LECTURE 1
INTRODUCTION TO CODE-BASED CRYPTOGRAPHY
DECODING A RANDOM CODE

Summer School: *Introduction to Quantum-Safe Cryptography*

Thomas Debris-Alazard

July 01, 2024

Inria, École Polytechnique

- Maxime Bombar (Post-doc at CWI, Netherland)

  maxime.bombar@cwi.nl

- Thomas Debris-Alazard (Researcher at Inria, France)

  thomas.debris@inria.fr

**Course Content:**

1. An Intractable Problem Related to Codes, Decoding

2. Random Codes

3. Information Set Decoding (ISD) Algorithms and Duals Attacks

4. Duality, Fourier Theory and Decoding Self-Reducibility (Worst-to-Average Case Reduction)

5. McEliece and Alekhnovitch Encryption' Schemes (From Original Propositions to Instantiations)
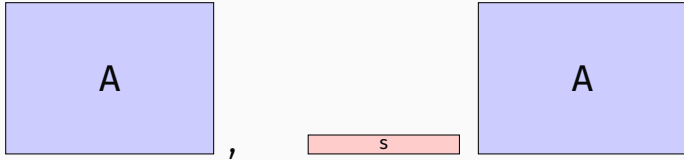
$\longrightarrow$ 3 lectures notes (long, for further reading): https://arxiv.org/pdf/2304.03541
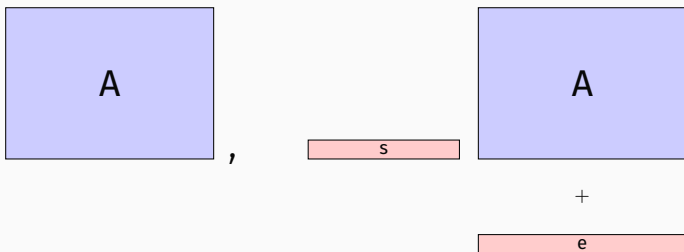
**Exercise Sessions:**

1. Starting Exercises to Get Familiar with Linear Codes & Crypto

2. Programming Session: Implement Basic ISDs and Breaking Challenges

3. Advanced Exercises About Code-Based Cryptography and Duality

$\longrightarrow$ 2 long exercise sheets: cryptanalyses of code-based encryption schemes

# Code-Based Cryptography?

A , s A

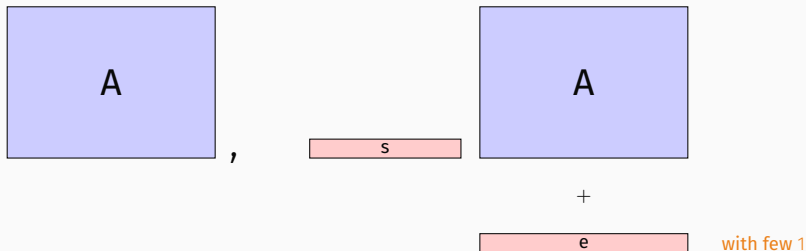Shannon (1948/1949) introduced the following problem (decoding),



**Aim:**

Recover    s

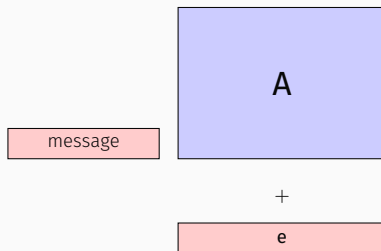Shannon (1948/1949) introduced the following problem (decoding),



$+$

$\mathbf{e}$     with few 1

Aim:

Recover $\mathbf{s}$

$\longrightarrow$ Matrix $\mathbf{A}$ and vectors $\mathbf{s}$, $\mathbf{e}$ are binary ($\in \mathbb{F}_2$)
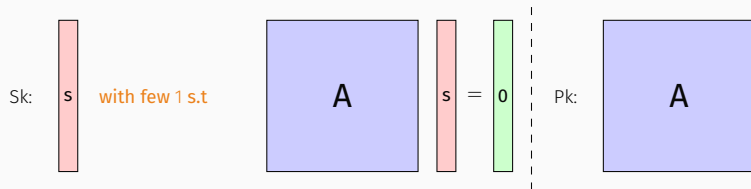
4

McEliece (1978):

$$A \leftarrow \text{Trapdoor}(): \text{public-key}$$



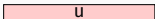- With the trapdoor: easy to recover message if **e** "short" (with few 1, a lot of 0),

- Without: hard

Alekhnovich (2003):



Sk: $s$ with few 1 s.t $\quad A \cdot s = 0 \qquad$ Pk: $A$

- To encrypt $b = 1$, send $\boxed{u} \quad \longleftarrow$ Unif

- To encrypt $b = 0$, send

$$\boxed{m}$$
$$A$$
$$+$$
$$\boxed{e} \quad \text{with few 1}$$

But how to decrypt?

$$\boxed{s} \cdot \boxed{A} : \boxed{s} \in \mathbb{F}_2^k$$

is known as a linear code!

Understanding what is a linear code: useful to

1. build trapdoors

2. understand the hardness of decoding

*The first purpose of linear codes was not cryptography. . .*

*It was telecommunication!*

$\longrightarrow$ Codes are at the core of information theory (and friends)

How to transmit $k$ bits over a noisy channel?

How to transmit $k$ bits over a noisy channel?

1. Fix $\mathcal{C}$ subspace $\subseteq \mathbb{F}_2^n$ of dimension $k$

2. Map $(m_1, \ldots, m_k) \longrightarrow \mathbf{c} = (c_1, \ldots, c_n) \in \mathcal{C}$ $\left(\text{adding } n - k \text{ bits redundancy}\right)$

3. Send $\mathbf{c}$ across the noisy channel



$\longrightarrow$ from $\mathbf{c} \oplus \mathbf{e}$: how to recover $\mathbf{e}$ and then $\mathbf{c}$?

$\left(\text{Decoding Problem}\right)$

Real life scenario: $\mathbf{c} + \mathbf{e}$ with $\mathbf{e} = (e_1, \ldots, e_n)$ *such that,*

$$\forall i \in [1, n], \quad \mathbb{P}(e_i = 1) = p \ \text{ and } \ \mathbb{P}(e_i = 0) = 1 - p$$

$\longrightarrow$ Each bit of $\mathbf{c}$ is flipped with probability $p$

Given a received corrupted word y:

$$\mathbb{P}\left(\mathbf{c} \text{ was sent} \mid \mathbf{y} \text{ is received}\right) = p^{d_H(\mathbf{c},\mathbf{y})}(1 - p)^{n - d_H(\mathbf{c},\mathbf{y})}$$

where $d_H(\mathbf{c}, \mathbf{y}) \overset{\text{def}}{=} \sharp \{i \in [1, n] \ : \ c_i \neq y_i\}$ (Hamming distance)

Real life scenario: $\mathbf{c} + \mathbf{e}$ with $\mathbf{e} = (e_1, \ldots, e_n)$ *such that,*

$$\forall i \in [1, n], \quad \mathbb{P}(e_i = 1) = p \text{ and } \mathbb{P}(e_i = 0) = 1 - p$$

$\longrightarrow$ Each bit of $\mathbf{c}$ is flipped with probability $p$

Given a received corrupted word y:

$$\mathbb{P}\left(\mathbf{c} \text{ was sent } | \mathbf{y} \text{ is received}\right) = p^{d_H(\mathbf{c},\mathbf{y})}(1 - p)^{n - d_H(\mathbf{c},\mathbf{y})}$$

where $d_H(\mathbf{c}, \mathbf{y}) \stackrel{\text{def}}{=} \sharp \{i \in [1, n] : c_i \neq y_i\}$ (Hamming distance)

Any decoding candidate $\mathbf{c} \in \mathcal{C}$ is even more likely

as it is close to the received message y for the Hamming distance.

$\longrightarrow$ It explains why historically the Hamming distance has been the considered metric

when dealing with codes. . .

# BASICS ON LINEAR CODES

$\mathbb{F}_q$: finite field with $q$ elements

**Linear Code:**

A linear code $\mathcal{C}$ of length $n$ and dimension $k$ $\left([n, k]_q\text{-code}\right)$:

subspace of $\mathbb{F}_q^n$ of dimension $k$

**First Examples:**

1. $\left\{ (f(x_1), \ldots, f(x_n)) : f \in \mathbb{F}_q[X] \text{ and } \deg(f) < k \right\}$ where the $x_i$'s are distinct elements of $\mathbb{F}_q$

   is an $[n, k]_q$-code

2. $\left\{ (\mathbf{u}, \mathbf{u} + \mathbf{v}) : \mathbf{u} \in U \text{ and } \mathbf{v} \in V \right\}$ where $U$ (*resp.* $V$) is an $[n, k_U]_q$-code (*resp.* $[n, k_V]_q$-code)

   is an $[2n, k_U + k_V]_q$-code

**Hamming Weight:**

Given $\mathbf{x} \in \mathbb{F}_q^n$, its Hamming weight is:

$$|\mathbf{x}| \overset{\text{def}}{=} \sharp \left\{ i \in [1, n] : x_i \neq 0 \right\}$$

**Minimum Distance:**

The minimum distance of $\mathcal{C}$ is:

$$d_{\min}(\mathcal{C}) \overset{\text{def}}{=} \min \left\{ |\mathbf{c}| : \mathbf{c} \in \mathcal{C}, \mathbf{c} \neq \mathbf{0} \right\}$$

$d_{\min}(\mathcal{C})$ is an important quantity:

"geometry" of $\mathcal{C}$ ; "efficiency" of $\mathcal{C}$ ; "security" of $\mathcal{C}$

$\mathcal{C}$ be an $[n, k]_q$-code

Basis representation: $\mathbf{g}_1, \ldots, \mathbf{g}_k$ basis of $\mathcal{C}$,

$$\mathcal{C} = \left\{ \mathbf{m}\mathbf{G} : \ \mathbf{m} \in \mathbb{F}_q^k \right\} \text{ where the rows of } \mathbf{G} \in \mathbb{F}_q^{k \times n} \text{ are the } \mathbf{g}_i$$

Reciprocally, any $\mathbf{G} \in \mathbb{F}_q^{k \times n}$ of rank $k$ defines the $[n, k]_q$-code,

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

Generator Matrix:

$\mathbf{G}$ is called a generator matrix

Dual Code:

Given $\mathcal{C}$, its dual $\mathcal{C}^\perp$ is the $[n, n-k]_q$-code,

$$\mathcal{C}^\perp \stackrel{\text{def}}{=} \left\{ \mathbf{c}^\perp \in \mathbb{F}_q^n : \ \forall \mathbf{c} \in \mathcal{C}, \ \mathbf{c} \cdot \mathbf{c}^\perp \stackrel{\text{def}}{=} \sum_{i=1}^n c_i \, c_i^\perp = 0 \in \mathbb{F}_q \right\}$$

$\longrightarrow$ Wait Lecture 4 to understand the rational behind this definition!

**Dual Code:**

Given $\mathcal{C}$, its dual $\mathcal{C}^\perp$ is the $[n, n-k]_q$-code,

$$\mathcal{C}^\perp \stackrel{\text{def}}{=} \left\{ \mathbf{c}^\perp \in \mathbb{F}_q^n : \forall \mathbf{c} \in \mathcal{C}, \ \mathbf{c} \cdot \mathbf{c}^\perp \stackrel{\text{def}}{=} \sum_{i=1}^n c_i \, c_i^\perp = 0 \in \mathbb{F}_q \right\}$$

$\longrightarrow$ Wait Lecture 4 to understand the rational behind this definition!

**Parity-check representation:** $\mathbf{h}_1, \ldots, \mathbf{h}_{n-k}$ basis of $\mathcal{C}^\perp$,

$$\mathcal{C} = \left\{ \mathbf{c} \in \mathbb{F}_q^n : \mathbf{H}\mathbf{c}^\mathsf{T} = \mathbf{0} \right\} \text{ where the rows of } \mathbf{H} \in \mathbb{F}_q^{(n-k) \times n} \text{ are the } \mathbf{h}_i$$

Reciprocally, any $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$ of rank $n-k$ defines the $[n, k]_q$-code,

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

**Parity-Check Matrix:**

$\mathbf{H}$ is called a parity-check matrix

- $G \in \mathbb{F}_q^{k \times n}$ generator matrix of $\mathcal{C}$ $\left( i.e., \mathcal{C} = \left\{ mG : \ m \in \mathbb{F}_q^k \right\} \right)$, $S \in \mathbb{F}_q^{k \times k}$ non-singular,

$$\longrightarrow SG \text{ still generator matrix of } \mathcal{C}$$

- $H \in \mathbb{F}_q^{(n-k) \times n}$ parity-check matrix of $\mathcal{C}$ $\left( i.e., \mathcal{C} = \left\{ c \in \mathbb{F}_q^n : \ Hc^\mathsf{T} = 0 \right\} \right)$, $S \in \mathbb{F}_q^{(n-k) \times (n-k)}$ non-singular,

$$\longrightarrow SH \text{ still parity-check matrix of } \mathcal{C}$$

$\mathbf{G} \in \mathbb{F}_q^{k \times n}$ generator matrix $\xleftrightarrow{\text{easy to compute?}}$ $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$ parity-check matrix

$G \in \mathbb{F}_q^{k \times n}$ generator matrix $\xleftrightarrow{\text{easy to compute?}}$ $H \in \mathbb{F}_q^{(n-k) \times n}$ parity-check matrix

Yes!

1. Show that if $H \in \mathbb{F}_q^{(n-k) \times n}$ has rank $n-k$ and $GH^\mathsf{T} = 0$, then $H$ parity-check (exercise)

2. Perform a Gaussian elimination: $SG = \left(I_k \mid A\right)$, then $H = \left(-A^\top \mid I_{n-k}\right)$ is a parity-check matrix

*Would you rather choose generator or parity-check representation?*

*Would you rather choose generator or parity-check representation?*

Sorry for the team generator matrix :(

Usually, the parity-check representation is more convenient

Let $\mathcal{C}_{\text{Ham}}$ be the $[7, 4]_2$-code of generator matrix:

$$\mathbf{G} \overset{\text{def}}{=} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{pmatrix}$$

$$\mathbf{H} \overset{\text{def}}{=} \begin{pmatrix} 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{pmatrix}$$

has rank 3 and verifies $\mathbf{G}\mathbf{H}^{\mathsf{T}} = \mathbf{0}$.

Let $\mathbf{c} + \mathbf{e}$ where $\begin{cases} \mathbf{c} \in \mathcal{C}_{\text{Ham}} \\ |\mathbf{e}| = 1 \end{cases}$ : how to easily recover $\mathbf{e}$?

Given $\mathbf{c} + \mathbf{e}$: recover $\mathbf{e}$

$\longrightarrow$ Make modulo $\mathcal{C}$ to extract the information about $\mathbf{e}$

Coset Space: $\mathbb{F}_q^n / \mathcal{C}$

Given an $[n, k]_q$-code $\mathcal{C}$, $\quad \sharp \, \mathbb{F}_q^n / \mathcal{C} = q^{n-k}$ and $\mathbb{F}_q^n / \mathcal{C} = \left\{ \mathbf{x}_i + \mathcal{C} \; : \; 1 \leq i \leq q^{n-k} \right\}$

A natural set of representatives via a parity-check $\mathbf{H}$: syndromes

$\mathbf{x}_i + \mathcal{C} \in \mathbb{F}_q^n / \mathcal{C} \longmapsto \mathbf{H} x_i^{\mathsf{T}} \in \mathbb{F}_q^{n-k}$ (called a syndrome)

is an isomorphism

$\mathcal{C}$ be an $[n, k]_q$-code of parity-check matrix $H$

| Noisy codeword | Syndrome |
|:---:|:---:|
| $c + e$ | $He^{\mathsf{T}}$ |

- From $c + e$: $\quad H(c + e)^{\mathsf{T}} = Hc^{\mathsf{T}} + He^{\mathsf{T}} = He^{\mathsf{T}}$

- From $He^{\mathsf{T}}$: compute with linear algebra $y$ s.t
$$H y^{\mathsf{T}} = H e^{\mathsf{T}} \iff H(y - e)^{\mathsf{T}} = 0 \iff y - e \in \mathcal{C} \iff y = c + e$$

# THE WORST-CASE DECODING PROBLEM

*Two* formulations for the worst-case decoding:

Problem (Noisy Codeword Decoding):

- Given: $G \in \mathbb{F}_q^{k \times n}$ of rank $k$, $t \in [0, n]$, $y \in \mathbb{F}_q^n$ where $y = c + e$ with $c = mG$ for some $m \in \mathbb{F}_q^k$ and $|e| = t$

- Find: $e$ $\left(\text{or equivalently } m\right)$

Problem (Syndrome Decoding):

- Given: $H \in \mathbb{F}_q^{(n-k) \times n}$ of rank $n - k$, $t \in [0, n]$, $s \in \mathbb{F}_q^{n-k}$ where $He^T = s^T$ with $|e| = t$

- Find: $e$

$\longrightarrow$ These problems are equivalent!

$n$ length   ;   $k$ dimension   ;   $t$ decoding distance

Let, $\mathcal{A}$ be an algorithm such that $\mathcal{A}(\mathbf{G}, \mathbf{mG} + \mathbf{e}) \longmapsto \mathbf{e}$

Given $(\mathbf{H}, \mathbf{He}^{\mathsf{T}})$: our aim, recover $\mathbf{e}$ using $\mathcal{A}$

1. Compute with linear algebra $\mathbf{G}$ (rank $k$) such that $\mathbf{GH}^{\mathsf{T}} = \mathbf{0}$

2. Compute (again) with linear algebra $\mathbf{y}$ such that $\mathbf{Hy}^{\mathsf{T}} = \mathbf{He}^{\mathsf{T}}$

3. Notice that $\mathbf{H}(\mathbf{y} - \mathbf{e})^{\mathsf{T}} = \mathbf{0} \iff \mathbf{y} - \mathbf{e} = \mathbf{mG}$ for some $\mathbf{m} \in \mathbb{F}_q^k$

4. Feed $(\mathbf{G}, \mathbf{y})$ to $\mathcal{A}$: it recovers $\mathbf{e}$

Exercise: show that the reciprocal holds

In what follows, we will mainly keep the parity-check representation!

Worst-Case Decisional Decoding Problem

- Input: $\mathsf{H} \in \mathbb{F}_q^{(n-k) \times n}$, $\mathsf{s} \in \mathbb{F}_q^{n-k}$ where $n, k \in \mathbb{N}$ with $k \leq n$ and an integer $t \leq n$.

- Decision: it exists $\mathsf{e} \in \mathbb{F}_q^n$ of Hamming weight $t$ such $\mathsf{He}^{\mathsf{T}} = \mathsf{s}^{\mathsf{T}}$?

This problem is NP-complete

*Is it useful?*

Be careful of the input set!

The above NP-completeness shows that (if $P \neq NP$)

We cannot easily solve the decoding problem for all codes and all decoding distances...

$\longrightarrow$ There are codes for which decoding is hard!

Not a safety guarantee for cryptographic applications!

*Is decoding hard for all codes?*

No! (remember Hamming code...)

**Generalized Reed-Solomon (GRS) Codes:**

Given $z \in (\mathbb{F}_q^\star)^n$ and $x \in \mathbb{F}_q^n$ s.t $x_i \neq x_j$ (in particular $n \leq q$) and $k \leq n$.

The code $\text{GRS}_k(x, z)$ is defined as:

$$\text{GRS}_k(x, z) \stackrel{\text{def}}{=} \left\{ \Big( z_1 f(x_1), \ \ldots, \ z_n f(x_n) \Big) \, : \, f \in \mathbb{F}_q[X] \ \text{and} \ \ \deg(f) < k \right\}$$

$\longrightarrow$ GRS are used in QR-codes!

**Exercise**: $\text{GRS}_k(x, z)$ has generator matrix:

$$G \stackrel{\text{def}}{=} \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \\ x_1^2 & x_2^2 & \cdots & x_n^2 \\ \vdots & \vdots & \vdots & \vdots \\ x_1^k & x_2^k & \cdots & x_n^k \end{pmatrix} \begin{pmatrix} z_1 & & & 0 \\ & z_2 & & \\ & & \ddots & \\ 0 & & & z_n \end{pmatrix}$$

Decoding Algorithm:

Given, $\mathrm{GRS}_k(\mathbf{x}, \mathbf{z})$ and $\mathbf{c} + \mathbf{e}$ such that $\begin{cases} \mathbf{c} \in \mathrm{GRS}_k(\mathbf{x}, \mathbf{z}) \\ |\mathbf{e}| \leq \left\lfloor \frac{n-k}{2} \right\rfloor \end{cases}$

Then, we can recover $(\mathbf{c}, \mathbf{e})$ in polynomial time in the size of inputs, *i.e.,* $O\left(n^\ell\right)$ for some $\ell$.

$\longrightarrow$ See Exercise Session

- There are codes for which decoding is hard (NP-Completeness)

- Decoding is easy for some family of codes (for instance Generalized-Reed-Solomon codes)

Is decoding hard for almost all codes?

# AVERAGE DECODING PROBLEM

DP$(n, q, R, \tau)$, $k \stackrel{\text{def}}{=} Rn$ and $t \stackrel{\text{def}}{=} \tau n$

*Sample:* $\boxed{H} \longleftarrow \text{Unif}\left(\mathbb{F}_q^{(n-k) \times n}\right)$, $\boxed{x} \longleftarrow \text{Unif}\left(z \ : \ |z| = t\right)$

**Input:** $\boxed{H}$ , $\boxed{s} = \boxed{H} \boxed{x}$

**Recover:** $\boxed{e}$ s.t $\boxed{H} \boxed{e} = \boxed{s}$ and $\boxed{e} \in \left\{z \ : \ |z| = t\right\}$

For a fixed $R = k/n$, with respect to $\tau = t/n$, the solution will be unique or not!

Let, $\varepsilon = \mathbb{P}_{H,x}\left(\mathcal{A}(H, s = xH^{T}) = e \text{ such that } |e| = t \text{ and } eH^{T} = s\right)$

Using the law of total probability:

$$\varepsilon = \frac{1}{q^{k \times (n-k)} \times (q-1)^t \binom{n}{t}} \sum_{\substack{x_0 \in \mathbb{F}_q^n, \ |x_0|=t \\ H_0 \in \mathbb{F}_q^{(n-k) \times n}}} \mathbb{P}\left(\mathcal{A}(H_0, s = x_0 H^{T}) = e \text{ s.t } |e| = t \text{ and } eH^{T} = s\right)$$

$\longrightarrow \varepsilon$ is the average success probability of $\mathcal{A}$ over all fixed possible inputs

$\left(\text{above probabilities are computed over the internal randomness of } \mathcal{A}\right)$

Consequence:

If $\varepsilon$ is negligible, then $\mathcal{A}$ fails to decode almost all codes

**Exponential Complexity for Decoding in Average:**

For all known algorithms $\mathcal{A}$ ($T$ running time of one iteration $\mathcal{A}$)

$$\frac{T}{\varepsilon} = 2^{\alpha(q,R,\tau)\, n(1+o(1))} \text{ for some } \alpha(q,R,\tau) \geq 0$$
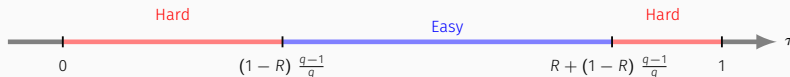


Figure 1: Hardness of DP($n, q, R, \tau$) as function of $\tau$

**Exponential Complexity for Decoding in Average:**

For all known algorithms $\mathcal{A}$ ($T$ running time of one iteration $\mathcal{A}$)

$$\frac{T}{\varepsilon} = 2^{\alpha(q,R,\tau)\, n(1+o(1))} \text{ for some } \alpha(q,R,\tau) \geq 0$$
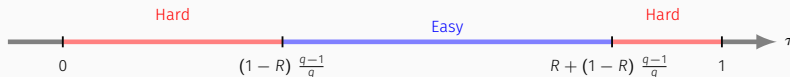


Figure 1: Hardness of DP($n, q, R, \tau$) as function of $\tau$

▶ McEliece encryption: $t = \tau n = \Theta\left(\frac{n}{\log n}\right)$

▶ Other encryptions: $t = \tau n = \Theta\left(\sqrt{n}\right)$

▶ Authenticated protocols: $t = \tau n = Cn$ where $C$ constant quite small

▶ Wave Signature: $t = \tau n = Cn$ where $C$ large constant, $C \approx 0.95$

DP$'$ $(n, q, R, \tau)$. Let $k \stackrel{\text{def}}{=} \lfloor Rn \rfloor$ and $t \stackrel{\text{def}}{=} \lfloor \tau n \rfloor$

- **Input:** $(\mathbf{G}, \mathbf{y} \stackrel{\text{def}}{=} \mathbf{s}\mathbf{G} + \mathbf{x})$ where $\mathbf{G}$, $\mathbf{s}$ and $\mathbf{x}$ are uniformly distributed over $\mathbb{F}_q^{k \times n}$, $\mathbb{F}_q^k$ and words of Hamming weight $t$ in $\mathbb{F}_q^n$.

- **Output:** an error $\mathbf{e} \in \mathbb{F}_q^n$ of Hamming weight $t$ such that $\mathbf{y} - \mathbf{e} = \mathbf{m}\mathbf{G}$ for some $\mathbf{m} \in \mathbb{F}_q^k$.

**Exercise Session:**

For any algorithm $\mathcal{A}$ solving DP$'$ with probability $\varepsilon$ and time $T$:

Describe an algorithm $\mathcal{B}$ solving DP in the $\approx$ same time with probability $\geq \varepsilon - O\left(q^{-\min(k, n-k)}\right)$

(and the reciprocal)

$\longrightarrow$ Same average hardness with syndromes or noisy codewords formalism!