

## A Distinguisher for High Rate McEliece Cryptosystems

Jean-Charles Faugère, Valérie Gauthier-Umana, Ayoub Otmani, Ludovic Perret, Jean-Pierre Tillich

► **To cite this version:**

Jean-Charles Faugère, Valérie Gauthier-Umana, Ayoub Otmani, Ludovic Perret, Jean-Pierre Tillich. A Distinguisher for High Rate McEliece Cryptosystems. IEEE Transactions on Information Theory, Institute of Electrical and Electronics Engineers, 2013, 59 (10), pp.6830-6844. <<http://dx.doi.org/10.1109/TIT.2013.2272036>>. <10.1109/TIT.2013.2272036>. <hal-00776068>

**HAL Id: hal-00776068**

**<https://hal.inria.fr/hal-00776068>**

Submitted on 14 Jan 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A Distinguisher for High Rate McEliece Cryptosystems

Jean-Charles Faugère, Valérie Gauthier-Umaña, Ayoub Otmani, Ludovic Perret, Jean-Pierre Tillich

## Abstract

The Goppa Code Distinguishing (GD) problem consists in distinguishing the matrix of a Goppa code from a random matrix. The hardness of this problem is an assumption to prove the security of code-based cryptographic primitives such as McEliece’s cryptosystem. Up to now, it is widely believed that the GD problem is a hard decision problem. We present the first method allowing to distinguish alternant and Goppa codes over any field. Our technique can solve the GD problem in polynomial-time provided that the codes have sufficiently large rates. The key ingredient is an algebraic characterization of the key-recovery problem. The idea is to consider the rank of a linear system which is obtained by linearizing a particular polynomial system describing a key-recovery attack. Experimentally it appears that this dimension depends on the type of code. Explicit formulas derived from extensive experimentations for the rank are provided for “generic” random, alternant, and Goppa codes over any alphabet. Finally, we give theoretical explanations of these formulas in the case of random codes, alternant codes over any field of characteristic two and binary Goppa codes.

## Index Terms

McEliece cryptosystem, CFS signature, Algebraic cryptanalysis, Goppa Code Distinguishing problem.

## I. INTRODUCTION

**T**HIS paper investigates the difficulty of the Goppa Code Distinguishing (GD) problem which first appeared in [1]. This is a decision problem that aims at recognizing a generator matrix of a binary Goppa code from a randomly drawn binary matrix. Up to now, it is assumed that no polynomial time algorithm exists that distinguishes a generator matrix of a Goppa code from a randomly picked generator matrix.

The main motivation for introducing the GD problem is to formally relate the problem of decoding a random linear code to the security of the McEliece public-key cryptosystem [2]. Since its apparition, this cryptosystem has withstood many attacks and after more than thirty years now, it still belongs to the very few unbroken public key cryptosystems. This situation substantiates the claim that inverting the encryption function, and in particular recovering the private key from public data, is intractable. The classical methods for inverting the McEliece encryption function without finding a trapdoor all resort to the use of the best general decoding algorithms [3]–[10]. All these algorithms, whose time complexity is exponential (in the length), attempt to solve the long-standing problem of decoding random linear code [11]. They also assume (implicitly or explicitly) that there does not exist an algorithm that is able to decode more efficiently McEliece public keys. Note that if ever such an algorithm exists, it would permit to solve the GD problem.

On the other hand, no significant breakthrough has been observed with respect to the problem of recovering the private key [12], [13]. This has led to state that the generator matrix of a binary Goppa code does not disclose any visible structure that an attacker could exploit. This is strengthened by the fact that Goppa codes share many characteristics with random codes. For instance they asymptotically meet the Gilbert-Varshamov bound. They also have a trivial permutation group, *etc.* Hence, the hardness of the GD problem has become a classical belief, and as a consequence, a *de facto* assumption to prove the semantic security in the standard model (IND-CPA in [14] and IND-CCA2 in [15]), and the security in the random oracle model against existential forgery [1], [16] of the signature scheme [1].

We present a deterministic polynomial-time distinguisher for codes whose rate is close to 1. This includes in particular codes encountered with the signature scheme CFS [1], [17]. However, we emphasize that our method can distinguish codes also used in McEliece’s encryption scheme. For instance, the binary Goppa code obtained with  $m = 13$  and  $r = 19$  corresponding to a 90-bit security McEliece public key is distinguishable. More precisely, when the length of the code goes to infinity an asymptotic formula can be derived for the smallest rate  $R_{\text{crit}}$  for which we can distinguish a  $q$ -random code from a  $q$ -ary

J.-C. Faugère and L. Perret are with INRIA, Paris-Rocquencourt Center, POLSYS Project, UPMC Univ Paris 06, UMR 7606, LIP6, F-75005, Paris, France, CNRS, UMR 7606, LIP6, F-75005, Paris, France, e-mail: jean-charles.faugere@inria.fr, ludovic.perret@lip6.fr

Valérie Gauthier-Umaña is with GREYC - Université de Caen - Ensicaen, Boulevard Maréchal Juin, 14050 Caen Cedex, France, e-mail: valerie.gauthier01@unicaen.fr

A. Otmani is with LITIS - Université de Rouen, Technopôle du Madrillet, Avenue de l’Université, F-76801 Saint-Étienne-du-Rouvray Cedex e-mail: ayoub.otmani@univ-rouen.fr

J.-P. Tillich is with SECRET Project - INRIA Rocquencourt, Domaine de Voluceau, B.P. 105, 78153 Le Chesnay Cedex - France, e-mail: jean-pierre.tillich@inria.fr

alternant or Goppa code (Theorem 3). For a given  $q$  and assuming that the length is  $q^m$  then when  $m$  tends to infinity, we have:

$$R_{\text{crit}} = 1 - \sqrt{\frac{2m \log q}{q^m \log m}} (1 + o(1)).$$

where all logarithms are taken to base 2.

Our distinguisher is based on the algebraic attack developed against compact variants of McEliece [18]. In this approach, the key-recovery problem is transformed into the one of solving an algebraic system. By using a linearization technique, we are able to derive a linear system whose rank is different from what one would expect in the random case. More precisely, we observe experimentally that this *defect* in the rank is directly related to the type of codes. We provide explicit formulas for “generic” random, alternant, and Goppa codes over any alphabet. We performed extensive experiments to confirm that the formulas are accurate. Eventually, we prove the formula in the random case and give explanations in the case of alternant codes over any field of characteristic two and binary Goppa codes. However, the existence of our distinguisher does not undermine the security of primitives based on Goppa codes, but basically, it proves that the GD assumption is false for some parameters, and consequently should be used with great care as an assumption for a security reduction.

The paper is organized as follows. In Section III, we introduce the algebraic system that any McEliece cryptosystem must satisfy. In Section IV, we construct a linear system deduced from this algebraic system. This defines an algebraic distinguisher. We then provide explicit formulas that predicts the behavior of the distinguisher coming from experimentations. In Section VI, we give a proof of its typical behavior in the random case. In Section VII and Section VIII, we give explanations of the formulas for alternant and binary Goppa codes. Lastly, we conclude over the cryptographic implications the distinguisher induces and we deduce an asymptotic formula for the smallest rate for which we can distinguish a random code from an alternant code or a Goppa code.

## II. CODE-BASED PUBLIC-KEY CRYPTOGRAPHY

The general problem of decoding random linear codes is a potential candidate for building public-key cryptographic primitives such as an encryption scheme. McEliece in [2] was the first to use this problem in public-key cryptography. The general idea is to start from a family of codes equipped with a polynomial-time decoding algorithm. The fundamental concept of this proposal is to consider two equivalent representations of a code: one should facilitate the decoding, whereas from the other one, the decoding should be infeasible. Although his design principle is general, he explicitly advocated to use binary Goppa codes [19].

### A. Coding theory background

Code-based public-key cryptography focuses on linear codes that have a polynomial time decoding algorithm. We recall that a  $q$ -ary (linear) *code*  $\mathcal{C}$  over the finite field  $\mathbb{F}_q$  of  $q$  elements defined by a  $k \times n$  matrix  $\mathbf{G}$  (with  $k \leq n$ ) whose entries belong to  $\mathbb{F}_q$  is the vector space spanned by its rows *i.e.*,

$$\mathcal{C} \stackrel{\text{def}}{=} \left\{ \mathbf{u}\mathbf{G} \mid \mathbf{u} \in \mathbb{F}_q^k \right\}.$$

The length of  $\mathcal{C}$  is  $n$  and its rate is the ratio  $R \stackrel{\text{def}}{=} k/n$ . The role of decoding algorithms is to correct errors of prescribed weight. We say that a decoding algorithm corrects  $r$  errors if it recovers  $\mathbf{u}$  from the knowledge of  $\mathbf{u}\mathbf{G} + \mathbf{e}$  for all possible  $\mathbf{e} \in \mathbb{F}_q^n$  of weight at most  $r$ .

One famous family of codes is the one of binary Goppa codes. It belongs to the more general class of alternant codes ([20, Chap. 12, p. 365]). The main well-known feature of an alternant code is the possibility of being decoded in polynomial time. It is more convenient to describe this class through a parity-check matrix over an extension field  $\mathbb{F}_{q^m}$  of  $\mathbb{F}_q$  over which the code is defined. We recall that a parity-check matrix  $\mathbf{H}$  of a  $q$ -ary code  $\mathcal{C}$  is defined as a matrix such that:

$$\mathcal{C} \stackrel{\text{def}}{=} \left\{ \mathbf{c} \in \mathbb{F}_q^n \mid \mathbf{H}\mathbf{c}^T = 0 \right\}.$$

where the symbol  $T$  means the transpose operation. For  $q$ -ary alternant codes of length  $n \leq q^m$ , there exists a parity-check matrix with a very special form related to rectangular Vandermonde matrices:

$$\mathbf{V}_r(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \begin{pmatrix} y_1 & \cdots & y_n \\ y_1 x_1 & \cdots & y_n x_n \\ \vdots & & \vdots \\ y_1 x_1^{r-1} & \cdots & y_n x_n^{r-1} \end{pmatrix}$$

where  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_n)$  are in  $\mathbb{F}_{q^m}^n$ .

*Definition 1 (Alternant code):* A  $q$ -ary alternant code of order  $r$  associated to  $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{F}_{q^m}^n$  where all  $x_i$ 's are distinct and  $\mathbf{y} = (y_1, \dots, y_n) \in (\mathbb{F}_{q^m}^*)^n$  denoted by  $\mathcal{A}_r(\mathbf{x}, \mathbf{y})$  is  $\left\{ \mathbf{c} \in \mathbb{F}_q^n \mid \mathbf{V}_r(\mathbf{x}, \mathbf{y})\mathbf{c}^T = \mathbf{0} \right\}$ .

It is well-known that the dimension  $k$  of an alternant codes of degree  $r$  satisfies  $k \geq n - rm$ . Moreover, a key feature about them is the following property.

*Proposition 1:* An alternant codes of degree  $r$  can decode in polynomial time all errors of weight at most  $\frac{r}{2}$  whenever there exists a parity-check matrix in the form  $\mathbf{V}_r(\mathbf{x}^*, \mathbf{y}^*)$  for some vectors  $\mathbf{x}^*$  and  $\mathbf{y}^*$ .

*Definition 2 (Goppa codes):* A  $q$ -ary Goppa code  $\mathcal{G}(\mathbf{x}, \Gamma)$  associated to a polynomial  $\Gamma(z) \stackrel{\text{def}}{=} \sum_{i=0}^r \gamma_i z^i$  of degree  $r$  over  $\mathbb{F}_{q^m}$  and an  $n$ -tuple  $\mathbf{x} = (x_1, \dots, x_n)$  of distinct elements of  $\mathbb{F}_{q^m}$  satisfying  $\Gamma(x_i) \neq 0$  for all  $i, 1 \leq i \leq n$ , is the  $q$ -ary alternant code  $\mathcal{A}_r(\mathbf{x}, \mathbf{y})$  of order  $r$  with  $y_i = \Gamma(x_i)^{-1}$ .

Naturally Goppa codes, viewed as alternant codes, inherit a decoding algorithm that corrects up to  $\frac{r}{2}$  errors. But in the case of *binary* Goppa codes, it is possible to correct twice as many errors. The starting point is the following result given in [20, p. 341].

*Theorem 1:* A binary Goppa code  $\mathcal{G}(\mathbf{x}, \Gamma)$  associated to a Goppa polynomial  $\Gamma(z)$  of degree  $r$  without multiple roots is equal to the alternant code  $\mathcal{A}_{2r}(\mathbf{x}, \mathbf{y})$  with  $y_i = \Gamma(x_i)^{-2}$ .

*Corollary 1 ([21]):* There exists a polynomial time algorithm decoding all errors of weight at most  $r$  for any Goppa code  $\mathcal{G}(\mathbf{x}, \Gamma)$  where  $\Gamma(z)$  is of degree  $r$  and has no multiple roots.

It is worthwhile recalling that the only requirement for decoding a binary Goppa  $\mathcal{G}$  is either to know  $\mathbf{x}$  and  $\Gamma(z)$  or to know two vectors  $\mathbf{x}^*$  and  $\mathbf{y}^*$  such that:

$$\mathcal{G} = \mathcal{A}_{2r}(\mathbf{x}^*, \mathbf{y}^*). \quad (1)$$

## B. McEliece cryptosystem

We briefly recall here the general principle of McEliece's encryption scheme.

*Secret key:* the triplet  $(\mathbf{S}, \mathbf{G}_s, \mathbf{P})$  of matrices defined over a finite field  $\mathbb{F}_q$  over  $q$  elements, with  $q$  being a power of two, that is  $q = 2^s$ .  $\mathbf{G}_s$  is a full rank matrix of size  $k \times n$ , with  $k < n$ ,  $\mathbf{S}$  is of size  $k \times k$  and is invertible.  $\mathbf{P}$  is a permutation matrix of size  $n \times n$ . The generator matrix  $\mathbf{G}_s$  is chosen in such a way that its associated linear code has a decoding algorithm which corrects in polynomial time  $r$  errors.

*Public key:* the matrix  $\mathbf{G} = \mathbf{S}\mathbf{G}_s\mathbf{P}$ .

*Encryption:* A plaintext  $\mathbf{u} \in \mathbb{F}_q^k$  is encrypted by choosing a random vector  $\mathbf{e}$  in  $\mathbb{F}_q^n$  of weight at most  $r$ . The corresponding ciphertext is  $\mathbf{c} = \mathbf{u}\mathbf{G} + \mathbf{e}$ .

*Decryption:*  $\mathbf{c}' = \mathbf{c}\mathbf{P}^{-1}$  is computed from the ciphertext  $\mathbf{c}$ . Notice that  $\mathbf{c}' = (\mathbf{u}\mathbf{S}\mathbf{G}_s\mathbf{P} + \mathbf{e})\mathbf{P}^{-1} = \mathbf{u}\mathbf{S}\mathbf{G}_s + \mathbf{e}\mathbf{P}^{-1}$  and that  $\mathbf{e}\mathbf{P}^{-1}$  is of Hamming weight at most  $r$ . Therefore the aforementioned decoding algorithm can recover in polynomial time  $\mathbf{u}\mathbf{S}$  and therefore the plaintext  $\mathbf{u}$  by multiplication by  $\mathbf{S}^{-1}$ .

## C. CFS signatures

Another important code-based cryptographic primitive is the CFS signature scheme [1]. A user whose public key is  $\mathbf{G}$  and who wishes to sign a message  $\mathbf{x} \in \mathbb{F}_2^k$  has to compute a string  $\mathbf{u}$  such that the Hamming weight of  $\mathbf{x} - \mathbf{u}\mathbf{G}$  is at most  $r$ . Anyone (a *verifier*) can publicly check the validity of a signature. Unfortunately, this approach can only provide signatures for messages  $\mathbf{x}$  that are within distance  $r$  from a codeword  $\mathbf{u}\mathbf{G}$ . The CFS scheme prompts to modify the message by appending a counter incremented until the decoding algorithm can find such a signature. The efficiency of this scheme heavily depends on the number of trials. With a binary Goppa codes of length  $n = 2^m$  and dimension  $k = n - mr$ , the number of trials is of order  $r!$ . So one has to choose a very small  $r$  and therefore to take a very large  $n$  in order to be secure. The code rate is then equal to  $\frac{2^m - rm}{2^m} = 1 - \frac{mr}{2^m}$  which is quite close to 1 for large  $n$  (that is for large values of  $2^m$ ) and moderate values of  $r$ . For instance, a 80-bit security CFS scheme requires to take  $n = 2^{21}$  and  $r = 10$  whereas the McEliece cryptosystem for the same security needs to choose  $n = 2^{11}$  and  $r = 32$  ([17]). Thus one major difference between the McEliece cryptosystem and the CFS scheme lies in the choice of the parameters.

## D. Goppa Code Distinguishing Problem

The minimum requirement for an encryption function is that it should be infeasible from a given ciphertext  $\mathbf{c}$  and public data<sup>1</sup> like the public key  $pk$ , ciphertexts, *etc.* to recover the corresponding plaintext  $\mathbf{x}$ . This issue is directly linked to the following computational problem.

<sup>1</sup>This kind of attack is called a *Chosen Plaintext Attack* (CPA).

*Definition 3 (McEliece Problem):* Let  $\mathbf{G}$  be a generator matrix of a binary Goppa code of length  $n \leq 2^m$  and dimension  $k = n - rm$  where  $m$  and  $t$  are positive integers. Let  $\mathbf{x}$  be a vector from  $\mathbb{F}_{2^m}^k$  and let  $\mathbf{e}$  be a vector from  $\mathbb{F}_{2^m}^n$  of weight  $t$ . Finally, we set  $\mathbf{c} \stackrel{\text{def}}{=} \mathbf{x}\mathbf{G} + \mathbf{e}$ . Then the *McEliece Problem* asks to find  $\mathbf{x}$  and  $\mathbf{e}$  only from  $\mathbf{G}$  and  $\mathbf{c}$ .

One obvious way of solving this problem consists in devising a method that recovers the private key. But, it is also possible to recover a plaintext from a specific ciphertext without resorting to a key-recovery attack. In particular, an attacker against the McEliece scheme would find the plaintext by applying general decoding methods like [2]–[10], [22]–[26] on the public matrix  $\mathbf{G}$ . Such attacks are called *decoding attacks*.

The only known methods that aim to solve the McEliece problem are based either on an exhaustive search of the private key or on applying very general decoding methods. Both approaches run in exponential time on some of the parameters. But this situation is still unsatisfactory because there is no certitude that there does not exist a better way to solve it.

A classical stance is to claim that binary Goppa codes look like random linear codes. It amounts to say that there does not exist a polynomial-time computable quantity which behaves differently depending on whether the code is a Goppa or a random code. Currently, it is an open problem to establish a formal proof that would substantiate the claim that a binary Goppa code is *indistinguishable* from a random code.

This assumption is attractive because it enables to rely on the hardness of decoding a random linear code to prove the security of the McEliece function. This reasoning does make sense because binary Goppa codes share several common aspects<sup>2</sup> with a randomly picked linear code. Furthermore, all the general decoding algorithms do not exploit the information, even partially, that a matrix describes a ‘hidden’ Goppa code. Based on this, the authors of [1] defined the *Goppa code distinguishing problem*. Before formalizing this problem, we introduce some notation. For any integers  $n$  and  $k$  such that  $k \leq n$ . We denote by  $\text{Goppa}(n, k)$  the set of  $k \times n$  generator matrices of binary Goppa codes. Similarly,  $\text{Random}(n, k)$  is the set of binary  $k \times n$  random generator matrices.

*Definition 4 (Goppa Code Distinguishing (GD) Problem):* A *distinguisher*  $\mathcal{D}$  is an algorithm that takes as input a matrix  $\mathbf{G}$  and returns a bit.  $\mathcal{D}$  solves the GD problem if it wins the following game:

- $b \leftarrow \{0, 1\}$
- If  $b = 0$  then  $\mathbf{G} \leftarrow \text{Goppa}(n, k)$  else  $\mathbf{G} \leftarrow \text{Random}(n, k)$
- If  $\mathcal{D}(\mathbf{G}) = b$  then  $\mathcal{D}$  wins else  $\mathcal{D}$  loses.

*Definition 5:* The advantage  $\text{Adv}^{GD}(\mathcal{D})$  of a GD distinguisher  $\mathcal{D}$  is defined by  $\text{Adv}^{GD}(\mathcal{D}) \stackrel{\text{def}}{=} \left| \Pr[\mathcal{D}(\mathbf{G}) = 1 : \mathbf{G} \leftarrow \text{Goppa}(n, k)] - \Pr[\mathcal{D}(\mathbf{G}) = 1 : \mathbf{G} \leftarrow \text{Random}(n, k)] \right|$ ,

where  $\Pr[\mathcal{D}(\mathbf{G}) = 1 : \mathbf{G} \leftarrow \text{Goppa}(n, k)]$  is the probability that  $\mathcal{D}$  outputs 1 when  $\mathbf{G}$  is a random binary generator matrix of a Goppa code, and  $\Pr[\mathcal{D}(\mathbf{G}) = 1 : \mathbf{G} \leftarrow \text{Random}(n, k)]$  is defined similarly for a binary random matrix.

*Definition 6:* A function  $\varepsilon(k)$  is *negligible* if for any integer  $a > 0$ , there exists an integer  $k_a > 0$  such that:

$$\forall k \geq k_a, \quad \varepsilon(k) < \frac{1}{k^a}.$$

The interest of negligible function is to keep a probability of an event negligible even after polynomially many tries. We are now able to state an important assumption<sup>3</sup>.

*Assumption 1 ([1]):*  $\text{Adv}^{GD}(\mathcal{D})$  is negligible for any polynomial-time algorithm  $\mathcal{D}$  that solves the GD problem.

Until our recent work in [27] and this paper, the only known algorithm that solves the GD problem enumerates binary Goppa codes and tests the code equivalence thanks to the *Support Splitting* algorithm [28]. This approach runs in time  $\mathcal{O}\left(\frac{n^{r-1}}{mr}\right)$  for binary Goppa codes of degree  $r$  and length  $n$  with  $m \leq \log_2 n$  and  $r = \frac{1}{m}(1 - R)n$  where  $R$  is the code rate.

### E. Semantically Secure Conversions

The fundamental issue when dealing with cryptographic primitives is to prove its security. Several approaches are possible. The most natural one is to show that the primitive resists to the best known attacks. However, this does not guarantee that there will not appear one day a better attack that renders the primitive insecure. The methodology of *security proof by reduction* appeared to remedy this question by linking a security notion that a cryptographic primitive should verify to an algorithmic problem widely considered as hard. The approach is similar to the one that proves the NP-Completeness of a given problem. Such a ‘‘security proof’’ proves that if an attacker exists then it can be used as a subroutine to solve a hard problem. In other words, such an attacker has little chances to exist.

<sup>2</sup>Similarly to random codes, Goppa codes asymptotically meet the Gilbert-Varshamov bound. They have also a trivial permutation group like random codes.

<sup>3</sup>According to [1], proving or disproving the hardness of the GD problem will have a significant impact: ‘‘Classification issues are in the core of coding theory since its emergence in the 50’s. So far nothing significant is known about Goppa codes, more precisely there is no known property invariant by permutation and computable in polynomial time which characterizes Goppa codes. Finding such a property or proving that none exists would be an important breakthrough in coding theory and would also probably seal the fate, for good or ill, of Goppa code-based cryptosystems’’.

These simple facts prompt to design conversions that would lead to an IND-CCA secure encryption scheme. The first article to propose such a conversion for the McEliece cryptosystem is [29] which proposes a conversion resulting into an IND-CCA2 in the *Random Oracle Model* under the assumption that the problem of decoding random linear codes is difficult. This work was then followed by [30] which proposes another modification while providing an IND-CPA secure encryption scheme in the standard model<sup>4</sup> under the assumptions that both decoding random linear codes *and* distinguishing Goppa codes are difficult problems. Finally, under the same assumptions, [31] proposed (a modified) McEliece cryptosystem that is IND-CCA2 in the *standard model*.

### III. ALGEBRAIC CRYPTANALYSIS OF MCELIECE-LIKE CRYPTOSYSTEMS

The McEliece cryptosystem relies on binary Goppa codes which belong to the class of *alternant codes*. We are now able to construct an algebraic system as explained in [18] for a key-recovery. This algebraic system will be the main ingredient for building a distinguisher. We assume that the public matrix is a  $k \times n$  generator matrix  $\mathbf{G}$  where by assumption  $k = n - rm$ . We know that the knowledge of a matrix  $\mathbf{V}_r(\mathbf{x}^*, \mathbf{y}^*)$  for some vectors  $\mathbf{x}^*$  and  $\mathbf{y}^*$  allows to efficiently decode the public code defined by  $\mathbf{G}$ . Furthermore, from the definition of  $\mathbf{G}$ , we also know that:

$$\mathbf{V}_r(\mathbf{x}^*, \mathbf{y}^*)\mathbf{G}^T = \mathbf{0}.$$

Let  $X_1, \dots, X_n$  and  $Y_1, \dots, Y_n$  be  $2n$  variables corresponding to the  $x_i^*$ 's and the  $y_i^*$ 's. Observe that such  $x_i^*$ 's and  $y_i^*$ 's are a particular solution [18] of the following polynomial system:

$$\left\{ g_{i,1}Y_1X_1^j + \dots + g_{i,n}Y_nX_n^j = 0 \mid 1 \leq i \leq k \text{ and } 0 \leq j \leq r-1 \right\} \quad (2)$$

where the  $g_{i,j}$ 's are the entries of the known matrix  $\mathbf{G}$ .

Clearly, solving this system would lead to a possibly equivalent private key. For compact variants [32], [33] of [2], additional structures permit to drastically reduce the number of variables allowing to solve (2) for a large set of parameters in polynomial-time using dedicated Gröbner bases techniques [18]. But the general case is currently a major *open question*. However, we describe a simple way for *partially* solving (2). It basically consists in deriving a linear system from the polynomial system (2). Note that this operation is actually the first step performed during the computation of Gröbner bases algorithms such as by F4 or F5 [34], [35]. From now on, we will always assume that  $q = 2^s$  with  $s \geq 1$ . We can assume that  $\mathbf{G} = (g_{ij})$  with  $1 \leq i \leq k$  and  $1 \leq j \leq n$  is in reduced row echelon form over its  $k$  first positions:

$$\mathbf{G} = (\mathbf{I}_k \mid \mathbf{P})$$

where  $\mathbf{P} = (p_{ij})$  for  $1 \leq i \leq k$  and  $k+1 \leq j \leq n$  is the submatrix of  $\mathbf{G}$  formed by its last  $n - k = mr$  columns. Next, for any  $i \in \{1, \dots, k\}$  and  $e \in \{0, \dots, r-1\}$ , we can rewrite (2) as

$$Y_iX_i^e = \sum_{j=k+1}^n p_{i,j}Y_jX_j^e. \quad (3)$$

In particular, it follows that for all  $i$  in  $\{1, \dots, k\}$ :

$$\begin{cases} Y_i &= \sum_{j=k+1}^n p_{i,j}Y_j \\ Y_iX_i &= \sum_{j=k+1}^n p_{i,j}Y_jX_j \\ Y_iX_i^2 &= \sum_{j=k+1}^n p_{i,j}Y_jX_j^2. \end{cases}$$

Then, thanks to the trivial identity  $Y_i(Y_iX_i^2) = (Y_iX_i)^2$  for all  $i$  in  $\{1, \dots, k\}$ , we get:

$$\sum_{j=k+1}^n p_{i,j}Y_j \sum_{j=k+1}^n p_{i,j}Y_jX_j^2 = \left( \sum_{j=k+1}^n p_{i,j}Y_jX_j \right)^2.$$

It is possible to reorder this to obtain:

$$\sum_{j=k+1}^{n-1} \sum_{j'=j+1}^n p_{i,j}p_{i,j'} (Y_jY_{j'}X_{j'}^2 + Y_{j'}Y_jX_j^2) = 0.$$

<sup>4</sup>There is no hash function in this model

We thus obtain a linear system  $\mathcal{L}_{\mathbf{P}}$  of  $k$  equations involving  $\binom{mr}{2}$  variables  $Z_{jj'} \stackrel{\text{def}}{=} Y_j Y_{j'} X_j^2 + Y_{j'} Y_j X_{j'}^2$

$$\mathcal{L}_{\mathbf{P}} \stackrel{\text{def}}{=} \begin{cases} \sum_{j=k+1}^{n-1} \sum_{j'>j}^n p_{1,j} p_{i,j'} Z_{jj'} = 0 \\ \vdots \\ \sum_{j=k+1}^{n-1} \sum_{j'>j}^n p_{k,j} p_{i,j'} Z_{jj'} = 0 \end{cases} \quad (4)$$

*Definition 7:* For any integer  $r \geq 1$  and  $m \geq 1$ , the number of variables  $\binom{mr}{2}$  in the linear system  $\mathcal{L}_{\mathbf{P}}$  as defined in (4) is denoted by  $N$  and its rank by  $\text{rank}(\mathcal{L}_{\mathbf{P}})$ . We denote by  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$  the kernel of  $\mathcal{L}_{\mathbf{P}}$  and its dimension as a  $\mathbb{F}_q$ -vector space is denoted by  $D$ .

Let us recall that  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$  is necessarily a  $\mathbb{F}_q$ -vector space since the linear system (4) have coefficients in  $\mathbb{F}_q$  but the solutions of (2) are sought in the extension field  $\mathbb{F}_{q^m}$ . Furthermore, we obviously have:

$$D = N - \text{rank}(\mathcal{L}_{\mathbf{P}}).$$

Hence, in order to recover the solutions of (2), it is necessary that  $\text{rank}(\mathcal{L}_{\mathbf{P}})$  is almost equal to the number of variables  $N = \binom{mr}{2}$ . For a random system, this is likely to happen when the number  $k$  of equations in (4) is greater than the number of unknowns, that is to say:

$$k \geq N.$$

It appears experimentally that  $D$  is amazingly large even in the case where  $k \geq N$ . It even depends on whether or not the code with generator matrix  $\mathbf{G}$  is chosen as a (generic) alternant code or as a Goppa code. Interestingly enough, when  $\mathbf{G}$  is chosen at random,  $\text{rank}(\mathcal{L}_{\mathbf{P}})$  is equal to  $\min\{k, N\}$  with very high probability. In particular, the dimension of the solution space is typically 0 when  $k$  is larger than the number of variables  $N$  as one would expect. This will be proved in Section VI. Although this *defect* in the rank is an obstacle to break the McEliece cryptosystem, it can be used to distinguish the public generator of a *structured* code from a random code.

#### IV. A DISTINGUISHER OF ALTERNANT AND GOPPA CODES

We consider three cases: when the  $p_{ij}$ 's are chosen uniformly and independently at random in  $\mathbb{F}_q$  then we denote by  $D_{\text{random}}$  the dimension of  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$ . When  $\mathbf{G}$  is chosen as a generator matrix of a random alternant (*resp.* Goppa) code of degree  $r$ , we denote it by  $D_{\text{alternant}}$  (*resp.*  $D_{\text{Goppa}}$ ). We carried out intensive computations with Magma [36] by randomly generating alternant and Goppa codes over the field  $\mathbb{F}_q$  with  $q \in \{2, 4, 8, 16, 32\}$  for  $r$  in the range  $\{3, \dots, 50\}$  and several values of  $m$ . Furthermore, in our probabilistic model, a random alternant code is obtained by picking uniformly and independently at random two vectors  $(x_1, \dots, x_n)$  and  $(y_1, \dots, y_n)$  from  $(\mathbb{F}_{q^m})^n$  such that the  $x_i$ 's are all different and the  $y_i$ 's are all nonzero. A random Goppa code is obtained by taking a random vector  $(x_1, \dots, x_n)$  in  $(\mathbb{F}_{q^m})^n$  with all the  $x_i$ 's different and a random *irreducible* polynomial  $\Gamma(z) = \sum_i \gamma_i z^i$  of degree  $r$ . Our experiments have revealed that the dimension of  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$  is *predictable* and follows formulas.

*Experimental Fact 1 (Alternant Case):* As long as

$$N - D_{\text{alternant}} < k,$$

then with high probability,  $D_{\text{alternant}}$  is equal to  $T_{\text{alternant}}$  where:

$$T_{\text{alternant}} \stackrel{\text{def}}{=} \frac{1}{2} m(r-1) \left( (2e+1)r - 2 \frac{q^{e+1} - 1}{q-1} \right) \quad (5)$$

where  $e \stackrel{\text{def}}{=} \lceil \log_q(r-1) \rceil$ .

*Experimental Fact 2 (Goppa Case):* As long as

$$N - D_{\text{Goppa}} < k,$$

then with high probability,  $D_{\text{Goppa}}$  is equal to  $T_{\text{Goppa}}$  where:

$$T_{\text{Goppa}} \stackrel{\text{def}}{=} \begin{cases} \frac{1}{2} m(r-1)(r-2) = T_{\text{alternant}} & \text{for } r < q-1, \\ \frac{1}{2} mr \left( (2e+1)r - 2q^e + 2q^{e-1} - 1 \right) & \text{for } r \geq q-1, \end{cases} \quad (6)$$

with  $e$  being the unique integer such that:

$$(q-1)^2 q^{e-2} < r \leq (q-1)^2 q^{e-1}.$$

## V. EXPERIMENTAL RESULTS

We gathered in Table II-XI in the appendix some samples of the results obtained through intensive computations with the Magma system [36]. We randomly generated alternant and Goppa codes over the field  $\mathbb{F}_q$  with  $q \in \{2, 4, 8, 16, 32\}$  for values of  $r$  in the range  $\{3, \dots, 50\}$  and several  $m$ . The Goppa codes are generated by means of an irreducible  $\Gamma(z)$  of degree  $r$  and hence  $\Gamma(z)$  has no multiple roots. In particular, we can apply Theorem 1 in the binary case. We compare the dimensions of the solution space against the dimension  $D_{\text{random}}$  of the system derived from a random linear code. Table II and Table III give figures for the binary case with  $m = 14$ . We can check that  $D_{\text{random}}$  is equal to 0 for  $r \in \{3, \dots, 12\}$  and  $D_{\text{random}} = N - k$  as expected. We remark that  $D_{\text{alternant}}$  is different from  $D_{\text{random}}$  whenever  $r \leq 15$ , and  $D_{\text{Goppa}}$  is different from  $D_{\text{random}}$  as long as  $r \leq 25$ . Finally we observe that our formulas for  $T_{\text{alternant}}$  fit as long as  $k \geq N - D_{\text{alternant}}$  which correspond to  $r \leq 15$ . This is also the case for binary Goppa codes since we have  $T_{\text{Goppa}} = D_{\text{Goppa}}$  as long as  $k \geq N - D_{\text{Goppa}}$  i.e.,  $r \leq 25$ . We also give in Table X and Table XI the examples we obtained for  $q = 4$  and  $m = 6$  to check that the arguments also apply. We also compare binary Goppa codes and random linear codes for  $m = 15$  in Table IV-VI and  $m = 16$  in Table VII-IX. We see that  $D_{\text{random}}$  and  $D_{\text{Goppa}}$  are different for  $r \leq 33$  when  $m = 15$  and for  $m = 16$  they are different even beyond our range of experiment ( $r \leq 50$ ).

## VI. RANDOM CASE

The purpose of this section is to study the behavior of  $D_{\text{random}}$ , namely the dimension of  $\text{Ker}(\mathcal{L}_P)$  as  $\mathbb{F}_q$ -vector space when the entries of the matrix  $P$  are drawn independently from the uniform distribution over  $\mathbb{F}_q$ . In this case, we can show that:

*Theorem 2:* Assume that  $N \leq k$  and that the entries of  $P$  are drawn independently from the uniform distribution over  $\mathbb{F}_q$ . Then for any function  $\omega(x)$  tending to infinity as  $x$  goes to infinity, we have

$$\text{prob}\left(D_{\text{random}} \geq mr\omega(mr)\right) = o(1),$$

as  $mr$  goes to infinity.

Notice that if we choose  $\omega(x) = \log(x)$  for instance, then asymptotically the dimension  $D_{\text{random}}$  of the solution space is with very large probability smaller than  $mr \log(mr)$ . When  $m$  and  $r$  are of the same order (which is generally chosen in practice) this quantity is smaller than  $D_{\text{alternant}}$  or  $D_{\text{Goppa}}$  which are of the form  $\Omega(mr^2)$ .

The main ingredient for proving Theorem 2 consists in analyzing a certain (partial) Gaussian elimination process on the matrix:

$$M \stackrel{\text{def}}{=} \left( p_{ij} p_{ij'} \right)_{\substack{1 \leq i \leq k \\ k+1 \leq j < j' \leq n}}.$$

We can see the matrix  $M$  in block form, each block consists of the matrix:

$$B_j = \left( p_{i,k+j} p_{i,k+j'} \right)_{\substack{1 \leq i \leq k \\ 1 \leq j < j' \leq n-k}}.$$

Each block  $B_j$  is of size  $k \times (rm - j)$ . Notice that in  $B_j$ , the rows for which  $p_{i,k+j} = 0$  consist only of zeros. To start the Gaussian elimination process with  $B_1$ , we will therefore choose  $rm - 1$  rows for which  $p_{i,k+1} \neq 0$ . This gives a square matrix  $M_1$ . We perform Gaussian elimination on  $M$  by adding rows involved in  $M_1$  to put the first block  $B_1$  in standard form. We continue this process with  $B_2$  by picking now  $rm - 2$  rows which have not been chosen before and which correspond to  $p_{i,k+2} \neq 0$ . This yields a square submatrix  $M_2$  of size  $rm - 2$  and we continue this process until we reach the last block. The key observation is that:

$$\text{rank}(M) \geq \text{rank}(M_1) + \dots + \text{rank}(M_{rm-1}).$$

A rough analysis of this process yields Theorem 2. The important point is that what happens for different blocks are independent processes and it corresponds to looking at different rows of the matrix  $P$ . We give all the previous results that we need in order to prove Theorem 2.

It will be convenient to assume that the columns of  $M$  are ordered lexicographically. The index of the first column is  $(j, j') = (k + 1, k + 2)$ , the second one is  $(j, j') = (k + 1, k + 3)$ , while the last one is  $(j, j') = (n - 1, n)$ . The matrices  $M_i$ 's which are involved in the Gaussian elimination process mentioned above are defined inductively as follows. Let  $E_1$  be the subset of  $\{1, \dots, k\}$  of indices  $s$  such that  $p_{s,k+1} \neq 0$ . Let  $F_1$  be the subset of  $E_1$  formed by its first  $rm - 1$  elements (if these elements exist). Now, we set

$$M_1 \stackrel{\text{def}}{=} \left( p_{s,k+1} p_{s,j} \right)_{\substack{s \in F_1 \\ k+1 < j \leq n}}. \quad (7)$$

Let  $r_1$  be the rank of  $M_1$ . To simplify the discussion, we assume that:

- 1)  $F_1 = \{1, 2, \dots, rm - 1\}$ ,
- 2) the submatrix  $N_1$  of  $M_1$  formed by its first  $r_1$  rows and columns is of full rank.



Note that we can always assume this by performing suitable row and column permutations. In other words  $M$  has the following block structure:

$$M = \begin{pmatrix} N_1 & B_1 \\ A_1 & C_1 \end{pmatrix}.$$

We denote:

$$M^{(1)} \stackrel{\text{def}}{=} \begin{pmatrix} N_1^{-1} & O \\ -A_1 N_1^{-1} & I \end{pmatrix} M,$$

where  $O$  is a matrix of size  $r_1 \times (k - r_1)$  with only zero entries and  $I$  is the identity matrix of size  $k - r_1$ . Notice that  $M^{(1)}$  takes the block form:

$$M^{(1)} = \begin{pmatrix} I & B'_1 \\ O & C'_1 \end{pmatrix}.$$

This is basically performing Gaussian elimination on  $M$  in order to have the first  $r_1$  columns in standard form. We then define inductively the  $E_i, F_i, M_i, M^{(i)}$  and  $N_i$  as follows:

$$E_i \stackrel{\text{def}}{=} \left\{ s \mid 1 \leq s \leq k, p_{s,k+i} \neq 0 \right\} \setminus \bigcup_{u=1}^{i-1} F_{i-u},$$

$$F_i \stackrel{\text{def}}{=} \text{the first } rm - i \text{ elements of } E_i.$$

$M_i$  is the submatrix of  $M^{(i-1)}$  obtained from the rows in  $F_i$  and the columns associated to the indices of the form  $(k + i, j')$  where  $j'$  ranges from  $k + i + 1$  to  $n$ .  $M^{(i)}$  is obtained from  $M^{(i-1)}$  by first choosing a square submatrix  $N_i$  of  $M_i$  of full rank and with the same rank as  $M_i$  and then by performing Gaussian elimination on the rows in order to put the columns of  $M^{(i-1)}$  involved in  $N_i$  in standard form (i.e., the submatrix of  $M^{(i-1)}$  corresponding to  $N_i$  becomes the identity matrix while the other entries in the columns involved in  $N_i$  become zero). It is clear that the whole process leading to  $M^{(rm-1)}$  amounts to perform (partial) Gaussian elimination to  $M$ . Hence:

*Lemma 1:* When  $|E_i| \geq rm - i$ , for all  $i \in \{1, \dots, rm - 1\}$ , we have:

$$\text{rank}(M) \geq \sum_{i=1}^{rm-1} \text{rank}(M_i).$$

Another observation is that  $M_i$  is equal to the sum of the submatrix  $(p_{s,k+i} p_{s,j})_{\substack{s \in F_i \\ k+i < j \leq n}}$  of  $M$  and a certain matrix which is some function on the entries  $p_{t,k+i} p_{t,j}$  where  $t$  belongs to  $F_1 \cup \dots \cup F_{i-1}$  and  $j$  ranges over  $\{k + i + 1, n\}$ . Since by definition of  $F_i$ ,  $p_{s,k+i}$  is different from 0 for  $s$  in  $F_i$ . In addition, the rank of  $M_i$  does not change by multiplying each row of index  $s$  by  $p_{s,k+i}^{-1}$ . Then, it turns out that the rank of  $M_i$  is equal to the rank of a matrix which is the sum of the matrix  $(p_{s,j})_{\substack{s \in F_i \\ k+i < j \leq n}}$ , another matrix depending on the  $p_{t,k+i} p_{t,j}$ 's (where  $t$  ranges over  $F_1 \cup \dots \cup F_{i-1}$ ) and the  $p_{s,k+i}$ 's with  $s \in F_i$ . This proves that:

*Lemma 2:* Assume that  $|E_i| \geq rm - i$  for all  $i \in \{1, \dots, rm - 1\}$ . Then, the random variables  $\text{rank}(M_i)$  are independent and  $\text{rank}(M_i)$  is distributed as the rank of a square matrix of size  $rm - i$  with entries drawn independently from the uniform distribution on  $\mathbb{F}_q$ .

Another essential ingredient for proving Theorem 2 is the following well known lemma (see for instance [37][Theorem 1])

*Lemma 3:* There exist two positive constants  $A$  and  $B$  depending on  $q$  such that the probability  $p(s, \ell)$  that a random  $\ell \times \ell$  matrix over  $\mathbb{F}_q$  is of rank  $\ell - s$  (where the coefficients are drawn independently from each other from the uniform distribution on  $\mathbb{F}_q$ ) satisfies

$$\frac{A}{q^{s^2}} \leq p(s, \ell) \leq \frac{B}{q^{s^2}}.$$

This enables to control the exponential moments of the defect of a random matrix. For a square matrix  $M$  of size  $\ell \times \ell$ , we define the defect  $d(M)$  by  $d(M) \stackrel{\text{def}}{=} \ell - \text{rank}(M)$ .

*Lemma 4:* If  $M$  is random square matrix whose entries are drawn independently from the uniform distribution over  $\mathbb{F}_q$ , then there exists some constant  $K$  such that for every  $\lambda > 0$ ,

$$\mathbb{E} \left( q^{\lambda d(M)} \right) \leq K q^{\frac{\lambda^2}{4}},$$

$\mathbb{E}(\cdot)$  denoting the expectation.

*Proof:* By using Lemma 3, we obtain:

$$\mathbb{E} \left( q^{\lambda d(M)} \right) \leq \sum_{d=0}^{\infty} q^{\lambda d} \frac{B}{q^{d^2}} \leq B \sum_{d=0}^{\infty} q^{\lambda d - d^2}.$$

Observe that the maximum of the function  $d \mapsto q^{\lambda d - d^2}$  is reached for  $d_0 = \frac{\lambda}{2}$  and is equal to  $q^{\frac{\lambda^2}{4}}$ . Then, we can write the sum above as:

$$\sum_{d=0}^{\infty} q^{\lambda d - d^2} = \sum_{d \leq d_0} q^{\lambda d - d^2} + \sum_{d > d_0} q^{\lambda d - d^2}$$

Finally, we notice that:

$$\begin{aligned} \frac{q^{\lambda(d+1) - (d+1)^2}}{q^{\lambda d - d^2}} &\leq \frac{q^{\lambda(d_0+1) - (d_0+1)^2}}{q^{\lambda d_0 - d_0^2}} = \frac{1}{q} \text{ for } d > d_0, \\ \frac{q^{\lambda(d-1) - (d-1)^2}}{q^{\lambda d - d^2}} &\leq \frac{q^{\lambda(d_0-1) - (d_0-1)^2}}{q^{\lambda d_0 - d_0^2}} = \frac{1}{q} \text{ for } d \leq d_0. \end{aligned}$$

This leads to:

$$\begin{aligned} \sum_{d=0}^{\infty} q^{\lambda d - d^2} &\leq \sum_{d \leq d_0} q^{d - \lfloor d_0 \rfloor} q^{\frac{\lambda^2}{4}} + \sum_{d > d_0} q^{\lceil d_0 \rceil - d} q^{\frac{\lambda^2}{4}} \\ &= O\left(q^{\frac{\lambda^2}{4}}\right). \end{aligned}$$

We can use now the previous lemma together with Lemma 1 and Lemma 2 to derive the following lemma. ■

*Lemma 5:* Assuming that  $|E_i| \geq rm - i$  for all  $i \in \{1, \dots, t\}$ , we get:

$$\mathbf{prob}\left(\sum_{i=1}^t d(M_i) \geq u\right) \leq K^t q^{-\frac{u^2}{t}}$$

where  $K$  is the constant appearing in Lemma 4.

*Proof:* Let  $D \stackrel{\text{def}}{=} \sum_{i=1}^t d(\mathbf{M}_i)$ . Using Markov's inequality:

$$\mathbf{prob}(D \geq u) \leq \frac{\mathbb{E}(q^{\lambda D})}{q^{\lambda u}} \tag{8}$$

for some well chosen  $\lambda > 0$ . The exponential moment appearing at the numerator is upper-bounded with the help of the previous lemma and by using the independence of the random variables  $q^{\lambda d(\mathbf{M}_i)}$ , i.e.,:

$$\begin{aligned} \mathbb{E}(q^{\lambda D}) &= \mathbb{E}\left(q^{\lambda \sum_{i=1}^t d(\mathbf{M}_i)}\right) \\ &= \prod_{i=1}^t \mathbb{E}\left(q^{\lambda d(\mathbf{M}_i)}\right) \\ &\leq K^t q^{\frac{t\lambda^2}{4}}. \end{aligned} \tag{9}$$

Using now (9) in (8), we obtain

$$\mathbf{prob}(D \geq \alpha t) \leq K^t \frac{q^{\frac{t\lambda^2}{4}}}{q^{\lambda \alpha t}} = K^t q^{\frac{t\lambda^2}{4} - \lambda \alpha t}.$$

We choose  $\lambda = \frac{2u}{t}$  to minimize this upper-bound, leading to:

$$\mathbf{prob}(D \geq u) \leq K^t q^{-\frac{u^2}{t}}. \tag{10}$$

The last ingredient for proving Theorem 2 is a bound on the probability that  $E_i$  is too small to construct  $F_i$ . ■

*Lemma 6:* Let  $u_i \stackrel{\text{def}}{=} \binom{mr}{2} - \frac{(2rm-i)(i-1)}{2}$  and  $F$  be the event " $|F_j| = rm - j$  for  $j \in \{1, \dots, i-1\}$ " then

$$\mathbf{prob}(|E_i| < rm - i \mid F) \leq e^{-2 \frac{\left(\frac{q-1}{q} u_i - rm - i + 1\right)^2}{u_i}}.$$

*Proof:* When all the sets  $F_j$  are of size  $rm - j$  for  $j$  in  $\{1, \dots, i-1\}$ , it remains

$$N - \sum_{j=1}^{i-1} (rm - j) = N - \frac{(2rm - i)(i - 1)}{2} = u_i$$

rows which can be picked up for  $E_i$ . Let  $S_t$  be the sum of  $t$  Bernoulli variables of parameter  $\frac{q-1}{q}$ . We obviously have

$$\mathbf{prob}(|E_i| < rm - i \mid F) = \mathbf{prob}(S_{u_i} < rm - i).$$

It remains to use the Hoeffding inequality on the binomial tails to finish the proof.  $\blacksquare$

We are ready now to prove Theorem 2:

*Proof of Theorem 2:* Let  $u = \lceil \sqrt{mr\omega(mr)} \rceil$ . We observe now that if all  $E_j$ 's are of size at least  $rm - j$  for  $j \in \{1, \dots, u\}$ , we can write

$$\begin{aligned}
D &= N - \text{rank}(M) \\
&\leq N - \sum_{i=1}^{rm-u} \text{rank}(M_i) \text{ (by Lemma 1)} \\
&= \sum_{i=1}^{rm-1} (rm - i) - \sum_{i=1}^{rm-u} \text{rank}(M_i) \\
&= \sum_{i=1}^{rm-u} d(M_i) + \sum_{i=rm-u+1}^{rm-1} (rm - i) \\
&= \sum_{i=1}^{rm-u} d(M_i) + \frac{u(u-1)}{2} \\
&< \sum_{i=1}^{rm-u} d(M_i) + \frac{mr\omega(mr)}{2}.
\end{aligned}$$

From this we deduce that

$$\begin{aligned}
\mathbf{prob}(D_{\text{random}} \geq mr\omega(mr)) &\leq \mathbf{prob}(A \cup B) \\
&\leq \mathbf{prob}(A) + \mathbf{prob}(B)
\end{aligned}$$

where  $A$  is the event “ $\sum_{i=1}^{rm-u} d(M_i) \geq \frac{mr\omega(mr)}{2}$ ” and  $B$  is the event “for at least one  $E_j$  with  $j \in \{1, \dots, rm - u\}$  we have  $|E_j| < rm - j$ ”. We use now Lemma 5 to prove that  $\mathbf{prob}(A) = o(1)$  as  $rm$  goes to infinity. We finish the proof by noticing that the probability of the complementary set of  $B$  satisfies

$$\begin{aligned}
\mathbf{prob}(\bar{B}) &= \mathbf{prob}\left(\bigcap_{i=1}^{rm-u} |E_i| \geq rm - i\right) \\
&= \prod_{i=1}^{rm-u} \mathbf{prob}(|E_i| \geq rm - i \mid F) \\
&= 1 - o(1) \text{ (by Lemma 6)}.
\end{aligned}$$

$\blacksquare$

## VII. ALTERNANT CASE

The goal of this section is to explain the value of the dimension  $D_{\text{alternant}}$  of  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$  for  $q$ -ary alternant codes of degree  $r$ . We shall see that this dimension will be obtained by first identifying a  $\mathbb{F}_{q^m}$ -basis of  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$  when viewed as a linear system with coefficients in  $\mathbb{F}_{q^m}$ . To set up the linear system  $\mathcal{L}_{\mathbf{P}}$  as defined in (4), we have used the trivial identity  $Y_i Y_i X_i^2 = (Y_i X_i)^2$ . The fundamental remark is that we can use any identity  $Y_i X_i^a Y_i X_i^b = Y_i X_i^c Y_i X_i^d$  with  $a, b, c, d \in \{0, 1, \dots, r-1\}$  such that  $a + b = c + d$ . Such identities lead to the same algebraic system  $\mathcal{L}_{\mathbf{P}}$ :

$$\sum_{(j,j') \in J} p_{i,j} p_{i,j'} \left( Y_j X_j^a Y_{j'} X_{j'}^b + Y_{j'} X_{j'}^a Y_j X_j^b + Y_j X_j^c Y_{j'} X_{j'}^d + Y_{j'} X_{j'}^c Y_j X_j^d \right) = 0 \quad (10)$$

where by definition

$$J \stackrel{\text{def}}{=} \left\{ (j, j') \in \mathbb{N} \times \mathbb{N} \mid k+1 \leq j < j' \leq n \right\}. \quad (11)$$

Consequently, the fact that *there are many different ways of combining the equations of the algebraic system together yielding the same linearized system*  $\mathcal{L}_{\mathbf{P}}$  explains why the dimension  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$  is large. Indeed, if a non-trivial solution of  $X_1, \dots, X_n, Y_1, \dots, Y_n$  exists then all the expressions of the form

$$Y_j X_j^a Y_{j'} X_{j'}^b + Y_{j'} X_{j'}^a Y_j X_j^b + Y_j X_j^c Y_{j'} X_{j'}^d + Y_{j'} X_{j'}^c Y_j X_j^d$$

for all  $(j, j') \in J$ , and for all  $a, b, c, d \in \{0, 1, \dots, r-1\}$  such that  $a + b = c + d$  lead to a solution of the linear system  $\mathcal{L}_{\mathbf{P}}$ .

In what follows, we exhibit further elements of  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$ . To this end, we use the automorphisms  $x \mapsto x^{q^\ell}$  where  $\ell$  is in  $\{0, \dots, m-1\}$ . Indeed, we can also consider the identity:

$$(Y_i X_i^a)^{q^\ell} (Y_i X_i^b)^{q^\ell} = (Y_i X_i^c)^{q^\ell} (Y_i X_i^d)^{q^\ell} \quad (12)$$

for any integers  $a, b, c, d, \ell$  and  $\ell'$  such that:

$$aq^{\ell'} + bq^\ell = cq^{\ell'} + dq^\ell.$$

We get again the linear system  $\mathcal{L}_P$ . However, assuming that  $\ell' \leq \ell$ , solutions obtained from such equations are exactly those coming from the identity:

$$Y_i X_i^a Y_i^{q^{\ell-\ell'}} X_i^{bq^{\ell-\ell'}} = Y_i X_i^c Y_i^{q^{\ell-\ell'}} X_i^{dq^{\ell-\ell'}}. \quad (13)$$

We can focus on vectors that satisfy equations obtained with  $0 \leq a, b, c, d < r, 0 \leq \ell < m$  and  $a + q^\ell b = c + q^\ell d$  i.e., the equations obtained from:

$$Y_i X_i^a (Y_i X_i^b)^{q^\ell} = Y_i X_i^c (Y_i X_i^d)^{q^\ell}.$$

We now try to determine the number of linearly independent solutions induced by such identities.

*Definition 8:* Let  $a, b, c$  and  $d$  be integers in  $\{0, \dots, r-1\}$  and an integer  $\ell$  in  $\{0, \dots, m-1\}$  such that  $a + q^\ell b = c + q^\ell d$ . We define  $\mathbf{U}^{a,b,c,d,\ell} = \left( \mathbf{U}_{j,j'}^{a,b,c,d,\ell} \right)_{(j,j') \in J}$  with

$$\mathbf{U}_{j,j'}^{a,b,c,d,\ell} \stackrel{\text{def}}{=} Y_j X_j^a Y_{j'}^{q^\ell} X_{j'}^{q^\ell b} + Y_j X_j^c Y_{j'}^{q^\ell} X_{j'}^{q^\ell d}.$$

In the same way, we set  $\bar{\mathbf{U}}^{a,b,c,d,\ell} = \left( \bar{\mathbf{U}}_{j,j'}^{a,b,c,d,\ell} \right)_{(j,j') \in J}$  where

$$\bar{\mathbf{U}}_{j,j'}^{a,b,c,d,\ell} \stackrel{\text{def}}{=} Y_{j'} X_{j'}^a Y_j^{q^\ell} X_j^{q^\ell b} + Y_{j'} X_{j'}^c Y_j^{q^\ell} X_j^{q^\ell d}.$$

Finally, we define  $\mathbf{Z}^{a,b,c,d,\ell} = \left( \mathbf{Z}_{j,j'}^{a,b,c,d,\ell} \right)_{(j,j') \in J}$  such that:

$$\mathbf{Z}_{j,j'}^{a,b,c,d,\ell} \stackrel{\text{def}}{=} \mathbf{U}_{j,j'}^{a,b,c,d,\ell} + \bar{\mathbf{U}}_{j,j'}^{a,b,c,d,\ell}.$$

The linear system  $\mathcal{L}_P$  rewrites then as:

$$\sum_{k+1 \leq j < j' \leq n} p_{i,j} p_{i,j'} \mathbf{Z}_{j,j'}^{a,b,c,d,\ell} = 0.$$

By construction, there exists relations occurring between some  $\mathbf{Z}^{a,b,c,d,\ell}$ 's. For instance we always have  $\mathbf{Z}^{a,b,a,b,\ell} = 0$  for any integers  $a, b$  and  $\ell$ . We also have other basic relations:

*Lemma 7:* Let integers  $a, b, c, d, e, f$  be in  $\{0, \dots, r-1\}$ , and an integer  $\ell$  be in  $\{0, \dots, m-1\}$  such that  $a + q^\ell b = c + q^\ell d$ . We have:

$$\mathbf{Z}^{a,b,c,d,\ell} + \mathbf{Z}^{c,d,e,f,\ell} = \mathbf{Z}^{a,b,e,f,\ell}. \quad (14)$$

Without loss of generality, we can assume that  $d > b$  and let us set  $\delta = d - b$ . Moreover, as we have  $a + q^\ell b = c + q^\ell d$ , it implies that  $a = c + q^\ell \delta$ . Thus, any vector  $\mathbf{Z}^{a,b,c,d,\ell}$  is uniquely described by the tuple  $(b, c, \delta, \ell)$  by setting  $d = b + \delta$  and  $a = c + q^\ell \delta$  provided that  $1 \leq \delta \leq r-1-b$  and  $0 \leq c + q^\ell \delta \leq r-1$ .

The next proposition shows that some vectors  $\mathbf{Z}^{c+q^\ell \delta, b, c, b+\delta, \ell}$  can be expressed as a linear combination of vectors defined with  $\delta = 1$ .

*Proposition 2:* Let  $\ell, \delta, b$  and  $c$  be integers such that  $\ell \geq 0, \delta \geq 1, 1 \leq b + \delta \leq r-1$  and  $1 \leq c + q^\ell \delta \leq r-1$ . For all  $1 \leq i \leq \delta$ , we set  $b_i \stackrel{\text{def}}{=} b + i - 1$  and  $c_i \stackrel{\text{def}}{=} c + q^\ell (\delta - i)$ . We have

$$\mathbf{Z}^{c+q^\ell \delta, b, c, b+\delta, \ell} = \sum_{i=1}^{\delta} \mathbf{Z}^{c_i+q^\ell, b_i, c_i, b_i+1, \ell}. \quad (15)$$

*Proof:* When  $\delta = 1$ , the equality (15) is obviously verified. We now assume that  $\delta \geq 2$ . Let  $b^* \stackrel{\text{def}}{=} b + 1, \delta^* \stackrel{\text{def}}{=} \delta - 1$  and  $c^* \stackrel{\text{def}}{=} c + q^\ell \delta^*$ . Then  $c^*$  is the integer such that  $c^* + q^\ell = c + q^\ell \delta$ , one can see that  $c + q^\ell \delta^* = c + q^\ell (\delta - 1) = c^*$ . By Lemma 7 we have:

$$\begin{aligned} \mathbf{Z}^{c^*+q^\ell, b, c^*, b+1, \ell} + \mathbf{Z}^{c+q^\ell \delta^*, b^*, c, b^*+\delta^*, \ell} &= \mathbf{Z}^{c^*+q^\ell, b, c, b^*+\delta^*, \ell} \\ &= \mathbf{Z}^{c+q^\ell \delta, b, c, b+\delta, \ell}. \end{aligned}$$

This means that

$$\mathbf{Z}^{c+q^\ell \delta, b, c, b+\delta, \ell} = \mathbf{Z}^{c^*+q^\ell, b, c^*, b+1, \ell} + \mathbf{Z}^{c+q^\ell \delta^*, b^*, c, b^*+\delta^*, \ell}.$$

The proof follows by induction. ■

*Remark 1:* Note that if we have  $t \ 0 \leq a, b, c, d < r, 0 \leq \ell < m$  be such that  $a + q^\ell b = c + q^\ell d$  then we also have  $0 \leq \ell \leq \lfloor \log_q(r-1) \rfloor$ .

From Proposition 2, we can deduce that the set of vectors  $\mathbf{Z}^{c+q^\ell \delta, b, c, b+\delta, \ell}$  obtained with tuples  $(\delta, b, c, \ell)$  such that  $\delta = 1$  form a spanning set. Actually, we can characterize more precisely this set.

*Definition 9:* Let  $\mathcal{B}_r$  be the set of *nonzero* vectors  $\{\mathbf{Z}^{c+q^\ell, b, c, b+1, \ell} \mid b, c, \ell\}$ . We then have:

$$\mathcal{B}_r = \left\{ \mathbf{Z}^{c+1, b, c, b+1, 0} \mid 0 \leq b < c \leq r-2 \right\} \cup \left\{ \mathbf{Z}^{c+q^\ell, b, c, b+1, \ell} \mid 0 \leq b \leq r-2, 1 \leq \ell \leq \lfloor \log_q(r-1) \rfloor, 0 \leq c \leq r-1-q^\ell \right\}$$

We are now in position to conclude for the alternant case.

*Proposition 3:* Let  $r$  be an integer such that  $r \geq 3$  and let us denote by  $|\mathcal{B}_r|$  the cardinality of  $\mathcal{B}_r$ . Then:

$$T_{\text{alternant}} = m|\mathcal{B}_r|.$$

*Proof:* Let us set  $e \stackrel{\text{def}}{=} \lfloor \log_q(r-1) \rfloor$ . Then the number of elements in  $\mathcal{B}_r$  is given by the number of tuples  $(b, c, \ell)$ . Therefore we get:

$$\begin{aligned} |\mathcal{B}_r| &= \frac{1}{2}(r-1)(r-2) + \sum_{\ell=1}^e \sum_{b=0}^{r-2} (r-q^\ell) \\ &= \frac{1}{2}(r-1) \left( r-2 + 2er - 2 \sum_{\ell=1}^e q^\ell \right) \\ &= \frac{1}{2}(r-1) \left( (2e+1)r - 2 \sum_{\ell=0}^e q^\ell \right) \\ &= \frac{1}{m} T_{\text{alternant}}. \end{aligned}$$

Proposition 3 gives an explanation of the value of  $D_{\text{alternant}}$ . Indeed, it shows that  $\mathcal{B}_r$  is a  $\mathbb{F}_{q^m}$ -basis that provides a  $\mathbb{F}_q$ -basis of  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$  by the following heuristic. ■

*Heuristic 1:* Consider a certain decomposition of the elements of  $\mathbb{F}_{q^m}$  in a  $\mathbb{F}_q$  basis. Let  $\pi_i : \mathbb{F}_{q^m} \rightarrow \mathbb{F}_q$  be the function giving the  $i$ -th coordinate in this decomposition with  $1 \leq i \leq m$ . By extension we denote for  $\mathbf{z} = (z_j)_{1 \leq j \leq n} \in (\mathbb{F}_{q^m})^n$  by  $\pi_i(\mathbf{z})$  the vector  $(\pi_i(z_j))_{1 \leq j \leq n} \in \mathbb{F}_q^n$ . Then, for any  $j$  such that  $1 \leq j \leq n$  and for random choices of  $x_j$ 's and  $y_j$ 's, the set  $\{\pi_i(\mathbf{Z}) \mid 1 \leq i \leq m \text{ and } \mathbf{Z} \in \mathcal{B}_r\}$  forms a basis of  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$ .

In the next section, we will investigate binary Goppa codes. Note that for  $q$ -ary Goppa codes of degree  $r < q-1$  we observed that  $T_{\text{alternant}} = T_{\text{Goppa}}$ . In this case, it is easy to see that (5) simplifies to:

$$\frac{1}{2}m(r-1)(r-2) \stackrel{\text{def}}{=} T_{\text{Goppa}}.$$

This is due to the fact that  $e = 0$  when  $r < q-1$ . We leave as an open question the proof that  $q$ -ary Goppa codes of degree  $r < q-1$  behave for our distinguisher as alternant codes. We focus now on the classical case – in code-based cryptography – of binary Goppa codes.

## VIII. BINARY GOPPA CASE

The goal of this section is to identify a basis of  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$  for binary Goppa codes of degree  $r$ . We assume therefore that  $q = 2$ . In that special case, the theoretical expression  $T_{\text{Goppa}}$  (Experimental Fact 2) has a simpler expression.

*Proposition 4:* Let  $e = \lceil \log_2 r \rceil + 1$ . When  $q = 2$ , Equation (6) can be simplified to

$$T_{\text{Goppa}} = \frac{1}{2}mr \left( (2e+1)r - 2^e - 1 \right).$$

Theorem 1 shows that a binary Goppa code of degree  $r$  can be regarded as a binary alternant code of degree  $2r$ . This seems to indicate that we should have

$$D_{\text{Goppa}}(r) = T_{\text{alternant}}(2r).$$

This is not the case though. It turns out that  $D_{\text{Goppa}}(r)$  is significantly smaller than this. In our experiments, we have found out that the vectors of  $\mathcal{B}_{2r}$  still form a generating set for  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$ . Unfortunately, they are not independent anymore. Our goal is therefore to identify the additional dependencies occurring in  $\mathcal{B}_{2r}$ . We will see that many of them come from  $\mathbb{F}_{2^m}$ -relations induced by the Goppa polynomial  $\Gamma(z)$ . Recall that by definition  $Y_i = \Gamma(X_i)^{-2}$ . This fact will allow to derive two types of linear dependencies. The first type of linear relations is rather natural, whilst the second type is more subtle.

### A. Goppa Polynomials: First Linear Dependencies

We will derive the first linear dependencies by means of the Goppa polynomial  $\Gamma(X)$ .

*Proposition 5:* Let  $t, \ell$  and  $c$  be integers such that  $0 \leq t \leq r-2$ ,  $1 \leq \ell \leq \lfloor \log_2(2r-1) \rfloor$  and  $0 \leq c \leq 2r-2^\ell-1$ . We also set  $c^* \stackrel{\text{def}}{=} c+2^{\ell-1}$ . It holds that :

$$\sum_{b=0}^r \gamma_b^{2^\ell} \mathbf{Z}^{c+2^\ell, t+b, c, t+b+1, \ell} = \mathbf{Z}^{c^*+2^{\ell-1}, 2t, c^*, 2t+1, \ell-1} + \mathbf{Z}^{c+2^{\ell-1}, 2t+1, c, 2t+2, \ell-1}. \quad (16)$$

*Proof:* Let  $\ell, \delta, b$  and  $c$  be integers such that  $\ell \geq 0, \delta \geq 1, 1 \leq b+\delta \leq r-1, 1 \leq c+2^\ell \delta \leq r-1$ . Let also  $d = b+\delta$  and  $a = c+2^\ell \delta$ . We can write that for any  $(j, j')$  in  $J$  (notation being as in (11)):

$$\begin{aligned} \mathbf{U}_{j, j'}^{a, b, c, d, \ell} &= \mathbf{U}_{j, j'}^{c+2^\ell \delta, b, c, b+\delta, \ell} = Y_j Y_{j'}^{2^\ell} \left( X_j^a X_{j'}^{2^\ell b} + X_j^c X_{j'}^{2^\ell d} \right) \\ &= Y_j Y_{j'}^{2^\ell} X_{j'}^{2^\ell b} \left( X_j^a + X_j^c X_{j'}^{2^\ell \delta} \right) = Y_j Y_{j'}^{2^\ell} X_{j'}^{2^\ell b} X_j^c \left( X_j^{2^\ell \delta} + X_{j'}^{2^\ell \delta} \right) \\ &= \left( X_j^\delta + X_{j'}^\delta \right)^{2^\ell} Y_j X_j^c \left( Y_{j'} X_{j'}^b \right)^{2^\ell}. \end{aligned}$$

In the same way, we can prove that for all  $(j, j')$  in  $J$ :

$$\overline{\mathbf{U}}_{j, j'}^{c+2^\ell \delta, b, c, b+\delta, \ell} = \left( X_j^\delta + X_{j'}^\delta \right)^{2^\ell} Y_{j'} X_{j'}^c \left( Y_j X_j^b \right)^{2^\ell}.$$

Observe now that:

$$Y_j X_j^c \left( Y_{j'} X_{j'}^{t+b} \right)^{2^\ell} = Y_j X_j^c Y_{j'}^{2^{\ell-1}} X_{j'}^{2^\ell t} \left( Y_{j'} X_{j'}^{2b} \right)^{2^{\ell-1}}.$$

Let  $\gamma_b$  be the coefficient of  $z^b$  in  $\Gamma(z)$ , that is  $\Gamma(z) = \sum_{b=0}^r \gamma_b z^b$ . Since  $Y_j \Gamma(X_{j'})^2 = 1$ , we have  $Y_{j'} \sum_{b=0}^r \gamma_b^2 X_{j'}^{2b} = 1$ . This

implies that  $\sum_{b=0}^r \gamma_b^2 Y_{j'}^2 X_{j'}^{2(b+t)} = Y_{j'} X_{j'}^{2t}$  and then:

$$\sum_{b=0}^r \gamma_b^{2^\ell} Y_j X_j^c \left( Y_{j'} X_{j'}^{t+b} \right)^{2^\ell} = Y_j X_j^c \left( Y_{j'} X_{j'}^{2t} \right)^{2^{\ell-1}}.$$

Therefore, we get:

$$\sum_{b=0}^r \gamma_b^{2^\ell} \mathbf{U}^{c+2^\ell \delta, t+b, c, t+b+\delta, \ell} = \mathbf{U}^{c'+2^{\ell'} \delta', b', c', b'+\delta', \ell'} \quad (17)$$

with  $\ell' = \ell-1$ ,  $\delta' = 2\delta$ ,  $b' = 2t$ , and  $c' = c$ . Since  $c'+2^{\ell'} \delta' = c+2^\ell \delta$  and  $c+2^\ell \delta \leq 2r-1$  we have  $c'+2^{\ell'} \delta' \leq 2r-1$ . Moreover, we require  $b'+\delta' \leq 2r-1$  which means  $2(t+\delta) \leq 2r-1$ . This last inequality implies  $t+\delta \leq r-1$ .

Similarly, we can prove:

$$\sum_{b=0}^r \gamma_b^{2^\ell} \overline{\mathbf{U}}^{c+2^\ell \delta, t+b, c, t+b+\delta, \ell} = \overline{\mathbf{U}}^{c'+2^{\ell'} \delta', b', c', b'+\delta', \ell'}.$$

We are now in position to conclude the proof. Thanks to the previous results:

$$\sum_{b=0}^r \gamma_b^{2^\ell} \mathbf{Z}^{c+2^\ell, t+b, c, t+b+1, \ell} = \mathbf{Z}^{c+2^\ell, 2t, c, 2(t+1), \ell-1}$$

Moreover by Proposition 2, we also have:

$$\mathbf{Z}^{c+2^\ell, 2t, c, 2(t+1), \ell-1} = \mathbf{Z}^{c^*+2^{\ell-1}, 2t, c^*, 2t+1, \ell-1} + \mathbf{Z}^{c+2^{\ell-1}, 2t+1, c, 2t+2, \ell-1}$$

where by definition  $c^*$  is equal to  $c+2^{\ell-1}$ . ■

As a consequence,  $\mathcal{B}_{2r}$  can not be a basis of  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$ . We count the number of linear dependencies predicted by Proposition 5.

*Proposition 6:* Let  $N_L$  be the number of equations of the form (16) and let us set  $u \stackrel{\text{def}}{=} \lfloor \log_2(2r-1) \rfloor$ . Then, the following equality holds:

$$N_L = 2(r-1)(ru+1-2^u).$$

*Proof:* Each equation is defined by a triple  $(t, c, \ell)$ . As  $0 \leq t \leq r-2$ ,  $1 \leq \ell \leq u$  and  $0 \leq c \leq 2r-2^\ell-1$ . Thus:

$$N_L = \sum_{t=0}^{r-2} \sum_{\ell=1}^u (2r-2^\ell).$$
■

### B. Goppa Polynomials: Additional Linear Dependencies

It turns out that there are still other dependencies. To see this, we define the vector  $\mathbf{Q}^{a,b,c,d,\ell} \stackrel{\text{def}}{=} \left( \left( \mathbf{Z}_{j,j'}^{a,b,c,d,\ell} \right)^2 \right)_{(j,j') \in J}$ . It immediately follows that:

$$\mathbf{Q}_{j,j'}^{a,b,c,d,\ell} = \left( \mathbf{U}_{j,j'}^{a,b,c,d,\ell} \right)^2 + \left( \overline{\mathbf{U}}_{j,j'}^{a,b,c,d,\ell} \right)^2.$$

We exhibit new linear dependencies between some elements of  $\mathcal{B}_{2r}$  and the vectors  $\mathbf{Q}_{j,j'}^{a,b,c,d,\ell}$ .

*Proposition 7:* For any integers  $b \geq 0, t \geq 0, \delta \geq 1$  and  $\ell$  such that  $0 \leq \ell \leq \lfloor \log_2(2r-1) \rfloor - 1, b + \delta \leq 2r - 1$  and  $t + 2^\ell \delta \leq r - 1$ , we have

$$\mathbf{Z}^{2t+2^{\ell+1}\delta, b, 2t, b+\delta, \ell+1} = \sum_{c=0}^r \gamma_c^2 \mathbf{Q}^{c+2^\ell\delta, b, t+c, b+\delta, \ell}. \quad (18)$$

*Proof:* We recall that  $Y_j \Gamma(X_j)^2 = 1$  implies that  $\sum_{c=0}^r \gamma_b^2 Y_j^2 X_j^{2(c+t)} = Y_j X_j^{2t}$ . Thus, for any  $(j, j')$  in  $J$ , we have:

$$\begin{aligned} W &\stackrel{\text{def}}{=} \sum_{c=0}^r \gamma_c^2 \left( \mathbf{U}_{j,j'}^{c+2^\ell\delta, b, t+c, b+\delta, \ell} \right)^2 \\ &= (X_j^\delta + X_{j'}^\delta)^{2^{\ell+1}} (Y_{j'} X_{j'}^b)^{2^{\ell+1}} X_j^{2t} Y_j^2 \sum_{c=0}^r \gamma_c^2 X_j^{2c} = (X_j^\delta + X_{j'}^\delta)^{2^{\ell+1}} (Y_{j'} X_{j'}^b)^{2^{\ell+1}} Y_j X_j^{2t} \\ &= \mathbf{U}_{j,j'}^{c'+2^{\ell'}\delta, b', c', b'+\delta', \ell'} \end{aligned}$$

with  $\ell' = \ell + 1, \delta' = \delta, b' = b, c' = 2t$  and  $c' + 2^{\ell'}\delta' = 2t + 2^{\ell+1}\delta$ . In particular, one can easily check that the necessary conditions are  $b + \delta \leq 2r - 1$  and  $t + 2^\ell \delta \leq r - 1$  in order for this equation to hold. We finish the proof by doing the same procedure for

$$\sum_{c=0}^r \gamma_c^2 \left( \overline{\mathbf{U}}_{j,j'}^{c+2^\ell\delta, b, t+c, b+\delta, \ell} \right)^2 = \overline{\mathbf{U}}_{j,j'}^{c'+2^{\ell'}\delta, b', c', b'+\delta', \ell'}.$$

We can count the number of linearly dependencies predicted by Proposition 7. ■

*Proposition 8:* Let  $N_Q$  be the number of vectors of  $\mathcal{B}_{2r}$  satisfying Equation (18) and let us set  $u \stackrel{\text{def}}{=} \lfloor \log_2(2r-1) \rfloor$ . Then we have that

$$N_Q = (2r-1)(ru - 2^u + 1).$$

*Proof:* By Proposition 7 we know that  $N_Q$  is the number of vectors  $\mathbf{Z}^{2t+2^{\ell+1}\delta, b, 2t, b+\delta, \ell+1}$  obtained with  $\delta = 1, b \geq 0, t \geq 0$  and satisfying  $0 \leq \ell \leq u - 1, b + \delta \leq 2r - 1$  and  $t + 2^\ell \delta \leq r - 1$ . Therefore:

$$N_Q = \sum_{l=0}^{u-1} \sum_{t=0}^{r-1-2^\ell} (2r-1). \quad (19)$$

We now want to count the number of linear dependencies induced by Proposition 7 and Proposition 5. The difficulty is that some of the  $N_Q$  vectors of  $\mathcal{B}_{2r}$  are counted twice because they appear both in linear relations of the form (16) and ‘‘quadratic’’ equations of the form (18). Let  $N_{L \cap Q}$  be the number of such vectors. More precisely, let  $\mathcal{B}_{2r}^{\text{quad}}$  be the subset of vectors of  $\mathcal{B}_{2r}$  which are involved in an Equation of type (18). There are equations of type (16) which involve only vectors of  $\mathcal{B}_{2r}^{\text{quad}}$ . Let  $N_1$  be their numbers. Moreover, it is possible by adding two equations of type (16) involving at least one vector which is not in  $\mathcal{B}_{2r}^{\text{quad}}$  to obtain an equation which involves only vectors of  $\mathcal{B}_{2r}^{\text{quad}}$ . Let  $N_0$  be the number of such sums. Finally, let  $N_{L \cap Q} \stackrel{\text{def}}{=} N_1 + N_0$ . It is possible to count such equations. ■

*Proposition 9:*  $N_{L \cap Q} = (r-1) \left( \left( u - \frac{1}{2} \right) r - 2^u + 2 \right)$  where  $u \stackrel{\text{def}}{=} \lfloor \log_2(2r-1) \rfloor$ .

*Proof:* We will consider vectors  $\mathbf{Z}^{c+2^\ell, b, c, b+1, \ell}$  of  $\mathcal{B}_{2r}$  that satisfy Equation (18) and such that there exists a linear relation that link them. In other words, we consider all the linear relations of the form

$$\sum_i \alpha_i \mathbf{Z}^{c_i+2^{\ell_i}, b_i, c_i, b_i+1, \ell_i} = 0$$

with  $\alpha_i$  in  $\mathbb{F}_{2^m}$  and where each  $\mathbf{Z}^{c_i+2^{\ell_i}, b_i, c_i, b_i+1, \ell_i}$  is equal to a linear relation of the form (18). We will see that the number of *independent* equations is equal to  $N_{L \cap Q}$ . First, one can observe that for any such vectors we necessary have  $c_i$  even and

$1 \leq \ell_i \leq u$ . We also know by Proposition 5 that for any integers  $t, \ell$  and  $c$  such that  $0 \leq t \leq r-2$ ,  $1 \leq \ell \leq u$  and  $0 \leq c \leq 2r-2^\ell-1$ , we have the following linear relation:

$$\sum_{b=0}^r \gamma_b^{2^\ell} \mathbf{Z}^{c+2^\ell, t+b, c, t+b+1, \ell} = \mathbf{Z}^{c^*+2^{\ell-1}, 2t, c^*, 2t+1, \ell-1} + \mathbf{Z}^{c+2^{\ell-1}, 2t+1, c, 2t+2, \ell-1}$$

where by definition  $c^* = c + 2^{\ell-1}$ . Note in particular that whenever  $c$  is even then  $c^*$  is also even and if  $\ell \geq 2$  then we obtain a linear relation between some vectors that also satisfy quadratic equations of the form (18). Each equation enables to remove one quadratic equation. So if we denote by  $N_1$  the number of equations of the form of (16) with  $c$  even and  $\ell \geq 2$ , we then have:

$$\begin{aligned} N_1 &= \sum_{t=0}^{r-2} \sum_{\ell=2}^u \left( \frac{1}{2} (2r - 2^\ell) \right) \\ &= (r-1) \sum_{\ell=1}^{u-1} (r - 2^\ell) \\ &= (r-1) \left( (u-1)r - 2^u + 2 \right). \end{aligned} \quad (20)$$

Moreover in the case  $\ell = 1$ , Equation (17) becomes

$$\sum_{b=0}^r \gamma_b^2 \mathbf{Z}^{c+2, t+b, c, t+b+1, 1} = \mathbf{Z}^{c+2, 2t, c, 2t+2, 0}.$$

In particular, when  $c$  is even, say for instance  $c = 2t'$  for some integer, then this last equation can be rewritten as:

$$\sum_{b=0}^r \gamma_b^2 \mathbf{Z}^{2t'+2, t+b, 2t', t+b+1, 1} = \mathbf{Z}^{2t'+2, 2t, 2t', 2t+2, 0}. \quad (21)$$

We know that when  $t' = t$  then  $\mathbf{Z}^{2t'+2, 2t, 2t', 2t+2, 0}$  is zero. In that case we obtain new relations between vectors satisfying quadratic equations that are independent even from those obtained with  $\ell \geq 2$ . As for the case when  $t \neq t'$  we also have

$$\mathbf{Z}^{2t'+2, 2t, 2t', 2t+2, 0} = \mathbf{Z}^{2t+2, 2t', 2t, 2t'+2, 0}.$$

From this identity and from Equation (21) we then obtain new relations of the following form:

$$\sum_{b=0}^r \gamma_b^2 \mathbf{Z}^{2t'+2, t+b, 2t', t+b+1, 1} = \sum_{b=0}^r \gamma_b^2 \mathbf{Z}^{2t+2, t'+b, 2t, t'+b+1, 1}. \quad (22)$$

This last equation involves only vectors that satisfy also quadratic equations. So the number  $N_0$  of equations of the form (22) is given by the number of sets  $\{t, t'\}$ . But by assumption  $t$  and  $t'$  should satisfy  $0 \leq t \leq r-2$  and  $c = 2t'$  with  $0 \leq c \leq 2r-3$ , which implies that  $0 \leq t' \leq r-2$ . Therefore,  $N_0$  is equal to the number  $(t, t')$  such that  $t \leq t'$  and thus we get:

$$N_0 = \sum_{t=0}^{r-2} \sum_{t'=t}^{r-2} = \frac{1}{2} (r-1)r. \quad (23)$$

Finally, by gathering all the cases we therefore obtain that:

$$N_{L \cap Q} = N_1 + N_0 = (r-1) \left( (u-1)r - 2^u + 2 \right) + \frac{1}{2} (r-1)r.$$

■

*Proposition 10:* For any integer  $r \geq 2$ , we have

$$\frac{1}{m} T_{\text{Goppa}}(r) = |\mathcal{B}_{2r}| - N_L - N_Q + N_{L \cap Q}.$$

*Proof:* Set  $u \stackrel{\text{def}}{=} \lfloor \log_2(2r-1) \rfloor$ . From Equation (5), we have

$$|\mathcal{B}_{2r}| = (2r-1) \left( (2u+1)r - 2^{u+1} + 1 \right)$$

which implies from Proposition 8

$$\begin{aligned} |\mathcal{B}_{2r}| - N_Q &= (2r-1) \left( (2u+1)r - 2^{u+1} + 1 \right) \\ &\quad - (2r-1) \left( (ru - 2^u + 1) \right) \\ &= (2r-1) \left( (u+1)r - 2^u \right). \end{aligned}$$



TABLE I  
A BINARY GOPPA CODE OF LENGTH  $n = 2^m$  AND DEGREE  $r < r_{\max}$  IS DISTINGUISHABLE FROM A RANDOM CODE.

$m$	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
$r_{\max}$	5	8	8	11	16	20	26	34	47	62	85	114	157	213	290	400
$r_{\text{crit}}$	5	6	8	11	14	19	25	34	46	62	84	114	156	214	293	402

Moreover, from Proposition 6 and Proposition 9, we can write:

$$\begin{aligned} N_L - N_{L \cap Q} &= (r-1)(2ur + 2 - 2^{u+1}) \\ &\quad - (r-1)\left(ur - \frac{r}{2} - 2^u + 2\right) \\ &= (r-1)\left(\left(u + \frac{1}{2}\right)r - 2^u\right). \end{aligned}$$

Therefore by gathering all these equalities we obtain:

$$|\mathcal{B}_{2r}| - (N_L + N_Q - N_{L \cap Q}) = r \left( \left(u + \frac{3}{2}\right)r - 2^u - \frac{1}{2} \right).$$

On the other hand from Proposition 4, we have

$$\frac{1}{m} T_{\text{Goppa}}(r) = \frac{1}{2} r ((2e+1)r - 2^e - 1)$$

where  $e = \lceil \log_2 r \rceil + 1$ . Using the basic inequality  $2r - 1 < 2r < 2(2r - 1)$ , we have therefore

$$\log_2(2r - 1) < \log_2(r) + 1 < \log_2(2r - 1) + 1$$

which finally implies  $\lceil \log_2 r \rceil = u$ . Thus,  $\frac{1}{m} T_{\text{Goppa}}(r) = \frac{1}{2} r ((2u+3)r - 2^{u+1} - 1)$ . ■

## IX. CONCLUSION AND CRYPTOGRAPHIC IMPLICATIONS

Based upon these experimental observations, it is straightforward to define a *distinguisher* between random codes, alternant codes and Goppa codes.

*Definition 10:* Let  $m$  and  $r$  be integers such that  $m \geq 1$  and  $r \geq 1$ . Let  $\mathbf{G}$  be a  $k \times n$  matrix whose entries are in  $\mathbb{F}_q$  with  $n \leq q^m$  and  $k \stackrel{\text{def}}{=} n - rm$ . Without loss of generality, we assume that  $\mathbf{G}$  is systematic i.e.,  $\mathbf{G} = (\mathbf{I}_k \mid \mathbf{P})$ . Let  $\mathcal{L}_{\mathbf{P}}$  be the linear system associated to  $\mathbf{G}$  as defined in (4), and  $D$  be the dimension of  $\text{Ker}(\mathcal{L}_{\mathbf{P}})$ . We define the *Random Code Distinguisher*  $\mathcal{D}$  as the mapping which takes as input  $\mathbf{G}$  and outputs  $\mathcal{D}(G)$  in  $\{-1, 0, 1\}$  such that:

$$\mathcal{D}(G) = \begin{cases} -1 & \text{if } D = T_{\text{alternant}} \\ 0 & \text{if } D = T_{\text{Goppa}} \\ 1 & \text{otherwise.} \end{cases} \quad (24)$$

The existence of a distinguisher for the specific case of binary Goppa codes is not valid for any value of  $r$  and  $m$  but tends to be true for codes that have a rate  $\frac{n-mr}{n}$  very close to one. We will elaborate on this point below. This kind of codes are mainly encountered with the signature scheme [1]. If we assume that the length  $n$  is equal to  $2^m$  and we denote by  $r_{\max}$  the smallest integer  $r$  such that  $N - T_{\text{Goppa}} \geq 2^m - mr$  then any binary Goppa code of degree  $r < r_{\max}$  can be distinguished (Table I). For example, the binary Goppa code obtained with  $m = 13$  and  $r = 19$  corresponding to a 90-bit security McEliece public key is distinguishable. More interestingly, all the keys proposed in [17] for the CFS signature scheme can be distinguished.

*Asymptotic Behaviour:* When the length  $n$  of the code goes to infinity an asymptotic formula can be derived for the smallest rate  $R_{\text{crit}}$  allowing distinguish a random code from an alternant code or a Goppa code. We derive such a formula when we assume for simplicity that the cardinality  $q$  of the base field is fixed and  $n$  is chosen as  $n = q^m$  (in practice  $n$  is chosen either in this way or at least of the same order as  $q^m$ ). We also assume that the dimension  $k$  of the code satisfies  $k = n - rm$ . We denote the rate  $k/n$  of the code by  $R$ . Finally, we also make the assumption that the dimensions  $D_{\text{alternant}}$  and  $D_{\text{Goppa}}$  are given by their theoretical values  $T_{\text{alternant}}$  and  $T_{\text{Goppa}}$  respectively and that the dimension of  $D_{\text{random}}$  is given by  $T_{\text{random}} \stackrel{\text{def}}{=} \max(0, \binom{m}{2} - k)$ . This critical rate  $r_{\text{crit}}$  corresponds to the smallest value of  $r$  for which  $T_{\text{random}}$  becomes bigger than  $T_{\text{alternant}}$  (asymptotically there will be no difference between Goppa codes or alternant codes). It holds that:

$$r_{\text{crit}} \stackrel{\text{def}}{=} \min\{r > 0 : T_{\text{random}} \geq T_{\text{alternant}}\}.$$

We let  $R_{\text{crit}} \stackrel{\text{def}}{=} \frac{n - r_{\text{crit}}m}{n} = 1 - \frac{r_{\text{crit}}m}{n}$ . Our claim is that

*Theorem 3:* Let  $n = q^m$ . When  $q$  is fixed and  $m$  tends to infinity, we have

$$r_{\text{crit}} = \sqrt{\frac{2q^m \log q}{m \log m}} (1 + o(1)),$$

$$R_{\text{crit}} = 1 - \sqrt{\frac{2m \log q}{q^m \log m}} (1 + o(1)).$$

where all logarithms are taken to base 2.

The proof of this theorem is given in the appendix.

In Table I, we have computed the value of  $\left\lceil \sqrt{\frac{2q^m \log q}{m \log m}} \right\rceil$  for several  $m$  ( $q$  is equal to 2). This shows that our approximation is rather close to  $r_{\text{max}}$  computed in practice (even for small values of  $m$ ).

*Concluding remarks.* We emphasize that the existence of such a distinguisher does not undermine the security of [2] and [1]. It only shows that the GD assumption should be used in any security reduction with great care.

It has also been observed in [38] that our distinguisher is equivalent to consider the dimension of the square code of the dual of the public code. It should be added that this notion has been used recently to cryptanalyze successfully two cryptographic schemes both of them relying on modified generalized Reed-Solomon codes [39], [40].

Finally, we mention that [41] shows that the natural reduction of GD to a hidden subgroup problem yields negligible information. As a consequence, they rule out the direct analogue of a quantum attack using the so-called Quantum Fourier Sampling (QFS) which breaks number theoretic problems [42]. More precisely, [41] shows that QFS has a negligible advantage against GD when the rate is  $\geq 1 - \frac{\log_q(n)^{3/2}}{\sqrt{5n}} = R_{\text{QFS}}$ . Whilst our result is somewhat contradictory with [41], it is interesting to observe that  $R_{\text{crit}}$  and the critical rate  $R_{\text{QFS}}$  share some similarities.

## APPENDIX

### A. Experimental Results

TABLE II  
 $q = 2$  AND  $m = 14$ .

$r$	3	4	5	6	7	8	9	10	11	12	13	14	15	16
$N$	861	1540	2415	3486	4753	6216	7875	9730	11781	14028	16471	19110	21945	24976
$k$	16342	16328	16314	16300	16286	16272	16258	16244	16230	16216	16202	16188	16174	16160
$D_{\text{random}}$	0	0	0	0	0	0	0	0	0	0	269	2922	5771	8816
$D_{\text{alternant}}$	42	126	308	560	882	1274	1848	2520	3290	4158	5124	6188	7350	8816
$T_{\text{alternant}}$	42	126	308	560	882	1274	1848	2520	3290	4158	5124	6188	7350	8610
$D_{\text{Goppa}}$	252	532	980	1554	2254	3080	4158	5390	6776	8316	10010	11858	13860	16016
$T_{\text{Goppa}}$	252	532	980	1554	2254	3080	4158	5390	6776	8316	10010	11858	13860	16016

TABLE III  
 $q = 2$  AND  $m = 14$ .

$r$	17	18	19	20	21	22	23	24	25	26	27	28	29	30
$N$	28203	31626	35245	39060	43071	47278	51681	56280	61075	66066	71253	76636	82215	87990
$k$	16146	16132	16118	16104	16090	16076	16062	16048	16034	16020	16006	15992	15978	15964
$D_{\text{random}}$	12057	15494	19127	22956	26981	31202	35619	40232	45041	50046	55247	60644	66237	72026
$D_{\text{alternant}}$	12057	15494	19127	22956	26981	31202	35619	40232	45041	50046	55247	60644	66237	72026
$T_{\text{alternant}}$	10192	11900	13734	15694	17780	19992	22330	24794	27384	30100	32942	35910	39004	42224
$D_{\text{Goppa}}$	18564	21294	24206	27300	30576	34034	37674	41496	45500	50046	55247	60644	66237	72026
$T_{\text{Goppa}}$	18564	21294	24206	27300	30576	34034	37674	41496	45500	49686	54054	58604	63336	68250

TABLE IV  
 $q = 2$  AND  $m = 15$ .

$r$	3	4	5	6	7	8	9	10	11	12	13	14	15	16
$N$	990	1770	2775	4005	5460	7140	9045	11175	13530	16110	18915	21945	25200	28680
$k$	32723	32708	32693	32678	32663	32648	32633	32618	32603	32588	32573	32558	32543	32528
$D_{\text{random}}$	0	0	0	0	0	0	0	0	0	0	0	0	0	0
$D_{\text{Goppa}}$	270	570	1050	1665	2415	3300	4455	5775	7260	8910	10725	12705	14850	17160
$T_{\text{Goppa}}$	270	570	1050	1665	2415	3300	4455	5775	7260	8910	10725	12705	14850	17160

TABLE V  
 $q = 2$  AND  $m = 15$ .

$r$	17	18	19	20	21	22	23	24	25	26	27	28	29	30
$N$	32385	36315	40470	44850	49455	54285	59340	64620	70125	75855	81810	87990	94395	101025
$k$	32513	32498	32483	32468	32453	32438	32423	32408	32393	32378	32363	32348	32333	32318
$D_{\text{random}}$	0	3817	7987	12382	17002	21847	26917	32212	37732	43477	49447	55642	62062	68707
$D_{\text{Goppa}}$	19890	22815	25935	29250	32760	36465	40365	44460	48750	53235	57915	62790	67860	73125
$T_{\text{Goppa}}$	19890	22815	25935	29250	32760	36465	40365	44460	48750	53235	57915	62790	67860	73125

TABLE VI  
 $q = 2$  AND  $m = 15$ .

$r$	31	32	33	34	35	36	37	38	39	40	41	42	43	44
$N$	107880	114960	122265	129795	137550	145530	153735	162165	170820	179700	188805	198135	207690	217470
$k$	32303	32288	32273	32258	32243	32228	32213	32198	32183	32168	32153	32138	32123	32108
$D_{\text{random}}$	75577	82672	89992	97537	105307	113302	121522	129967	138637	147532	156652	165997	175567	185362
$D_{\text{Goppa}}$	78585	84240	90585	97537	105307	113302	121522	129967	138637	147532	156652	165997	175567	185362
$T_{\text{Goppa}}$	78585	84240	90585	97155	103950	110970	118215	125685	133380	141300	149445	157815	166410	175230

TABLE VII  
 $q = 2$  AND  $m = 16$ .

$r$	3	4	5	6	7	8	9	10	11	12	13	14	15	16
$N$	1128	2016	3160	4560	6216	8128	10296	12720	15400	18336	21528	24976	28680	32640
$k$	65488	65472	65456	65440	65424	65408	65392	65376	65360	65344	65328	65312	65296	65280
$D_{\text{random}}$	0	0	0	0	0	0	0	0	0	0	0	0	0	0
$D_{\text{Goppa}}$	288	608	1120	1776	2576	3520	4752	6160	7744	9504	11440	13552	15840	18304
$T_{\text{Goppa}}$	288	608	1120	1776	2576	3520	4752	6160	7744	9504	11440	13552	15840	18304

TABLE VIII  
 $q = 2$  AND  $m = 16$ .

$r$	17	18	19	20	21	22	23	24	25	26	27	28	29	30
$N$	36856	41328	46056	51040	56280	61776	67528	73536	79800	86320	93096	100128	107416	114960
$k$	65264	65248	65232	65216	65200	65184	65168	65152	65136	65120	65104	65088	65072	65056
$D_{\text{random}}$	0	0	0	0	0	0	2360	8384	14664	21200	27992	35040	42344	49904
$D_{\text{Goppa}}$	21216	24336	27664	31200	34944	38896	43056	47424	52000	56784	61776	66976	72384	78000
$T_{\text{Goppa}}$	21216	24336	27664	31200	34944	38896	43056	47424	52000	56784	61776	66976	72384	78000

### B. Proof of Theorem 3

To prove Theorem 3 we will first use the following observation

*Lemma 8:* Let  $T_{\text{alternant}}$  be as defined in (5). Let also  $T_{\text{Goppa}}$  be as defined in (6). There exists constants  $K_1$  and  $K_2$  (resp.  $K'_1$  and  $K'_1$  and  $K'_2$ ) such that

$$mr^2(\log_q(r) + K_1) \leq T_{\text{alternant}} \leq mr^2(\log_q(r) + K_2), \quad (25)$$

$$mr^2(\log_q(r) + K'_1) \leq T_{\text{Goppa}} \leq mr^2(\log_q(r) + K'_2). \quad (26)$$

TABLE IX  
 $q = 2$  AND  $m = 16$ .

$r$	31	32	33	34	35	36	37	38	39	40	41	42	43
$N$	122760	130816	139128	147696	156520	165600	174936	184528	194376	204480	214840	225456	236328
$k$	65040	65024	65008	64992	64976	64960	64944	64928	64912	64896	64880	64864	64848
$D_{\text{random}}$	57720	65792	74120	82704	91544	100640	109992	119600	129464	139584	149960	160592	171480
$D_{\text{Goppa}}$	83824	89856	96624	103632	110880	118368	126096	134064	142272	150720	159408	168336	177504
$T_{\text{Goppa}}$	83824	89856	96624	103632	110880	118368	126096	134064	142272	150720	159408	168336	177504

TABLE X  
 $q = 4$  AND  $m = 6$ .

$r$	3	4	5	6	7	8	9	10	11	12	13	14	15	16
$N$	153	276	435	630	861	1128	1431	1770	2145	2556	3003	3486	4005	4560
$k$	4078	4072	4066	4060	4054	4048	4042	4036	4030	4024	4018	4012	4006	4000
$D_{\text{random}}$	0	0	0	0	0	0	0	0	0	0	0	0	0	560
$D_{\text{alternant}}$	6	18	60	120	198	294	408	540	690	858	1044	1248	1470	1710
$T_{\text{alternant}}$	6	18	60	120	198	294	408	540	690	858	1044	1248	1470	1710
$D_{\text{Goppa}}$	18	60	120	198	294	408	540	750	990	1260	1560	1890	2250	2640
$T_{\text{Goppa}}$	18	60	120	198	294	408	540	750	990	1260	1560	1890	2250	2640

TABLE XI  
 $q = 4$  AND  $m = 6$ .

$r$	17	18	19	20	21	22	23	24	25	26	27	28	29	30
$N$	5151	5778	6441	7140	7875	8646	9453	10296	11175	12090	13041	14028	15051	16110
$k$	3994	3988	3982	3976	3970	3964	3958	3952	3946	3940	3934	3928	3922	3916
$D_{\text{random}}$	1157	1790	2459	3164	3905	4682	5495	6344	7229	8150	9107	10100	11129	12194
$D_{\text{alternant}}$	2064	2448	2862	3306	3905	4682	5495	6344	7229	8150	9107	10100	11129	12194
$T_{\text{alternant}}$	2064	2448	2862	3306	3780	4284	4818	5382	5976	6600	7254	7938	8652	9396
$D_{\text{Goppa}}$	3060	3510	3990	4500	5040	5610	6210	6840	7500	8190	8910	9660	10440	11250
$T_{\text{Goppa}}$	3060	3510	3990	4500	5040	5610	6210	6840	7500	8190	8910	9660	10440	11250

*Proof:* We recall that

$$T_{\text{alternant}} \stackrel{\text{def}}{=} \frac{1}{2} m(r-1) \left( (2e+1)r - 2 \frac{q^{e+1} - 1}{q-1} \right)$$

where  $e \stackrel{\text{def}}{=} \lfloor \log_q(r-1) \rfloor$ . First, we remark that there exists some absolute constants  $K_3$  and  $K_4$  such that for all integers  $r \geq 2$

$$2r \log_q(r) + K_3 r \leq (2e+1)r - 2 \frac{q^{e+1} - 1}{q-1} \leq 2r \log_q(r) + K_4 r. \quad (27)$$

The upper bound is clear since:

$$(2e+1)r - 2 \frac{q^{e+1} - 1}{q-1} \leq 2r \log_q(r) + 2r.$$

For the lower bound, we remark that:

$$e \geq \log_q(r-1) - 1 = \log_q(r) + \log_q(1 - 1/r) - 1.$$

In addition:

$$\frac{q^{e+1} - 1}{q-1} \leq q \cdot q^e \leq r q.$$

As a consequence:

$$\left( (2e+1)r - 2 \frac{q^{e+1} - 1}{q-1} \right) \geq 2r \log_q(r) + r(2 \log_q(1 - 1/r) - 1 - 2q).$$

Finally, remark that  $\log_q(1 - 1/r)$  can be bounded from above by some (negative) constant. So, it holds that:

$$\left( (2e+1)r - 2 \frac{q^{e+1} - 1}{q-1} \right) \geq 2r \log_q(r) + K_3 r,$$

for some constant  $K_3$ .

Observe now that

$$\frac{1}{2}m(r-1)(2\log_q(r) + K_3r) = \frac{1}{2}(mr-m)(2\log_q(r) + K_3r) = \frac{1}{2}(2mr\log_q(r) + K_3mr^2 - 2m\log_q(r) - K_3mr).$$

The lower bound on  $T_{\text{alternant}}$  follows immediately from this. The expression can be lower bounded (resp. upper bounded) by a term of the form  $K_1mr^2$  (resp.  $K_2mr^2$ ) for some constant  $K_1$  (resp.  $K_2$ ). This holds for all positive integers  $r$ .

Finally, we recall that

$$T_{\text{Goppa}} \stackrel{\text{def}}{=} \begin{cases} \frac{1}{2}m(r-1)(r-2) = T_{\text{alternant}} & \text{for } r < q-1, \\ \frac{1}{2}mr((2e+1)r - 2q^e + 2q^{e-1} - 1) & \text{for } r \geq q-1, \end{cases}$$

with  $e$  being the unique integer such that:

$$(q-1)^2q^{e-2} < r \leq (q-1)^2q^{e-1}.$$

The bound (26) on can be proved in the same way. ■

From this lemma, we deduce that

*Lemma 9:* There exist two constants  $C_1$  and  $C_2$  such that for every  $r$  satisfying  $\binom{mr}{2} \geq n - mr$  we have

$$mr^2(m/2 - \log_q(r) + C_1) - q^m \leq T_{\text{random}} - T_{\text{alternant}} \quad , \quad (28)$$

$$mr^2(m/2 - \log_q(r) + C_2) - q^m \geq T_{\text{random}} - T_{\text{alternant}} \quad . \quad (29)$$

We also have the same inequalities when we replace  $T_{\text{random}} - T_{\text{alternant}}$  with  $T_{\text{random}} - T_{\text{Goppa}}$ .

*Proof:* For all positive integer values of  $r$  such that  $\binom{mr}{2} \geq n - mr$ , we have:

$$\begin{aligned} T_{\text{random}} &= N - k, \\ &= \binom{mr}{2} - q^m + mr, \\ &= m^2r^2/2 - q^m + mr/2. \end{aligned}$$

We can then conclude using Lemma 8. ■

From this lemma, we derive the following estimate for  $r_{\text{crit}}$ :

*Lemma 10:* When  $m$  goes to infinity we have

$$r_{\text{crit}} = \sqrt{\frac{2q^m \log q}{m \log m}} (1 + o(1)).$$

*Proof:* From Lemma 9, we know that:

$$T_{\text{random}} - T_{\text{alternant}} = mr^2(m/2 - \log_q(r)) - q^m + O(mr^2).$$

Let  $r_0 \stackrel{\text{def}}{=} \sqrt{\frac{2q^m \log q}{m \log m}}$ . It holds that:

$$\begin{aligned} 2\log_q(r_0) &= \log_q(2q^m \log q) - \log_q(m \log m), \\ &= m + \log_q(2 \log q) - \log_q(m) - \log_q(\log m). \end{aligned}$$

Thus:

$$\begin{aligned} mr_0^2(m/2 - \log_q r_0) - q^m &= \frac{2q^m \log q}{\log m} \left( m/2 - m/2 - \frac{\log(2 \log q)}{2 \log q} + \frac{\log m}{2 \log q} + \frac{\log \log m}{2 \log q} \right) - q^m, \\ &= \frac{q^m}{\log m} (\log \log m - \log(2 \log q)). \end{aligned}$$

We also observe that

$$mr_0^2 = \frac{2q^m \log q}{\log m}$$

is negligible compared to  $\frac{q^m}{\log m} (\log \log m - \log(2 \log q))$  when  $m$  goes to infinity. This can be used to show that  $T_{\text{random}} - T_{\text{alternant}}$  is positive for  $r = \lceil r_0 \rceil$  when  $m$  is large enough. Therefore for  $m$  large enough, we have  $r_{\text{crit}} \leq \lceil r_0 \rceil$ .

On the other hand, let  $\alpha$  be any positive constant  $< 1$ . We set:

$$r_\alpha \stackrel{\text{def}}{=} \sqrt{\frac{2\alpha q^m \log q}{m \log m}}.$$

Notice that the function  $f(x) = mx^2 (m/2 - \log_q x) - q^m$  can be shown to be increasing in the range  $(0, r_\alpha)$ . Therefore for every  $r \leq r_\alpha$ , we have

$$\begin{aligned} mr^2 (m/2 - \log_q r) - q^m &\leq \frac{2\alpha q^m \log q}{\log m} \left( m/2 - m/2 + \frac{\log m}{2 \log q} + \frac{\log \log m}{2 \log q} - \frac{\log(2\alpha \log q)}{2 \log q} \right) - q^m, \\ &\leq (\alpha - 1)q^m + \frac{2\alpha q^m \log q}{\log m} \left( \frac{\log \log m}{2 \log q} - \frac{\log(2\alpha \log q)}{2 \log q} \right), \\ &\leq (\alpha - 1)q^m + q^m \frac{\alpha \log \log m}{\log m}. \end{aligned}$$

Since  $mr^2 \leq \frac{2\alpha q^m \log q}{\log m}$ , it follows that any function of the form  $mr^2 (m/2 - \log_q r) - q^m + O(mr^2)$  will be negative for  $m$  large enough in the range  $(0, r_\alpha)$ .

This implies that  $r_{\text{crit}} \geq r_\alpha = \sqrt{\frac{2\alpha q^m \log q}{m \log m}}$  for  $m$  large enough. We deduce from this fact which holds for any  $0 < \alpha < 1$  and from the upper bound  $r_{\text{crit}} \leq \lceil r_0 \rceil = \left\lceil \sqrt{\frac{2q^m \log q}{m \log m}} \right\rceil$  that  $r_{\text{crit}} = \sqrt{\frac{2q^m \log q}{m \log m}} (1 + o(1))$  when  $m$  goes to infinity.

Finally, the proof of Theorem 3 is now obtained as follows. We have

$$\begin{aligned} R_{\text{crit}} &= \frac{q^m - mr_{\text{crit}}}{q^m}, \\ &= 1 - \sqrt{\frac{2m \log q}{q^m \log m}} (1 + o(1)). \end{aligned}$$

■

## REFERENCES

- [1] N. T. Courtois, M. Finiasz, and N. Sendrier, "How to achieve a McEliece-based digital signature scheme," in *ASIACRYPT*, vol. 2248, 2001, pp. 157–174.
- [2] R. J. McEliece, *A Public-Key System Based on Algebraic Coding Theory*. Jet Propulsion Lab, 1978, pp. 114–116, dSN Progress Report 44.
- [3] P. J. Lee and E. F. Brickell, "An observation on the security of McEliece's public-key cryptosystem," in *Advances in Cryptology - EUROCRYPT'88*, ser. Lecture Notes in Computer Science, vol. 330/1988. Springer, 1988, pp. 275–280.
- [4] J. S. Leon, "A probabilistic algorithm for computing minimum weights of large error-correcting codes," *IEEE Transactions on Information Theory*, vol. 34, no. 5, pp. 1354–1359, 1988.
- [5] J. Stern, "A method for finding codewords of small weight," in *Coding Theory and Applications*, ser. Lecture Notes in Computer Science, G. D. Cohen and J. Wolfmann, Eds., vol. 388. Springer, 1988, pp. 106–113.
- [6] A. Canteaut and F. Chabaud, "A new algorithm for finding minimum-weight words in a linear code: Application to McEliece's cryptosystem and to narrow-sense BCH codes of length 511," *IEEE Transactions on Information Theory*, vol. 44, no. 1, pp. 367–378, 1998.
- [7] D. J. Bernstein, T. Lange, and C. Peters, "Attacking and defending the McEliece cryptosystem," in *PQCrypto*, ser. LNCS, vol. 5299, 2008, pp. 31–46.
- [8] —, "Smaller decoding exponents: Ball-collision decoding," in *CRYPTO*, ser. Lecture Notes in Computer Science, P. Rogaway, Ed., vol. 6841. Springer, 2011, pp. 743–760.
- [9] A. May, A. Meurer, and E. Thome, "Decoding random linear codes in  $\tilde{O}(2^{0.054n})$ ," in *ASIACRYPT*, ser. Lecture Notes in Computer Science, D. H. Lee and X. Wang, Eds., vol. 7073. Springer, 2011, pp. 107–124.
- [10] A. Becker, A. Joux, A. May, and A. Meurer, "Decoding random binary linear codes in  $2^{n/20}$ ; How  $1 + 1 = 0$  improves information set decoding," in *EUROCRYPT*, ser. Lecture Notes in Computer Science, D. Pointcheval and T. Johansson, Eds., vol. 7237. Springer, 2012, pp. 520–536.
- [11] E. Berlekamp, R. McEliece, and H. van Tilborg, "On the inherent intractability of certain coding problems," *IEEE Transactions on Information Theory*, vol. 24, no. 3, pp. 384–386, May 1978.
- [12] J. Gibson, "Equivalent Goppa codes and trapdoors to McEliece's public key cryptosystem," in *Advances in Cryptology EUROCRYPT 91*, ser. Lecture Notes in Computer Science, D. Davies, Ed. Springer Berlin / Heidelberg, 1991, vol. 547, pp. 517–521.
- [13] P. Loidreau and N. Sendrier, "Weak keys in the McEliece public-key cryptosystem," *IEEE Transactions on Information Theory*, vol. 47, no. 3, pp. 1207–1211, 2001.
- [14] R. Nojima, H. Imai, K. Kobara, and K. Morozov, "Semantic security for the McEliece cryptosystem without random oracles," *Des. Codes Cryptography*, vol. 49, no. 1-3, pp. 289–305, 2008.
- [15] R. Dowsley, J. Müller-Quade, and A. C. A. Nascimento, "A CCA2 secure public key encryption scheme based on the McEliece assumptions in the standard model," in *CT-RSA*, 2009, pp. 240–251.
- [16] L. Dallot, "Towards a concrete security proof of Courtois, Finiasz and Sendrier signature scheme," in *WEWoRC*, 2007, pp. 65–77.
- [17] M. Finiasz and N. Sendrier, "Security bounds for the design of code-based cryptosystems," in *Asiacrypt 2009*, ser. LNCS, M. Matsui, Ed., vol. 5912. Springer, 2009, pp. 88–105.
- [18] J.-C. Faugère, A. Otmani, L. Perret, and J.-P. Tillich, "Algebraic cryptanalysis of McEliece variants with compact keys," in *EUROCRYPT*, ser. Lecture Notes in Computer Science, H. Gilbert, Ed., vol. 6110. Springer, 2010, pp. 279–298.
- [19] V. D. Goppa, "A new class of linear correcting codes," *Probl. Peredachi Inf.*, vol. 6, no. 3, 1970.
- [20] F. J. MacWilliams and N. J. A. Sloane, *The Theory of Error-Correcting Codes*, 5th ed. Amsterdam: North-Holland, 1986.
- [21] N. Patterson, "The algebraic decoding of Goppa codes," *IEEE Transactions on Information Theory*, vol. 21, no. 2, pp. 203–207, 1975.
- [22] J. van Tilburg, "On the McEliece public-key cryptosystem," in *CRYPTO '88: Proceedings of the 8th Annual International Cryptology Conference on Advances in Cryptology*. London, UK: Springer-Verlag, 1990, pp. 119–131.

- [23] A. Canteaut and H. Chabanne, "A further improvement of the work factor in an attempt at breaking McEliece's cryptosystem," in *EUROCODE 94*. INRIA, 1994, pp. 169–173.
- [24] A. Canteaut and F. Chabaud, "Improvements of the attacks on cryptosystems based on error-correcting codes," INRIA, Tech. Rep. 95–21, 1995.
- [25] I. Dumer, "Suboptimal decoding of linear codes : partition techniques," *IEEE Transactions on Information Theory*, vol. 42, no. 6, pp. 1971–1986, 1996.
- [26] A. Canteaut and N. Sendrier, "Cryptanalysis of the original McEliece cryptosystem," in *Advances in Cryptology - ASIACRYPT'98*, ser. Lecture Notes in Computer Science, no. 1514. Springer-Verlag, 1998, pp. 187–199.
- [27] J. Faugère, V. Gauthier-Umana, A. Otmani, L. Perret, and J. Tillich, "Distinguisher for high rate McEliece cryptosystems," in *Proceedings of the 2011 IEEE Information Theory Workshop (ITW 2011)*, Paraty, Brazil, Oct. 16-20 2011.
- [28] N. Sendrier, "Finding the permutation between equivalent linear codes: The support splitting algorithm," *IEEE Transactions on Information Theory*, vol. 46, no. 4, pp. 1193–1203, 2000.
- [29] K. Kobara and H. Imai, "Semantically secure McEliece public-key cryptosystems-conversions for McEliece PKC," in *Public Key Cryptography, 4th International Workshop on Practice and Theory in Public Key Cryptography, PKC 2001, Cheju Island, Korea, February 13-15, 2001, Proceedings*, ser. Lecture Notes in Computer Science, K. Kim, Ed., vol. 1992. Springer, 2001, pp. 19–35.
- [30] R. Nojima, H. Imai, K. Kobara, and K. Morozov, "Semantic security for the McEliece cryptosystem without random oracles," *Des. Codes Cryptography*, vol. 49, no. 1-3, pp. 289–305, 2008.
- [31] R. Dowsley, J. Müller-Quade, and A. C. A. Nascimento, "A CCA2 secure public key encryption scheme based on the McEliece assumptions in the standard model," in *CT-RSA*, 2009, pp. 240–251.
- [32] T. P. Berger, P. Cayrel, P. Gaborit, and A. Otmani, "Reducing key length of the McEliece cryptosystem," in *Progress in Cryptology - Second International Conference on Cryptology in Africa (AFRICACRYPT 2009)*, ser. Lecture Notes in Computer Science, B. Preneel, Ed., vol. 5580, Gammarrh, Tunisia, Jun. 21-25 2009, pp. 77–97.
- [33] R. Misoczki and P. S. L. M. Barreto, "Compact McEliece keys from Goppa codes," in *Selected Areas in Cryptography (SAC 2009)*, Calgary, Canada, Aug. 13-14 2009.
- [34] J.-C. Faugère, "A new efficient algorithm for computing gröbner bases (f4)," *Journal of Pure and Applied Algebra*, vol. 139(1-3), pp. 61–88, 1999.
- [35] —, "A new efficient algorithm for computing gröbner bases without reduction to zero : F5," in *ISSAC'02*. ACM press, 2002, pp. 75–83.
- [36] W. Bosma, J. J. Cannon, and C. Playoust, "The Magma algebra system I: The user language," *J. Symb. Comput.*, vol. 24, no. 3/4, pp. 235–265, 1997.
- [37] C. Cooper, "On the distribution of rank of a random matrix over a finite field," *Random Struct. Algorithms*, vol. 17, no. 3-4, pp. 197–212, 2000.
- [38] I. Marquez-Corbella and R. Pellikaan, "Error-correcting pairs for a public-key cryptosystem," in *CBC 2012, Code-based Cryptography Workshop 2012*, paper available on <http://www.win.tue.nl/~ruudp/paper/60.pdf>.
- [39] V. Gauthier, A. Otmani, and J.-P. Tillich, "A distinguisher-based attack on a variant of McEliece's cryptosystem based on Reed-Solomon codes," arxiv cs.CR, 2012, <http://arxiv.org/abs/1204.6459>.
- [40] —, "A distinguisher-based attack of a homomorphic encryption scheme relying on Reed-Solomon codes," *CoRR*, 2012, <http://arxiv.org/abs/1203.6686>.
- [41] H. Dinh, C. Moore, and A. Russell, "McEliece and Niederreiter cryptosystems that resist quantum Fourier sampling attacks," in *CRYPTO*, ser. Lecture Notes in Computer Science, P. Rogaway, Ed., vol. 6841. Springer, 2011, pp. 761–779.
- [42] P. W. Shor, "Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer," *SIAM J. Comput.*, vol. 26, no. 5, pp. 1484–1509, 1997.
- [43] P. Rogaway, Ed., *Advances in Cryptology - CRYPTO 2011 - 31st Annual Cryptology Conference, Santa Barbara, CA, USA, August 14-18, 2011. Proceedings*, ser. Lecture Notes in Computer Science, vol. 6841. Springer, 2011.