

A distinguisher for high-rate McEliece Cryptosystems

J.C. Faugère (INRIA, SALSA project),
Valérie Gauthier (Math. dep. Tech. Univ. of Denmark),
A. Otmani (Université Caen- INRIA, SECRET project),
L. Perret (INRIA, SALSA project),
J.-P. Tillich (INRIA, SECRET project)

May 12th, 2011

1. (Generalized) McEliece Cryptosystem $\text{McE}(\mathcal{K}_{n,k,t})$

C a q -ary, length n , dimension k , t -error correcting code

- **Public key:** G a $k \times n$ generator matrix of C in $\mathcal{K}(n, k, t)$
- **Secret key:** Ψ a t -error correcting procedure for C
- **Encryption:** $x \rightarrow xG + e$ with e of Hamming weight t
- **Decryption:** $y \rightarrow \Psi(y)G^{-1}$ with G^{-1} a right inverse of G .

Alternant codes/Goppa codes

▶ $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{F}_{q^m}^n$ with $x_i \neq x_j$ if $i \neq j$

▶ $\mathbf{y} = (y_1, \dots, y_n) \in \mathbb{F}_{q^m}^n$ with $y_i \neq 0$

For any $r < n$, let $\mathbf{H}_r(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \begin{pmatrix} y_1 & y_2 & \cdots & y_n \\ y_1 x_1 & y_2 x_2 & \cdots & y_n x_n \\ \vdots & \vdots & & \vdots \\ y_1 x_1^{r-1} & y_2 x_2^{r-1} & \cdots & y_n x_n^{r-1} \end{pmatrix}$

Definition 1. An *alternant* code is the kernel of an \mathbf{H} of this type

$$\mathcal{A}_r(\mathbf{x}, \mathbf{y}) = \{ \mathbf{v} \in \mathbb{F}_q^n \mid \mathbf{H}_r(\mathbf{x}, \mathbf{y}) \mathbf{v}^T = \mathbf{0}. \}.$$

Goppa code : $\exists \Gamma$, polynomial of degree r such that $y_i = \Gamma(x_i)^{-1}$.

Decoding Alternant and Goppa codes

Proposition 1. [decoding alternant codes] $r/2$ errors can be decoded in polynomial time as long as x and y are *known*.

Proposition 2. [The special case of binary Goppa codes] In the case of a binary Goppa code ($q = 2$), r errors can be decoded in polynomial time, if x and Γ are known and if Γ has only simple roots.

More generally a factor $\frac{q}{q-1}$ can be gained (exploited for instance in wild McEliece [**Bernstein-Lange-Peters 2010**]) by a suitable choice of Γ .

(public key) 2. Distinguisher problem

$\mathcal{K}^{\text{Goppa}}(n, k, t)$ the ensemble of generator matrices of t -error correcting Goppa codes of length n , dimension k

$\mathcal{K}^{\text{alt}}(n, k)$ the ensemble of generator matrices of alternant codes of length n , dimension k

$\mathcal{K}^{\text{lin}}(n, k)$ the ensemble of generator matrices of linear codes of length n and dimension k .

Can we distinguish between the cases

(i) $\mathbf{G} \in \mathcal{K}^{\text{Goppa}}(n, k, t)$

(ii) $\mathbf{G} \in \mathcal{K}^{\text{alt}}(n, k)$

(iii) $\mathbf{G} \in \mathcal{K}^{\text{lin}}(n, k)$?

Niederreiter Nied($\mathcal{K}_{n,k,t}$)

C a q -ary, length n , dimension k , t -error correcting code.

- **Public key:** \mathbf{H} a $(n - k) \times n$ parity check matrix of C , $\mathbf{H} \in \mathcal{K}_{n,k,t}$
- **Secret key:** Ψ a t -error correcting procedure for C
- **Encryption:** $e \rightarrow e\mathbf{H}^T$ with e of Hamming weight t
- **Decryption:** To decipher s , choose any \mathbf{y} of syndrome s , i.e. such that $s = \mathbf{y}\mathbf{H}^T$, and output $\mathbf{y} - \Psi(\mathbf{y})$.

A probabilistic model of an attacker

A (T, ϵ) adversary \mathcal{A} for $\mathbf{Nied}(\mathcal{K}_{n,k,t})$ is a program which runs in time at most T and is such that

$$\mathbf{Prob}_{H,e}(\mathcal{A}(H, eH^T) = e | H \in \mathcal{K}_{n,k,t}) \geq \epsilon$$

Most attacks actually deal with an adversary for $\mathbf{Nied}(\mathcal{K}^{\text{lin}}(n, k))$ instead of $\mathbf{Nied}(\mathcal{K}^{\text{Goppa}}(n, k, t))$.

How the distinguisher appears

$$\mathbf{Adv} \stackrel{\text{def}}{=} \mathbf{Prob}(\mathcal{A}(\mathbf{H}, \mathbf{eH}^T) = \mathbf{e} | \mathbf{H} \in \mathcal{K}_{n,k,t}^{\text{Goppa}}) - \mathbf{Prob}(\mathcal{A}(\mathbf{H}, \mathbf{eH}^T) = \mathbf{e} | \mathbf{H} \in \mathcal{K}_{n,k}^{\text{lin}})$$

Distinguisher D :

input $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$

Step 1 : pick a random $\mathbf{e} \in \mathbb{F}_q^n$ of weight t

Step 2: if $\mathcal{A}(\mathbf{H}, \mathbf{eH}^T) = \mathbf{e}$ then return 1, else return 0.

Advantage of $D = |\mathbf{Adv}|$.

Either a decoding algorithm on linear codes or a distinguisher for Goppa codes

Proposition 3. *If $\exists(T, \epsilon)$ -adversary against $\mathbf{Nied}(\mathcal{K}_{n,k,t}^{\text{Goppa}})$, then there exists either*

- (i) *a $(T, \epsilon/2)$ -adversary against $\mathbf{Nied}(\mathcal{K}^{\text{lin}}(n, k))$ (i.e. a **decoder** for general linear codes working in time T with success probability at $\geq \epsilon/2$).*
- (ii) *A distinguisher between $\mathbf{H} \in \mathcal{K}_{n,k,t}^{\text{Goppa}}$ and $\mathbf{H} \in \mathcal{K}_{n,k}^{\text{lin}}$ working in time $T + O(n^2)$ and with advantage at least $\epsilon/2$.*

3. Algebraic approach for attacking the McEliece cryptosystem

What is known: a basis of the code \rightarrow rows of a generator matrix $\mathbf{G} = (g_{ij})$ of size $k \times n$.

What we also know: $\exists \mathbf{x}, \mathbf{y} \in \mathbb{F}_{q^m}^n$ s.t.

$$\mathbf{H}_r(\mathbf{x}, \mathbf{y})\mathbf{G}^T = \mathbf{0}. \quad (1)$$

What we want to find: find in the case of an alternant code \mathbf{x}, \mathbf{y} , and in the special case of a binary Goppa code \mathbf{x} and Γ .

The algebraic system

$\mathbf{H}_r(\mathbf{x}, \mathbf{y})\mathbf{G}^T = \mathbf{0}$ translates to

$$\left\{ \begin{array}{l} g_{1,1}Y_1 + \cdots + g_{1,n}Y_n = 0 \\ \vdots \\ g_{k,1}Y_1 + \cdots + g_{k,n}Y_n = 0 \\ g_{1,1}Y_1X_1 + \cdots + g_{1,n}Y_nX_n = 0 \\ \vdots \\ g_{k,1}Y_1X_1 + \cdots + g_{k,n}Y_nX_n = 0 \\ \vdots \\ g_{1,1}Y_1X_1^{r-1} + \cdots + g_{1,n}Y_nX_n^{r-1} = 0 \\ \vdots \\ g_{k,1}Y_1X_1^{r-1} + \cdots + g_{k,n}Y_nX_n^{r-1} = 0 \end{array} \right. \quad (2)$$

where the $g_{i,j}$'s are known coefficients in \mathbb{F}_q and $k \geq n - r m$.

Freedom of choice in (2)

Proposition 4. *Theoretically, the system has $2n$ unknowns but we can take arbitrary values for one Y_i and for three X_i 's (as long as these values are different).*

Applications

When the number of unknowns is small, ex:

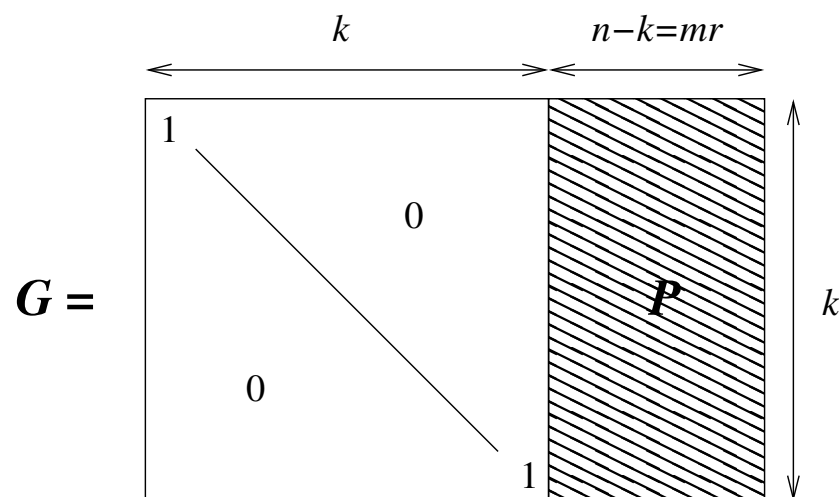
- Berger-Cayrel-Gaborit-Otmani proposal at AfricaCrypt'09 based on **quasi-cyclic alternant** codes
- Misoczki-Barreto at SAC'09 variant based on **quasi-dyadic Goppa** codes

⇒ algebraic system can be solved by (dedicated) Grobner basis techniques.

- ▶ breaks all parameters proposed in these articles ([Faugère-Otmani-Perret-Tillich;Eurocrypt 2010] with the exception of binary dyadic codes. Related to [Leander-Gauthier Umana; SCC2010])

4. A naive attack

W.l.o.g. we can assume that G is systematic in its k first positions.



Step 1 – expressing the $Y_i X_i^d$'s in terms of the $Y_j X_j^d$'s for $j \in \{k + 1, \dots, n\}$.

$\mathbf{P} = (p_{ij})_{\substack{1 \leq i \leq k \\ k+1 \leq j \leq n}}$. We can rewrite (2) as

$$\left\{ \begin{array}{l} Y_i \\ Y_i X_i \\ \dots \\ Y_i X_i^{r-1} \end{array} \right. = \begin{array}{l} \sum_{j=k+1}^n p_{i,j} Y_j \\ \sum_{j=k+1}^n p_{i,j} Y_j X_j \\ \dots \\ \sum_{j=k+1}^n p_{i,j} Y_j X_j^{r-1} \end{array} \quad (3)$$

for all $i \in \{1, \dots, k\}$.

Step 2.– Exploiting $Y_i(Y_iX_i^2) = (Y_iX_i)^2$

$$\begin{cases} Y_i & = \sum_{j=k+1}^n p_{i,j} Y_j \\ Y_i X_i & = \sum_{j=k+1}^n p_{i,j} Y_j X_j \\ Y_i X_i^2 & = \sum_{j=k+1}^n p_{i,j} Y_j X_j^2 \end{cases} \quad (4)$$

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

$$\Rightarrow \sum_{j=k+1}^n \sum_{j'>j} p_{i,j} p_{i,j'} (Y_j Y_{j'} X_{j'}^2 + Y_{j'} Y_j X_j^2) = 0$$

Step 3. – Linearization

$$Z_{jj'} \stackrel{\text{def}}{=} Y_j Y_{j'} X_{j'}^2 + Y_{j'} Y_j X_j^2$$

$$\sum_{j=k+1}^n \sum_{j'>j} p_{i,j} p_{i,j'} Z_{jj'} = 0.$$

▶ $\binom{n-k}{2} \approx \frac{m^2 r^2}{2}$ unknowns

▶ $k = n - mr$ equations

⇒ reveals $Z_{jj'}$ when $n - mr \geq \frac{m^2 r^2}{2}$?

▶ This happens for the Courtois-Finiasz-Sendrier scheme, ex: $n = 2^{21}$, $r = 10$, $m = 21$ which **has** to choose small values of r .

Linearized System

Definition 2. Assume that the public key G of the McEliece cryptosystem is in systematic form $(\mathbf{I}_k \mid \mathbf{P})$

The **linearized system** associated to G is

$$\left\{ \begin{array}{l} \sum_{j=k+1}^n \sum_{j'>j} p_{1,j} p_{1,j'} Z_{jj'} = 0 \\ \sum_{j=k+1}^n \sum_{j'>j} p_{2,j} p_{2,j'} Z_{jj'} = 0 \\ \vdots \\ \sum_{j=k+1}^n \sum_{j'>j} p_{k,j} p_{k,j'} Z_{jj'} = 0 \end{array} \right.$$

The **dimension** of the solution space is denoted by D .

Algebraic Distinguisher

Solving this system requires that

- Number of equations k is greater than the number of unknowns $\binom{n-k}{2}$
- rank is (almost) equal to the number of unknowns

If G is random then one **would expect** that the rank is $\min \left\{ k, \binom{n-k}{2} \right\}$

$$\implies D = \max \left\{ 0, \binom{n-k}{2} - k \right\}$$

But for **several structured** (Goppa, alternant) codes $\text{rank} < \min \left\{ k, \binom{n-k}{2} \right\}$
and this defect can be **quantified**

Example $q = 2$ and $m = 14$

r	3	4	5	6	7	8	9	10	11	12	13	14
$\binom{n-k}{2}$	861	1540	2415	3486	4753	6216	7875	9730	11781	14028	16471	19110
k	16342	16328	16314	16300	16286	16272	16258	16244	16230	16216	16202	16188
D_{rand}	0	0	0	0	0	0	0	0	0	0	269	292
$D_{\text{alternant}}$	42	126	308	560	882	1274	1848	2520	3290	4158	5124	6188
D_{Goppa}	252	532	980	1554	2254	3080	4158	5390	6776	8316	10010	11858

Example $q = 2$ and $m = 14$

r	15	16	17	18	19	20	21	22	23	24	25	26	27
$\binom{n-k}{2}$	21945	24976	28203	31626	35245	39060	43071	47278	51681	56280	61075	66066	71253
k	16174	16160	16146	16132	16118	16104	16090	16076	16062	16048	16034	16020	16006
D_{rand}	5771	8816	12057	15494	19127	22956	26981	31202	35619	40232	45041	50046	55247
$D_{\text{alternant}}$	7350	8816	12057	15494	19127	22956	26981	31202	35619	40232	45041	50046	55247
D_{Goppa}	13860	16016	18564	21294	24206	27300	30576	34034	37674	41496	45500	50046	55247

Alternant Case

Let $\ell \stackrel{\text{def}}{=} \lfloor \log_q(r-1) \rfloor$.

$$D_{\text{alternant}} = \frac{1}{2}m(r-1) \left((2\ell+1)r - 2\frac{q^{\ell+1}-1}{q-1} \right)$$

as long as $\binom{n-k}{2} - D_{\text{alternant}} < k$.

Goppa Case

Let ℓ the unique integer such that $q^\ell - 2q^{\ell-1} + q^{\ell-2} < r \leq q^{\ell+1} - 2q^\ell + q^{\ell-1}$

$$D_{\text{Goppa}} = \begin{cases} \frac{1}{2}m(r-1)(r-2) = D_{\text{alternant}} & \text{for } r < q-1 \\ \frac{1}{2}mr \left((2\ell+1)r - 2q^\ell + 2q^{\ell-1} - 1 \right) & \text{for } r \geq q-1 \end{cases}$$

as long as $\binom{n-k}{2} - D_{\text{Goppa}} < k$.

Example $q = 2$ and $m = 14$

r	3	4	5	6	7	8	9	10	11	12	13	14
$\binom{n-k}{2}$	861	1540	2415	3486	4753	6216	7875	9730	11781	14028	16471	19110
k	16342	16328	16314	16300	16286	16272	16258	16244	16230	16216	16202	16188
D_{rand}	0	0	0	0	0	0	0	0	0	0	269	2922
$D_{\text{alternant}}$	42	126	308	560	882	1274	1848	2520	3290	4158	5124	6188
$T_{\text{alternant}}$	42	126	308	560	882	1274	1848	2520	3290	4158	5124	6188
D_{Goppa}	252	532	980	1554	2254	3080	4158	5390	6776	8316	10010	11858
T_{Goppa}	252	532	980	1554	2254	3080	4158	5390	6776	8316	10010	11858

Example $q = 2$ and $m = 14$

r	15	16	17	18	19	20	21	22	23	24	25	26	27
$\binom{n-k}{2}$	21945	24976	28203	31626	35245	39060	43071	47278	51681	56280	61075	66066	71253
k	16174	16160	16146	16132	16118	16104	16090	16076	16062	16048	16034	16020	16006
D_{rand}	5771	8816	12057	15494	19127	22956	26981	31202	35619	40232	45041	50046	55247
$D_{\text{alternant}}$	7350	8816	12057	15494	19127	22956	26981	31202	35619	40232	45041	50046	55247
$T_{\text{alternant}}$	7350	8610	10192	11900	13734	15694	17780	19992	22330	24794	27384	30100	32942
D_{Goppa}	13860	16016	18564	21294	24206	27300	30576	34034	37674	41496	45500	50046	55247
T_{Goppa}	13860	16016	18564	21294	24206	27300	30576	34034	37674	41496	45500	49686	54054

Simplified Formulas for binary Goppa Codes

- ▶ Let $\ell \stackrel{\text{def}}{=} \lceil \log_2 r \rceil + 1$.

$$D_{\text{Goppa}} = \frac{1}{2}mr \left((2\ell + 1)r - 2^\ell - 1 \right)$$

as long as $\binom{mr}{2} - D_{\text{Goppa}} < n - mr$.

Binary Goppa Codes

In particular, assuming that $n = 2^m$, the binary Goppa code distinguishing problem is **solved** for any $r < r_{\max}$

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

- ▶ $m = 13$ and $r = 19$ corresponds to a 90-bit security McEliece public key.
- ▶ All CFS parameters fits in the range of validity of the algebraic distinguisher.

5. Explanation

- ▶ Formulas obtained through experimentations for random codes, alternant codes and **irreducible** Goppa codes over fields of size $q \in \{2, 4, 8, 16\}$.
- ▶ We have an explanation for alternant codes and binary Goppa codes by **guessing** a basis of the solution vector space over \mathbb{F}_q .
- ▶ It **does not** provide a proof.

Explanation for Alternant Codes – Step I

- ▶ Note that the **entries** of the system are in \mathbb{F}_q and **solutions** are sought in \mathbb{F}_q^m .
 - ▶ Let us view \mathbb{F}_q^m as a \mathbb{F}_q -vector space of dimension m , and let $\pi_i : \mathbb{F}_q^m \rightarrow \mathbb{F}_q$ be the function giving the i -th coordinate.
 - ▶ Hence, if a vector \mathbf{v} with $v_j \in \mathbb{F}_q^m$ is a solution then $\pi_i(\mathbf{v}) = \left(\pi_i(v_j) \right)_j$ whose entries are in \mathbb{F}_q is **also** a solution.
- \implies Any solution with entries over \mathbb{F}_q^m would **potentially** provide a basis of m solutions with entries over \mathbb{F}_q

Explanation for Alternant Codes – Step II

- ▶ We have used $Y_i Y_i X_i^2 = (Y_i X_i)^2$ which leads to:

$$\forall i \in \{1, \dots, k\}, \quad \sum_{j=k+1}^n \sum_{j'>j} p_{i,j} p_{i,j'} Y_j Y_{j'} (X_j^2 + X_{j'}^2) = 0$$

- ▶ But we can use any relation $Y_i X_i^a Y_i X_i^b = Y_i X_i^c Y_i X_i^d$ with a, b, c, d in $\{0, \dots, r-1\}$ such that $a + b = c + d$

$$\sum_{j=k+1}^n \sum_{j'>j} p_{i,j} p_{i,j'} Y_j Y_{j'} (X_j^a X_{j'}^b + X_j^b X_{j'}^a + X_j^c X_{j'}^d + X_j^d X_{j'}^c) = 0$$

Explanation for Alternant Codes – Step III

- ▶ For $r \geq q$, the automorphism $x \mapsto x^{q^\ell}$ for any $0 \leq \ell \leq m - 1$ can be used.

$$\forall e \in \{0, \dots, r - 1\}, \quad Y_i X_i^e = \sum_{j=k+1}^n p_{ij} Y_j X_j^e \implies Y_i^q X_i^{eq} = \sum_{j=k+1}^n p_{ij} Y_j^q X_j^{eq}$$

- ▶ We therefore can use the same trick, for instance $Y_i (Y_i X_i)^q = Y_i^q Y_i X_i^q$,

$$\sum_{j=k+1}^n \sum_{j' > j} p_{i,j} p_{i,j'} \left(Y_j Y_{j'}^q X_{j'}^q + Y_{j'} Y_j^q X_j^q + Y_j^q Y_{j'} X_{j'}^q + Y_{j'}^q Y_j X_j^q \right) = 0.$$

Explanation for Alternant Codes

- ▶ However the equations obtained $(Y_i X_i^a Y_i X_i^b)^q = (Y_i X_i^c Y_i X_i^d)^q$ do not provide new solutions after decomposition over \mathbb{F}_q since they are linearly dependent of those obtained from $Y_i X_i^a Y_i X_i^b = Y_i X_i^c Y_i X_i^d$.
- ▶ Hence, we only consider equations obtained from integers a, b, c, d, ℓ such that $a + bq^\ell = c + dq^\ell$

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

$$\mathbf{Z}_{a,b,c,d,\ell} \stackrel{\text{def}}{=} \left(Y_j X_j^a Y_{j'}^{q^\ell} X_{j'}^{bq^\ell} + Y_{j'} X_{j'}^a Y_j^{q^\ell} X_j^{bq^\ell} + Y_j X_j^c Y_{j'}^{q^\ell} X_{j'}^{dq^\ell} + Y_{j'} X_{j'}^c Y_j^{q^\ell} X_j^{dq^\ell} \right)_{1 \leq j < j' \leq n-k}$$

Explanation for Alternant Codes

- ▶ Let us assume that $d > b$ and set $\delta \stackrel{\text{def}}{=} d - b$ and then $a = c + q^\ell \delta$

$$\implies \mathbf{Z}_{a,b,c,d,\ell} = \mathbf{Z}_{c+q^\ell\delta,b,c,b+\delta,\ell}$$

- ▶ Let \mathcal{B}_r be the set $\mathbf{Z}_{c+q^\ell\delta,b,c,b+\delta,\ell}$ obtained with $\delta = 1$ and satisfying:

$$\begin{cases} 0 \leq b \leq r - 2 \text{ and } 0 \leq c \leq r - 1 - q^\ell & \text{if } 1 \leq \ell \leq \lfloor \log_q(r - 1) \rfloor \\ 0 \leq b < c \leq r - 2 & \text{if } \ell = 0. \end{cases}$$

Proposition 5. • Any $\mathbf{Z}_{c+q^\ell\delta,b,c,b+\delta,\ell}$ belongs to the \mathbb{F}_{q^m} -vector space generated by \mathcal{B}_r

- The cardinality of \mathcal{B}_r with $r \geq 3$ is equal to D/m .

Heuristic

For random choices of x_i 's and y_i 's defining the alternant code, the set $\left\{ \pi_i(\mathbf{Z}) \mid \mathbf{Z} \in \mathcal{B}_r \text{ and } 1 \leq i \leq m \right\}$ forms a **basis** of the vector space that is solution to the linearized system.

Conclusion

- ▶ Large dimension comes from **the many different ways** of combining the equations together yielding the same linearized system
- ▶ What happens for random generator is **proven** now.
- ▶ Binary Goppa codes can also be explained but **no** explanation for **non-binary** Goppa codes.
- ▶ The most difficult task is **identifying** a basis of the vector space of solutions.
- ▶ A slightly better distinguisher can be obtained by taking the subcode of codewords of even weights.
- ▶ Distinguisher \Rightarrow attack ?
- ▶ Approach requires $\frac{k}{n}$ very close to 1. Should very high rates be avoided in a McEliece like scheme ?