As a more general limitation and in contrast with lattices and the so-called Ring Learning With Errors problem, code-based cryptography does not benefit from search to decision reduction for structured codes.

References

- [1] Carlos Aguilar Melchor, Olivier Blazy, Jean Christophe Deneuville, Philippe Gaborit, and Gilles Zémor. Efficient encryption from random quasi-cyclic codes. *CoRR*, abs/1612.05572, 2016. <http://arxiv.org/abs/1612.05572>. 4
- [2] Benny Applebaum, Yuval Ishai, and Eyal Kushilevitz. Cryptography with constant input locality. In Alfred Menezes, editor, *CRYPTO 2007*, volume 4622 of *LNCS*, pages 92–110. Springer, Heidelberg, August 2007. <https://www.iacr.org/archive/crypto2007/46220092/46220092.pdf>. 5, 6, 7
- [3] Anja Becker, Antoine Joux, Alexander May, and Alexander Meurer. Decoding random binary linear codes in $2^{n/20}$: How $1 + 1 = 0$ improves information set decoding. In David Pointcheval and Thomas Johansson, editors, *EUROCRYPT 2012*, volume 7237 of *LNCS*, pages 520–536. Springer, Heidelberg, April 2012. <https://eprint.iacr.org/2012/026.pdf>. 19, 27, 28
- [4] Elwyn R Berlekamp, Robert J McEliece, and Henk CA van Tilborg. On the inherent intractability of certain coding problems. *IEEE Transactions on Information Theory*, 24(3):384–386, 1978. <http://authors.library.caltech.edu/5607/1/BERIEEEtit78.pdf>. 5
- [5] Daniel J Bernstein. Grover vs. mceliece. In *Post-Quantum Cryptography*, pages 73–80. Springer, 2010. <https://cr.yp.to/codes/grovercode-20091123.pdf>. 19, 28
- [6] Daniel J Bernstein, Tanja Lange, and Christiane Peters. Attacking and defending the mceliece cryptosystem. In *Post-Quantum Cryptography*, pages 31–46. Springer, 2008. <https://cr.yp.to/codes/mceliece-20080807.pdf>. 19, 28
- [7] Julien Bringer, Hervé Chabanne, Gérard Cohen, Bruno Kindarji, and Gilles Zémor. Optimal iris fuzzy sketches. In *Biometrics: Theory, Applications, and Systems, 2007. BTAS 2007. First IEEE International Conference on*, pages 1–6. IEEE, 2007. <https://arxiv.org/abs/0705.3740>. 13
- [8] Anne Canteaut and Florent Chabaud. A new algorithm for finding minimum weight words in a linear code: application to mceliece cryptosystem and to narrow-sense bch codes of length 511. *IEEE Transactions on Information Theory*, 44(1):367–378, 1998. <http://ieeexplore.ieee.org/document/651067/>. 19, 28

- [9] Rodolfo Canto Torres and Nicolas Sendrier. Analysis of information set decoding for a sub-linear error weight. In Tsuyoshi Takagi, editor, *Post-Quantum Cryptography - 7th International Workshop, PQCrypto 2016, Fukuoka, Japan, February 24-26, 2016, Proceedings*, volume 9606 of *Lecture Notes in Computer Science*, pages 144–161. Springer, 2016. <https://hal.inria.fr/hal-01244886>. 19, 27, 28
- [10] Robert Chien. Cyclic decoding procedures for bose-chaudhuri-hocquenghem codes. *IEEE Transactions on information theory*, 10(4):357–363, 1964. 16
- [11] Jean-Sébastien Coron, Helena Handschuh, Marc Joye, Pascal Paillier, David Pointcheval, and Christophe Tymen. Gem: A generic chosen-ciphertext secure encryption method. In *Cryptographers’ Track at the RSA Conference*, pages 263–276. Springer, 2002. http://www.di.ens.fr/~pointche/Documents/Papers/2002_rsa.pdf. 8
- [12] Jean-Christophe Deneuville, Philippe Gaborit, and Gilles Zémor. Ouroboros: A simple, secure and efficient key exchange protocol based on coding theory. In Tanja Lange and Tsuyoshi Takagi, editors, *Post-Quantum Cryptography - 8th International Workshop, PQCrypto 2017, Utrecht, The Netherlands, June 26-28, 2017, Proceedings*, volume 10346 of *Lecture Notes in Computer Science*, pages 18–34. Springer, 2017. http://www.unilim.fr/pages_perso/deneuville/files/ba43bf8d80cef2999dbf4308828213ec.pdf. 3
- [13] Ilya Dumer. On minimum distance decoding of linear codes. In *Proc. 5th Joint Soviet-Swedish Int. Workshop Inform. Theory*, pages 50–52, 1991. https://www.researchgate.net/publication/296573348_On_minimum_distance_decoding_of_linear_codes. 27
- [14] Matthieu Finiasz and Nicolas Sendrier. Security bounds for the design of code-based cryptosystems. In Mitsuru Matsui, editor, *ASIACRYPT 2009*, volume 5912 of *LNCs*, pages 88–105. Springer, Heidelberg, December 2009. <https://eprint.iacr.org/2009/414.pdf>. 19, 28
- [15] Eiichiro Fujisaki and Tatsuaki Okamoto. Secure integration of asymmetric and symmetric encryption schemes. In *Crypto*, volume 99, pages 537–554. Springer, 1999. https://link.springer.com/chapter/10.1007/3-540-48405-1_34. 8
- [16] Eiichiro Fujisaki and Tatsuaki Okamoto. Secure integration of asymmetric and symmetric encryption schemes. *Journal of cryptology*, pages 1–22, 2013. <https://link.springer.com/article/10.1007/s00145-011-9114-1>. 8
- [17] Philippe Gaborit. Shorter keys for code based cryptography. In *Proceedings of the 2005 International Workshop on Coding and Cryptography (WCC 2005)*, pages 81–91, 2005. http://www.unilim.fr/pages_perso/philippe.gaborit/shortIC.ps. 4

- [18] Philippe Gaborit and Marc Girault. Lightweight code-based identification and signature. In *2007 IEEE International Symposium on Information Theory*, pages 191–195. IEEE, 2007. https://www.unilim.fr/pages_perso/philippe.gaborit/isit_short_rev.pdf. 6
- [19] Shafi Goldwasser and Silvio Micali. Probabilistic encryption. *Journal of Computer and System Sciences*, 28(2):270–299, 1984. <http://groups.csail.mit.edu/cis/pubs/shafi/1984-jcss.pdf>. 7
- [20] Qian Guo, Thomas Johansson, and Carl Löndahl. A new algorithm for solving ring-lpn with a reducible polynomial. *IEEE Transactions on Information Theory*, 61(11):6204–6212, 2015. <https://arxiv.org/abs/1409.0472>. 27
- [21] Dennis Hofheinz, Kathrin Hövelmanns, and Eike Kiltz. A modular analysis of the fujisaki-okamoto transformation. Cryptology ePrint Archive, Report 2017/604, 2017. <http://eprint.iacr.org/2017/604>. 3, 8, 9, 10, 23, 28
- [22] W Cary Huffman and Vera Pless. *Fundamentals of error-correcting codes*. Cambridge university press, 2010. <https://www.amazon.fr/Fundamentals-Error-Correcting-Codes-Cary-Huffman/dp/0521131707>. 4
- [23] Laurie L Joiner and John J Komo. Decoding binary bch codes. In *Southeastcon’95. Visualize the Future., Proceedings., IEEE*, pages 67–73. IEEE, 1995. 16
- [24] Shu Lin and Daniel J Costello. *Error control coding*, volume 2. Prentice Hall Englewood Cliffs, 2004. 14, 15
- [25] Carl Löndahl, Thomas Johansson, Masoumeh Koochak Shooshtari, Mahmoud Ahmadian-Attari, and Mohammad Reza Aref. Squaring attacks on mceliece public-key cryptosystems using quasi-cyclic codes of even dimension. *Designs, Codes and Cryptography*, 80(2):359–377, 2016. <https://link.springer.com/article/10.1007/s10623-015-0099-x>. 27
- [26] Alexander May, Alexander Meurer, and Enrico Thomae. Decoding random linear codes in $\tilde{O}(2^{0.054n})$. In *Asiacrypt*, volume 7073, pages 107–124. Springer, 2011. https://link.springer.com/chapter/10.1007/978-3-642-25385-0_6. 27, 28
- [27] Alexander May and Ilya Ozerov. On computing nearest neighbors with applications to decoding of binary linear codes. In *EUROCRYPT (1)*, pages 203–228, 2015. <http://www.cits.rub.de/imperia/md/content/may/paper/codes.pdf>. 27, 28
- [28] Rafael Misoczki, Jean-Pierre Tillich, Nicolas Sendrier, and Paulo SLM Barreto. Mdp-mceliece: New mceliece variants from moderate density parity-check codes. In *Information Theory Proceedings (ISIT), 2013 IEEE International Symposium on*, pages 2069–2073. IEEE, 2013. <https://eprint.iacr.org/2012/409.pdf>. 5, 6

- [29] Tatsuaki Okamoto and David Pointcheval. React: Rapid enhanced-security asymmetric cryptosystem transform. *Topics in Cryptology—CT-RSA 2001*, pages 159–174, 2001. <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.150.5590&rep=rep1&type=pdf>. 8
- [30] Eugene Prange. The use of information sets in decoding cyclic codes. *IRE Transactions on Information Theory*, 8(5):5–9, 1962. <http://ieeexplore.ieee.org/document/1057777/>. 27
- [31] Nicolas Sendrier. Decoding one out of many. In *International Workshop on Post-Quantum Cryptography*, pages 51–67. Springer, 2011. <https://eprint.iacr.org/2011/367.pdf>. 6, 27, 28
- [32] Jacques Stern. A method for finding codewords of small weight. In *International Colloquium on Coding Theory and Applications*, pages 106–113. Springer, 1988. <https://link.springer.com/chapter/10.1007/BFb0019850>. 27