



**HAL**  
open science

## Using Reed-Solomon codes in the $(U | U + V)$ construction and an application to cryptography

Irene Márquez-Corbella, Jean-Pierre Tillich

### ► To cite this version:

Irene Márquez-Corbella, Jean-Pierre Tillich. Using Reed-Solomon codes in the  $(U | U + V)$  construction and an application to cryptography. International Symposium on Information Theory, Jul 2016, Barcelona, Spain. hal-01410201

**HAL Id: hal-01410201**

**<https://inria.hal.science/hal-01410201>**

Submitted on 6 Dec 2016

**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.

# Using Reed-Solomon codes in the $(U | U + V)$ construction and an application to cryptography

Irene Márquez-Corbella and Jean-Pierre Tillich

Inria, SECRET Project, 75589 Paris Cedex 12, France.

Email: irene.marquez-corbella@inria.fr and jean-pierre.tillich@inria.fr

**Abstract**—In this paper we present a modification of Reed-Solomon codes that beats the Guruswami-Sudan  $1 - \sqrt{R}$  decoding radius of Reed-Solomon codes at low rates  $R$ . The idea is to choose Reed-Solomon codes  $U$  and  $V$  with appropriate rates in a  $(U | U + V)$  construction and to decode them with the Koetter-Vardy soft information decoder. We suggest to use a slightly more general version of these codes (but which has the same decoding performance as the  $(U | U + V)$ -construction) for being used in code-based cryptography, namely to build a McEliece scheme. The point is here that these codes not only perform nearly as well (or even better in the low rate regime) as Reed-Solomon codes, but also that their structure seems to avoid the Sidelnikov-Shestakov attack which broke a previous McEliece proposal based on generalized Reed-Solomon codes.

## I. INTRODUCTION

*Improving upon the error correction performance of RS codes.* Reed-Solomon (RS) codes are among the most extensively used error correcting codes. It has long been known how to decode them up to half the minimum distance. This gives a decoding algorithm that corrects a fraction  $\frac{1-R}{2}$  of errors in an RS code of rate  $R$ . However it is only in the late nineties that a breakthrough was obtained in this setting with Sudan's algorithm [17] and its improvement in [8] who showed how to go beyond this barrier with an algorithm which in its [8] version decodes any fraction of errors smaller than  $1 - \sqrt{R}$ . Later on, it was shown that this decoding algorithm could also be modified a little bit in order to cope with soft information on the errors [9]. Then it was realized in [14] that by a slight modification of RS codes and by an increase of the alphabet size it was possible to beat the  $1 - \sqrt{R}$  decoding radius. Their new family of codes is list decodable beyond this radius for low rate. Then, [7] improved on these codes by presenting a new family of codes, namely *folded RS codes* with a polynomial time decoding algorithm achieving the list decoding capacity  $1 - R - \epsilon$  for every rate  $R$  and  $\epsilon > 0$ .

The first purpose of this paper is to present another modification of RS codes that improves the fraction of errors that can be corrected. It consists in using RS codes in a  $(U | U + V)$  construction. We will show that, in the low rate regime, this class of codes outperforms rather significantly a classical RS code decoded with the Guruswami and Sudan decoder [8]. The point is that this  $(U | U + V)$  code can be decoded in two steps :

- 1) First by subtracting the left part  $y_1$  to the right part  $y_2$  of the received vector  $(y_1|y_2)$  and decoding it with respect to  $V$ . In such a case, we are left with decoding a RS code with about twice as many errors.
- 2) Secondly, once we have recovered the right part  $v$  of the codeword, we can get a word  $(y_1, y_2 - v)$  which should match two copies of a same word  $u$  of  $U$ . We can model this decoding problem by having some soft information.

It turns out that the last channel error model is much less noisy than the original  $q$ -ary symmetric channel we started with. This soft information can be used in Koetter and Vardy's decoding algorithm. By this means we can choose  $U$  to be a RS code of much bigger rate than  $V$ . All in all, it turns out that by choosing  $U$  and  $V$  with appropriate rates we can beat the  $1 - \sqrt{R}$  bound in the low-rate regime.

It should be noted however that beating this  $1 - \sqrt{R}$  bound comes at the cost of having now an algorithm which does not work as for the aforementioned papers [7], [8], [14], [17] for every error of a given weight (the so called adversarial error model) but with probability  $1 - o(1)$  for errors of a given weight. However contrarily to [7], [14] which results in a significant increase of the alphabet size of the code, our alphabet size actually decreases when compared to a RS code: it can be half of the code length and can be even smaller when we apply this construction recursively. Indeed, we will show that we can even improve the error correction performances by applying this construction again to the  $U$  and  $V$  components, i.e. we can choose  $U$  to be a  $(U_1|U_1 + V_1)$  code and we replace in the same way the RS code  $V$  by a  $(U_2|U_2 + V_2)$  where  $U_1, U_2, V_1, V_2$  are RS codes.

*Application to cryptography.* In a second part of the paper we show how to use such codes (or codes derived by this approach) for cryptographic purposes, i.e. in a McEliece cryptosystem [11]. Recall that this public-key cryptosystem becomes more and more fashionable due to the threats on the most popular public key cryptosystems used today, namely RSA or DSA and ECDSA that would be completely broken by Shor's algorithm [15] if a large scale quantum computer could be built. Indeed, it is unlikely that a quantum computer would be able to threaten the security of the McEliece scheme because it is based on an NP-complete problem, namely decoding a random

linear code.

Probably one of the main drawbacks of McEliece when compared to RSA, DSA or ECDSA is its rather large key size. There have been several attempts to decrease the key size either by moving to more structured codes or to codes which have better error correction radius [2], [13]. Many of the structured algebraic proposals have been broken (see for instance [6]) but some of the quasi-cyclic code families that rely on modified LDPC codes or MDPC codes [1], [12] seem to resist cryptanalysis up to now. Relying on codes with better decoding performance met a similar fate, since here again many proposals of this kind have been broken. For instance [13] suggests to replace the binary Goppa codes of the original McEliece cryptosystem by Generalized RS codes (GRS) because of their much better decoding performance, but it got broken in [16].

There have been several attempts to repair GRS codes in this context either by adding random columns to the generator matrix of a GRS code [19] or by multiplying this generator matrix by the inverse of a sparse matrix with small average row weight  $m$  [2]. The [19] attempt got broken in [4] and the parameters of [2] got broken in [5] because  $m$  was chosen to be too small. The problem with the approach in [2] is that the attack of [5] fails when  $m = 2$ , but the solution is then no more competitive when compared to a Goppa code because the decoding radius gets also scaled down by a multiplicative factor of  $m$  when compared to a GRS code.

We suggest here to revive the approach in [2] with a generalized  $(U | U + V)$  scheme based on RS codes that has basically the same decoding capacity as a RS code and that looks in many respects like the [2] scheme with  $m = 2$ . This approach is also related to the approach pioneered by Wang in [18]. His code can be viewed as a certain subcode of our  $(U | U + V)$  construction. However the decrease of the code rate results in a significant deterioration of the key size when compared to a code with the same error correction capacity as an RS code.

Due to space reasons, proofs are omitted. For further details, we refer the readers to [10]. A linear code of length  $n$ , dimension  $k$  and distance  $d$  over a finite field  $\mathbb{F}_q$  is referred to as an  $[n, k, d]_q$ -code. We will also frequently use the following notation

**Notation 1.** For a vector  $\mathbf{x}$  we denote by  $x(i)$  the  $i$ -th coordinate of  $\mathbf{x}$ .

## II. $(U | U + V)$ -CONSTRUCTION

This Section is only stated for the  $q$ -ary symmetric channel model. In this section, we recall a few facts about the  $(U | U + V)$  construction and its decoding.

**Definition 1.** Let  $U$  be an  $[n, k_u, d_u]_q$  code and  $V$  be an  $[n, k_v, d_v]_q$  code. We define the  $(U | U + V)$ -construction of  $U$  and  $V$  as the linear code:

$$\mathcal{C} = \{(\mathbf{u} | \mathbf{u} + \mathbf{v}) | \mathbf{u} \in U \text{ and } \mathbf{v} \in V\}.$$

The code  $\mathcal{C}$  has parameters  $[2n, k_u + k_v, \min\{2d_u, d_v\}]_q$ .

### A. Soft-decision decoding of $(U | U + V)$ codes

Let  $U$  and  $V$  be two codes with parameters  $[n, k_u, d_u]_q$  and  $[n, k_v, d_v]_q$ , respectively and  $\mathcal{C} \stackrel{\text{def}}{=} (U | U + V)$ . Suppose we transmit the codeword  $(\mathbf{u} | \mathbf{u} + \mathbf{v}) \in \mathcal{C}$  over a noisy channel and we receive the vector:  $\mathbf{y} = (\mathbf{y}_1 | \mathbf{y}_2) = (\mathbf{u} | \mathbf{u} + \mathbf{v}) + (\mathbf{e}_1 | \mathbf{e}_2)$ .

Decoding proceeds in two steps:

- 1) We combine  $\mathbf{y}_1$  and  $\mathbf{y}_2$  to find  $\mathbf{v}$ . That is, we decode  $\mathbf{y}_2 - \mathbf{y}_1 = \mathbf{v} + \mathbf{e}_2 - \mathbf{e}_1$  with respect to  $V$ . In the case of a soft decoder for  $V$  we compute first the probability  $\text{prob}(v(i) = \alpha | y_1(i), y_2(i))$  for all  $\alpha \in \mathbb{F}_q$ . This information is then used in a soft decoder for  $V$ .
- 2) We subtract  $(\mathbf{0} | \mathbf{v})$  to  $(\mathbf{y}_1 | \mathbf{y}_2)$  to get  $(\mathbf{u} + \mathbf{e}_1 | \mathbf{u} + \mathbf{e}_2) = (\mathbf{z}_1 | \mathbf{z}_2)$ . This is a noisy version of  $(\mathbf{u} | \mathbf{u})$ . We compute now for all  $\alpha \in \mathbb{F}_q$  and all coordinates  $i$  the probabilities  $\text{prob}(u(i) = \alpha | z_1(i), z_2(i))$  which is then passed to a soft decoder for  $U$ .

Let us explain how these probabilities can be computed. We assume that the noise model is given by a discrete memoryless channel (DMC) with input alphabet  $\mathbb{F}_q$  and output alphabet  $\mathcal{Y}$ . The received vector is denoted by  $\mathbf{y} = (y_1, \dots, y_n) \in \mathcal{Y}^n$  and the channel model specifies the transition probabilities with the following matrix  $\Pi_{\mathbf{y}}$

$$\Pi_{\mathbf{y}}^i(\alpha) = \text{prob}(\alpha | y_i) \quad \text{for } i = 1, \dots, n \text{ and } \alpha \in \mathbb{F}_q.$$

$\Pi_{\mathbf{y}}^i$  denotes here the  $i$ -th column of  $\Pi_{\mathbf{y}}$  and  $\Pi_{\mathbf{y}}^i(\alpha)$  refers to the entry in the  $i$ -th column and row indexed by  $\alpha \in \mathbb{F}_q$ .

We will refer to  $\Pi$  as the  $q \times n$  **reliability matrix** of the codewords symbols. We will see below that this reliability matrix can also be obtained through the  $(U | U + V)$  decoding process. We will particularly be interested here in the  $q$ -ary symmetric channel model with crossover probability  $p$  ( $q$ -SC $_p$ ). The reliability matrix  $\Pi_{\mathbf{y}}$  for  $q$ -SC $_p$  is defined as follows:

$$\Pi_{\mathbf{y}}^i(\alpha) = \text{prob}(\alpha | y_i) = \begin{cases} 1 - p & \text{if } \alpha = y_i \\ \frac{p}{q-1} & \text{if } \alpha \neq y_i \end{cases}$$

Let us recall now how the reliability matrices for the decoder of  $U$  and  $V$  are computed from the initial reliability matrix.

a) *Reliability matrix for the  $V$ -decoder:* We call in what follows the error model for the  $V$ -decoder the **sum model** and denote the associated reliability matrix by  $\Pi \oplus \Pi$  when  $\Pi$  is the initial reliability matrix. Recall that before decoding, for each symbol  $X$  of  $V$  that we want to decode we compute the difference of two symbols  $X_1$  and  $X_2$  of the  $(U | U + V)$  code:  $X = X_2 - X_1$ . For each of these symbols we have a reliability information  $\text{prob}(X_1 = \alpha | Y_1)$  and  $\text{prob}(X_2 = \beta | Y_2)$  where  $Y_1$  and  $Y_2$  are random variables that are initially the received symbols corresponding to  $X_1$  and  $X_2$  after transmission on the noisy channel but that become sets of received symbols when we iterate the  $(U | U + V)$  construction as will be

seen. When  $X_1$  and  $X_2$  are uniformly distributed it can be verified that

$$\text{prob}(X = \alpha | Y_1, Y_2) = \sum_{\beta \in \mathbb{F}_q} \text{prob}(X_1 = \beta | Y_1) \cdot \text{prob}(X_2 = \alpha + \beta | Y_2)$$

This leads to the following definition.

$$(\Pi \oplus \Pi)_{\mathbf{y}}^i(\alpha) \stackrel{\text{def}}{=} \sum_{\beta \in \mathbb{F}_q} \Pi_{\mathbf{y}_1}^i(\beta) \cdot \Pi_{\mathbf{y}_2}^i(\alpha + \beta)$$

where  $\mathbf{y}_1$  and  $\mathbf{y}_2$  are the realizations of the channel transmission of  $u$  and  $u + v$  respectively.

b) *Reliability matrix for the U-decoder*: The computation of  $\text{prob}(u(i) = \alpha | z_1(i), z_2(i))$  can be performed by computing the probability that a uniformly distributed random variable over  $\mathbb{F}_q$  is equal to  $\alpha$  given two received symbols  $y_1$  and  $y_2$  for  $X$  sent over two memoryless channels. This probability is readily seen to be equal to

$$\frac{\text{prob}(X = \alpha | y_1) \cdot \text{prob}(X = \alpha | y_2)}{\sum_{\beta \in \mathbb{F}_q} \text{prob}(X = \beta | y_1) \cdot \text{prob}(X = \beta | y_2)}$$

We denote by  $\Pi \times \Pi$  the reliability matrix (the input) to a soft-decision decoding algorithm for the code  $U$ . Thus, each element of the reliability matrix  $\Pi \times \Pi$  related to the aforementioned quantities  $\mathbf{y}$  and  $\mathbf{v}$  is defined by:

$$(\Pi \times \Pi)_{\mathbf{y}, \mathbf{v}}^i(\alpha) \stackrel{\text{def}}{=} \frac{\Pi_{\mathbf{y}_1}^i(\alpha) \cdot \Pi_{\mathbf{y}_2}^i(\alpha + v(i))}{\sum_{\beta \in \mathbb{F}_q} \Pi_{\mathbf{y}_1}^i(\beta) \cdot \Pi_{\mathbf{y}_2}^i(\beta + v(i))}.$$

To simplify notation we will generally avoid the dependency on  $\mathbf{v}$  and simply write  $(\Pi \times \Pi)_{\mathbf{y}}$ .

### B. Algebraic-soft decision decoding of RS codes

Let us recall how the Koetter-Vardy soft decoder [9] can be analyzed. By [9, Theorem 12] their decoding algorithm outputs a list that contains the codeword  $\mathbf{c} \in C$  if

$$\frac{\langle \Pi, \lfloor \mathbf{c} \rfloor \rangle}{\sqrt{\langle \Pi, \Pi \rangle}} \geq \sqrt{k-1} + o(1)$$

as the codelength  $n$  tends to infinity, where  $\lfloor \mathbf{c} \rfloor$  represents a  $q \times n$  matrix with entries  $c_{i,\alpha} = 1$  if  $c_i = \alpha$ , and 0 otherwise; and  $\langle A, B \rangle \stackrel{\text{def}}{=} \sum_{i=1}^q \sum_{j=1}^n a_{i,j} b_{i,j}$ . We will consider here only discrete symmetric channel models that are defined below. Let us first introduce some notation.

**Notation 2** (Probability error vector of a DMC). *For a given DMC with  $q$ -ary inputs we denote by  $\pi$  the probability vector  $\pi = (\text{prob}(x = \alpha | y))_{\alpha \in \mathbb{F}_q}$  where  $x$  is the symbol that has been sent through the channel and  $y$  is the received vector. For a vector  $\mathbf{x} = (x_\beta)_{\beta \in \mathbb{F}_q}$  we denote by  $\mathbf{x}^{+\alpha}$  the vector  $\mathbf{x}^{+\alpha} = (x_{\beta+\alpha})_{\beta \in \mathbb{F}_q}$ .*

By viewing  $\pi$  as a random variable, we define as in [3] a symmetric channel by

**Definition 2** (Discrete symmetric channel with  $q$ -ary inputs). A DMC with  $q$ -ary inputs is said to be symmetric if and only if for any  $\alpha$  in  $\mathbb{F}_q$  we have

$$p_\alpha \text{prob}(\pi = \mathbf{p}) = p_0 \text{prob}(\pi = \mathbf{p}^{+\alpha}). \quad (1)$$

Note that this implies that for a discrete symmetric channel, for any possible realization  $\mathbf{p}$  of the probability vector  $\pi$  (i.e. when  $\text{prob}(\pi = \mathbf{p}) \neq 0$ ) we necessarily have  $p_0 \neq 0$ . It is proved in [3] that symmetric channels are closed under the  $+$  and  $\times$  operations on channels defined in Subsection II-A. We give now the asymptotic behavior for a symmetric channel of the Koetter-Vardy decoder. The proof of this theorem is found in the full version of this paper [10].

**Theorem 3.** *Let  $(\mathcal{C}_n)_{n \geq 1}$  be an infinite family of Reed-Solomon codes of rate  $\leq R$ . Denote by  $q_n$  the alphabet size of  $\mathcal{C}_n$  that is assumed to be a non decreasing sequence that goes to infinity with  $n$ . Consider an infinite family of  $q_n$ -ary symmetric channels with associated probability error vectors  $\pi_n$  such that  $\mathbb{E}(\|\pi_n\|^2)$  has a limit as  $n$  tends to infinity. Let*

$$C_{KV} \stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} \mathbb{E}(\|\pi_n\|^2).$$

*This infinite family of codes can be decoded correctly by the Koetter-Vardy decoding algorithm with probability  $1 - o(1)$  as  $n$  tends to infinity as soon as there exists  $\epsilon > 0$  such that*

$$R \leq C_{KV} - \epsilon.$$

*Remark 1.* Let us observe that for the  $q$ -SC $_p$  we have

$$\mathbb{E}(\|\pi\|^2) = (1-p)^2 + (q-1) \frac{p^2}{(q-1)^2} = (1-p)^2 + \mathcal{O}\left(\frac{1}{q}\right).$$

By letting  $q$  going to infinity, we recover in this way the performances of the Guruswami-Sudan algorithm which works as soon as  $R < (1-p)^2$ .

## III. CORRECTING ERRORS BEYOND THE GURUSWAMI-SUDAN BOUND

### A. The $(U | U + V)$ -construction

Now suppose we choose  $U$  and  $V$  as RS codes in a  $(U | U + V)$  construction. We start with a  $q$ -ary symmetric channel with error probability  $p$ . Recall that the reliability matrix for the  $U$ -decoder is  $\Pi_1 = \Pi \times \Pi$  whereas for the  $V$ -decoder it is  $\Pi_2 = \Pi \oplus \Pi$ .

**Lemma 4.** *Let  $\pi_U$  and  $\pi_V$  be the probability vectors corresponding to decoding the codes  $U$  and  $V$  respectively.*

- *The channel error model of the code  $V$  is a  $q$ -SC $_{p'}$  with  $p_1 \stackrel{\text{def}}{=} 2p - p^2$  and*

$$\mathbb{E}(\|\pi_V\|^2) = (1-p_1)^2 + \mathcal{O}\left(\frac{1}{q}\right) = (1-p)^4 + \mathcal{O}\left(\frac{1}{q}\right).$$

- *For the channel error model of the code  $U$  we have*

$$\mathbb{E}(\|\pi_U\|^2) = \frac{(2+p)(1-p)^2}{2-p} + \mathcal{O}\left(\frac{1}{q}\right).$$

By letting  $q$  going to infinity, we recover the performance of the  $(U | U + V)$  construction which works as soon as

$$R < \frac{(p^3 - 4p^2 + 4p - 4)(1-p)^2}{2(p-2)}$$

From Fig. 2 we deduce that the  $(U | U + V)$  decoder outperforms the RS decoder with Guruswami-Sudan as soon as  $R < 0.168$ .

### B. Recursive application of the $(U | U + V)$ construction

Now we will study what happens over the  $q$ -SC $_p$  if we apply recursively the  $(U | U + V)$  construction. So we start with a  $(U | U + V)$  code, we choose  $U$  to be a  $(U_1 | U_1 + V_1)$  code and  $V$  to be a  $(U_2 | U_2 + V_2)$  code, where  $U_1, U_2, V_1$  and  $V_2$  are RS codes over the same alphabet  $\mathbb{F}_q$  and of the same length. In other words, we look for a code of the form

$$(U_1 | U_1 + V_1 | U_1 + U_2 | U_1 + U_2 + V_1 + V_2) = \{(\mathbf{u}_1 | \mathbf{u}_1 + \mathbf{v}_1 | \mathbf{u}_1 + \mathbf{u}_2 | \mathbf{u}_1 + \mathbf{u}_2 + \mathbf{v}_1 + \mathbf{v}_2) : \mathbf{u}_i \in U_i, \mathbf{v}_i \in V_i\}$$

From now on, we will refer to this structure as the  $(U | U + V)$ -second level construction.

On the following we obtain the channel error models for decoding  $U_1, V_1, U_2$  and  $V_2$  respectively, their reliability matrices are given by  $\Pi_1 \times \Pi_1, \Pi_1 \oplus \Pi_1, \Pi_2 \times \Pi_2$  and  $\Pi_2 \oplus \Pi_2$  respectively (see Fig. 1).

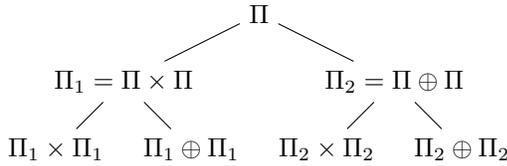


Fig. 1: The channel error models for decoding  $U_1$  ( $\Pi_1 \times \Pi_1$ ),  $V_1$  ( $\Pi_1 \oplus \Pi_1$ ),  $U_2$  ( $\Pi_2 \times \Pi_2$ ) and  $V_2$  ( $\Pi_2 \oplus \Pi_2$ ).

**Lemma 5.** Let  $\pi_{U_i}$  and  $\pi_{V_i}$  be the probability vectors corresponding to decoding the  $U_i$ 's and  $V_i$ 's.

- The channel error model of the code  $V_2$  is a  $q$ -SC $_{p_2}$  with  $p_2 \stackrel{\text{def}}{=} 2p_1 - p_1^2$  and

$$\mathbb{E} \left( \|\pi_{V_2}\|^2 \right) = (1 - p_2)^2 + \mathcal{O} \left( \frac{1}{q} \right) = (1 - p)^8 + \mathcal{O} \left( \frac{1}{q} \right);$$

- $\mathbb{E} \left( \|\pi_{U_2}\|^2 \right) = \frac{(2+p_1)(1-p_1)^2}{(2-p_1)} + \mathcal{O} \left( \frac{1}{q} \right);$
- $\mathbb{E} \left( \|\pi_{V_1}\|^2 \right) = \frac{(2+p)^2(1-p)^4}{(2-p)^2} + \mathcal{O} \left( \frac{1}{q} \right);$
- $\mathbb{E} \left( \|\pi_{U_1}\|^2 \right) = \frac{(5p^3 - 6p^2 - 5p - 4)(1-p)^2}{3p - 4} + \mathcal{O} \left( \frac{1}{q} \right).$

By letting  $q$  going to infinity, we recover the performance of this construction which works (in the asymptotic regime) as soon as

$$R < \lim_{q \rightarrow \infty} \frac{\sum_{i=1}^2 \mathbb{E} \left\{ \|\pi_{U_i}\|^2 \right\} + \mathbb{E} \left\{ \|\pi_{V_i}\|^2 \right\}}{4}$$

From Figure 2 we deduce that if we apply the  $(U | U + V)$ -second level construction we get better performance than decoding a classical RS code with the Guruswami-Sudan decoder for low rate codes, specifically for  $R < 0.326$ .

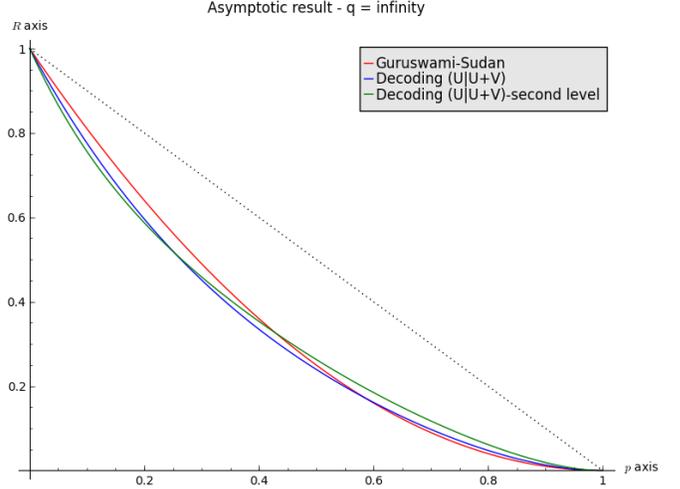


Fig. 2: Rate plotted against the crossover error probability  $p$  for several algorithms. The red line refers to the Guruswami-Sudan algorithm, the blue line to the  $(U | U + V)$ -construction and the green line to the  $(U | U + V)$ -second level construction.

### C. Runtime of the Algorithm

The runtime of Koetter-Vardy (KV) soft decoder is identical to the Sudan algorithm except for the “*soft interpolation step*” related with an optimal multiplicity matrix  $M$ . Asymptotically, for large code length, the optimal multiplicity matrix  $M$  becomes proportional to the reliability matrix  $\Pi$  of the decoder.

**Definition 3.** Given a matrix  $M = (m_{ij}) \in \mathbb{Z}_{\geq 0}^{q \times n}$ . We define its cost as:

$$\text{Cost}(M) = \frac{1}{2} \sum_{i=1}^q \sum_{j=1}^n m_{ij}(m_{ij} + 1) = \frac{1}{2} \langle M, M \rangle + \langle M, \mathbf{1} \rangle$$

where  $\mathbf{1}$  denotes the all-one matrix.

The cost of the chosen multiplicity matrix gives exactly the number of linear equations that we have to solve in the *soft interpolation step*. Thus, it is easy to check that the decoder associated with the  $(U | U + V)$ -construction has a similar runtime as the KV decoder (even lower for small values of the error rate  $p$ ). Moreover, the runtime of the  $(U | U + V)$ -second level construction is always lower than both mentioned decoders.

**Lemma 6.** The cost of the reliability matrix of the  $(U | U + V)$ -construction is:

$$\frac{\text{Cost}(\Pi \oplus \Pi) + \text{Cost}(\Pi \times \Pi)}{2} = \frac{(p^4 - 5p^3 + 8p^2 - 10p + 8)(1 - p)}{4(2 - p)}$$

## IV. A NEW MC-ELIECE SCHEME

As we have seen, this  $(U | U + V)$  construction gives codes which in the low rate regime have even better

error correction capacities than a standard RS code. This suggests to use such codes in a McEliece cryptosystem to replace the original Goppa codes. These  $(U | U + V)$  codes do not only have a better error correction capacity, they also allow to avoid the Sidelnikov-Shestakov attack [16] that broke a previous proposal based on GRS codes [13]. Furthermore we can even strengthen the security of this scheme by using instead of the  $(U | U + V)$  construction a generalized  $(U | U + V)$  code which has trivially the same error-correction capacity as the  $(U | U + V)$  construction but with better minimum distance properties which seems essential to avoid attacks based on finding minimum weight codewords in the code and trying to unravel the code structure from those minimum weight codewords. Analyzing precisely attacks of this kind needs however additional tools due to the peculiar structure of these generalized  $(U | U + V)$  codes (it is for instance inappropriate to use the analysis done for random codes) and is out of scope of this paper.

**Definition 4.** Let  $(U, V)$  be a pair of codes with parameters  $[n, k_u, d_u]_q$  and  $[n, k_v, d_v]_q$ , respectively. Consider the following matrix

$$\mathbf{D} = \left( \begin{array}{c|c} \mathbf{D}_1 & \mathbf{D}_3 \\ \hline \mathbf{D}_2 & \mathbf{D}_4 \end{array} \right) \in \mathbb{F}_q^{n \times n}$$

where the  $\mathbf{D}_i$ 's are diagonal matrices such that  $\mathbf{D}$  is non singular. We define the generalized  $(U | U + V)$ -construction of  $U$  and  $V$  with respect to  $\mathbf{D}$  as the matrix product code:

$$\{(\mathbf{u}\mathbf{D}_1 + \mathbf{v}\mathbf{D}_2 | \mathbf{u}\mathbf{D}_3 + \mathbf{v}\mathbf{D}_4) | \mathbf{u} \in U \text{ and } \mathbf{v} \in V\}.$$

It is denoted by  $[U, V] \cdot \mathbf{D}$ .

It is readily verified that this code has parameters  $[2n, k_u + k_v, d]$  with

$$\min\{2d_u, d_v\} \leq d \leq \min\{2d_u, 2d_v\}.$$

Note that the minimum distance of this generalized  $(U | U + V)$  can supersede the minimum distance of the standard  $(U | U + V)$  construction which is equal to  $\min\{2d_u, d_v\}$ .

Let  $U$  and  $V$  be codes with generator matrices  $\mathbf{G}_u$  and  $\mathbf{G}_v$ , respectively. It is a simple exercise to show that

$$\left( \begin{array}{c|c} \mathbf{G}_u\mathbf{D}_1 & \mathbf{G}_u\mathbf{D}_3 \\ \hline \mathbf{G}_v\mathbf{D}_2 & \mathbf{G}_v\mathbf{D}_4 \end{array} \right)$$

is a generator matrix for  $[U, V] \cdot \mathbf{D}$ .

These generalized  $(U | U + V)$  codes based on RS constituent codes have clearly an efficient decoding which is similar to the  $(U | U + V)$ -decoder. There are only a few differences: when we receive a word  $(\mathbf{y}_1, \mathbf{y}_2)$  we just compute the difference  $\mathbf{y}_1\mathbf{D}_3 - \mathbf{y}_2\mathbf{D}_1$  which should be a noisy version of  $\mathbf{v}(\mathbf{D}_2\mathbf{D}_3 - \mathbf{D}_4\mathbf{D}_1)$ . However the error correction capacity is the same as the original  $(U | U + V)$  with this kind of decoding algorithm. More precisely, the McEliece scheme we propose is the following

### Key generation:

- Choose  $U, V$  as RS codes of some length  $n$ .
- Construct a random matrix  $\mathbf{D}$  as described in Definition 4.
- Let  $G$  be a random generator matrix of the code  $\mathcal{C} = [U, V] \cdot \mathbf{D} \cdot \Sigma_{2n}$  where  $\Sigma_{2n}$  is a permutation matrix of size  $2n$  and  $\mathcal{A}_{\mathcal{C}}$  a decoding algorithm for  $\mathcal{C}$  that typically corrects  $t$  errors. It consists in applying  $\Sigma_{2n}^{-1}$  to the received word and then performing the aforementioned generalized  $(U | U + V)$ -decoder.

The *public key* and the *private key* are given respectively by:

$$\mathcal{K}_{\text{pub}} = (G, t) \quad \text{and} \quad \mathcal{K}_{\text{secret}} = \mathcal{A}_{\mathcal{C}}$$

**Encryption:**  $\mathbf{y} = \mathbf{m}G + \mathbf{e}$  where  $\mathbf{m}$  is the message and  $\mathbf{e}$  is a random error vector of weight at most  $t$ .

**Decryption:** Use  $\mathcal{K}_{\text{secret}}$  to retrieve  $\mathbf{m}$ .

Note that Wang proposed in [18] a very similar scheme, with the difference that  $U$  was a random code and  $V$  a RS code and that he took only a subcode of the generalized  $(U | U + V)$  code namely the code generated by  $(\mathbf{G}_u\mathbf{D}_1 + \mathbf{G}_v\mathbf{D}_2 | \mathbf{G}_u\mathbf{D}_3 + \mathbf{G}_v\mathbf{D}_4)$ . The code rate loss implied by this choice results in a significant loss in the key size (since we have to protect ourself against generic decoders for  $t$  errors for a code which is of much smaller dimension). The fact that  $U$  is random in his scheme however is a rather strong argument in favor of its security.

### CONCLUSION

This paper introduces a modification of RS codes that beats the Guruswami-Sudan decoding radius at low rates (specifically for rates  $R < 0.326$ ). Moreover, it seems to avoid Sidelnikov-Shestakov attack which makes this construction an interesting candidate to be used in Code-based Cryptography.

### REFERENCES

- [1] M. Baldi, M. Bianchi, and F. Chiaraluce. Security and complexity of the McEliece cryptosystem based on QC-LDPC codes. *IET Information Security*, 7(3):212–220, Sept. 2013.
- [2] M. Baldi, M. Bianchi, F. Chiaraluce, J. Rosenthal, and D. Schipani. Enhanced public key security for the McEliece cryptosystem. *J. Cryptology*, 2014.
- [3] A. Bennatan and D. Burshtein. Design and analysis of nonbinary LDPC codes over arbitrary discrete-memoryless channels. *IEEE Trans. Inform. Theory*, 52(2):549–583, Feb. 2006.
- [4] A. Couvreur, P. Gaborit, V. Gauthier-Umaña, A. Otmani, and J.-P. Tillich. Distinguisher-based attacks on public-key cryptosystems using Reed-Solomon codes. *Des. Codes Cryptogr.*, 73(2):641–666, 2014.
- [5] A. Couvreur, A. Otmani, J. Tillich, and V. Gauthier-Umaña. A polynomial-time attack on the BBCRS scheme. In J. Katz, editor, *Public-Key Cryptography - PKC 2015*, volume 9020 of *Lecture Notes in Comput. Sci.*, pages 175–193. Springer, 2015.
- [6] J.-C. Faugère, A. Otmani, L. Perret, F. de Portzamparc, and J.-P. Tillich. Folding alternant and Goppa Codes with non-trivial automorphism groups. *IEEE Trans. Inform. Theory*, 62(1):184–198, 2016.

- [7] V. Guruswami and A. Rudra. Error correction up to the information-theoretic limit. *Commun. ACM*, 52(3):87–95, Mar. 2009.
- [8] V. Guruswami and M. Sudan. Improved decoding of Reed-Solomon and algebraic-geometry codes. *IEEE Trans. Inform. Theory*, 45(6):1757–1767, 1999.
- [9] R. Koetter and A. Vardy. Algebraic soft-decision decoding of reed-solomon codes. *IEEE Trans. Inform. Theory*, 49(11):2809–2825, 2003.
- [10] I. Márquez-Corbella and J.-P. Tillich. Using Reed-Solomon codes in the  $(u|u + v)$  construction and an application in cryptography. preprint, 2016.
- [11] R. J. McEliece. *A Public-Key System Based on Algebraic Coding Theory*, pages 114–116. Jet Propulsion Lab, 1978.
- [12] R. Misoczki, J.-P. Tillich, N. Sendrier, and P. S. L. M. Barreto. MDPC-McEliece: New McEliece variants from moderate density parity-check codes. In *Proc. IEEE Int. Symposium Inf. Theory - ISIT*, pages 2069–2073, 2013.
- [13] H. Niederreiter. Knapsack-type cryptosystems and algebraic coding theory. *Problems of Control and Information Theory*, 15(2):159–166, 1986.
- [14] F. Parvaresh and A. Vardy. Correcting errors beyond the guruswami-sudan radius in polynomial time. In *Foundations of Computer Science, 2005. FOCS 2005. 46th Annual IEEE Symposium on*, pages 285–294, 2005.
- [15] P. W. Shor. Algorithms for quantum computation: Discrete logarithms and factoring. In S. Goldwasser, editor, *FOCS*, pages 124–134, 1994.
- [16] V. M. Sidelnikov and S. Shestakov. On the insecurity of cryptosystems based on generalized Reed-Solomon codes. *Discrete Math. Appl.*, 1(4):439–444, 1992.
- [17] M. Sudan. Decoding of Reed Solomon codes beyond the error-correction bound. *J. Complexity*, 13(1):180–193, 1997.
- [18] Y. Wang. Quantum resistant random linear code based public key encryption scheme RLCE, Dec. 2015.
- [19] C. Wieschebrink. Two NP-complete problems in coding theory with an application in code based cryptography. In *Proc. IEEE Int. Symposium Inf. Theory - ISIT*, pages 1733–1737, 2006.