# Code-based cryptography I

## Basic concepts and McElice system

Tanja Lange
with some slides by Tung Chou and Christiane Peters

Eindhoven University of Technology

SAC – Post-quantum cryptography

# Error correction

- Digital media is exposed to memory corruption.
- Many systems check whether data was corrupted in transit:
  - ISBN numbers have check digit to detect corruption.
  - ECC RAM detects up to two errors and can correct one error.
    64 bits are stored as 72 bits: extra 8 bits for checks and recovery.
- In general, $k$ bits of data get stored in $n$ bits, adding some redundancy.
- If no error occurred, these $n$ bits satisfy $n - k$ parity check equations; else can correct errors from the error pattern.
- Good codes can correct many errors without blowing up storage too much;
  offer guarantee to correct $t$ errors (often can correct or at least detect more).

# Linear codes

A binary linear code $C$ of length $n$ and dimension $k$ is a $k$-dimensional subspace of $\mathbb{F}_2^n$.
$C$ is usually specified as

- the row space of a generating matrix $G \in \mathbb{F}_2^{k \times n}$

$$C = \{\mathbf{m}G | \mathbf{m} \in \mathbb{F}_2^k\}$$

- the kernel space of a parity-check matrix $H \in \mathbb{F}_2^{(n-k) \times n}$

$$C = \{\mathbf{c} | H\mathbf{c}^\mathsf{T} = 0, \ \mathbf{c} \in \mathbb{F}_2^n\}$$

  Leaving out the $^\mathsf{T}$ from now on.

- Names: code word $\mathbf{c}$, error vector $\mathbf{e}$, received word $\mathbf{b} = \mathbf{c} + \mathbf{e}$.

# Example: Hamming code

Parity check matrix ($n = 7, k = 4$):

$$H = \begin{pmatrix} 1 & 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 & 1 \end{pmatrix}$$

An error-free string of 7 bits $\mathbf{b} = (b_0, b_1, b_2, b_3, b_4, b_5, b_6)$ satisfies these three equations:

$$\begin{array}{rcl} b_0 + b_1 + b_3 + b_4 & = & 0 \\ b_0 + b_2 + b_3 + b_5 & = & 0 \\ b_1 + b_2 + b_3 + b_6 & = & 0 \end{array}$$

If one error occurred at least one of these equations will not hold.
Failure pattern uniquely identifies the error location,
e.g., $1, 0, 1$ means

# Example: Hamming code

Parity check matrix ($n = 7, k = 4$):

$$H = \begin{pmatrix} 1 & 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 & 1 \end{pmatrix}$$

An error-free string of 7 bits $\mathbf{b} = (b_0, b_1, b_2, b_3, b_4, b_5, b_6)$ satisfies these three equations:

$$
\begin{array}{ccccccccc}
b_0 & +b_1 & & +b_3 & +b_4 & & & = & 0 \\
b_0 & & +b_2 & +b_3 & & +b_5 & & = & 0 \\
& b_1 & +b_2 & +b_3 & & & +b_6 & = & 0
\end{array}
$$

If one error occurred at least one of these equations will not hold.
Failure pattern uniquely identifies the error location,
e.g., $1, 0, 1$ means $b_1$ flipped.

# Example: Hamming code

Parity check matrix ($n = 7, k = 4$):

$$H = \begin{pmatrix} 1 & 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 & 1 \end{pmatrix}$$

An error-free string of 7 bits $\mathbf{b} = (b_0, b_1, b_2, b_3, b_4, b_5, b_6)$ satisfies these three equations:

$$
\begin{array}{rcl}
b_0 + b_1 + b_3 + b_4 & = & 0 \\
b_0 + b_2 + b_3 + b_5 & = & 0 \\
b_1 + b_2 + b_3 + b_6 & = & 0
\end{array}
$$

If one error occurred at least one of these equations will not hold.
Failure pattern uniquely identifies the error location,
e.g., $1, 0, 1$ means $b_1$ flipped.
In math notation, the failure pattern is $H \cdot \mathbf{b}$.

# Linear codes are linear

Example with generator matrix:

$$G = \begin{pmatrix} 1 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 \end{pmatrix}$$

$\mathbf{c} = (111)G = (10011)$ is a code word.

# Linear codes are linear

Example with generator matrix:

$$G = \begin{pmatrix} 1 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 \end{pmatrix}$$

$\mathbf{c} = (111)G = (10011)$ is a code word.

Linear codes are linear:
The sum of two code words is a code word:

# Linear codes are linear

Example with generator matrix:

$$G = \begin{pmatrix} 1 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 \end{pmatrix}$$

$\mathbf{c} = (111)G = (10011)$ is a code word.

Linear codes are linear:
The sum of two code words is a code word:

$$\mathbf{c}_1 + \mathbf{c}_2 = \mathbf{m}_1 G + \mathbf{m}_2 G = (\mathbf{m}_1 + \mathbf{m}_2)G.$$

Same with parity-check matrix:

# Linear codes are linear

Example with generator matrix:

$$G = \begin{pmatrix} 1 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 \end{pmatrix}$$

$\mathbf{c} = (111)G = (10011)$ is a code word.

Linear codes are linear:
The sum of two code words is a code word:

$$\mathbf{c}_1 + \mathbf{c}_2 = \mathbf{m}_1 G + \mathbf{m}_2 G = (\mathbf{m}_1 + \mathbf{m}_2)G.$$

Same with parity-check matrix:

$$H(\mathbf{c}_1 + \mathbf{c}_2) = H\mathbf{c}_1 + H\mathbf{c}_2 = 0 + 0 = 0.$$

# Hamming weight and distance

- The Hamming weight of a word is the number of nonzero coordinates.

$$\mathrm{wt}(1, 0, 0, 1, 1) = 3$$

- The Hamming distance between two words in $\mathbb{F}_2^n$ is the number of coordinates in which they differ.

$$d((1, 1, 0, 1, 1), (1, 0, 0, 1, 1)) =$$

# Hamming weight and distance

- The Hamming weight of a word is the number of nonzero coordinates.
$$\mathrm{wt}(1, 0, 0, 1, 1) = 3$$

- The Hamming distance between two words in $\mathbb{F}_2^n$ is the number of coordinates in which they differ.

$$d((1, 1, 0, 1, 1), (1, 0, 0, 1, 1)) = 1$$

# Hamming weight and distance

- The Hamming weight of a word is the number of nonzero coordinates.
$$\mathrm{wt}(1, 0, 0, 1, 1) = 3$$

- The Hamming distance between two words in $\mathbb{F}_2^n$ is the number of coordinates in which they differ.
$$d((1, 1, 0, 1, 1), (1, 0, 0, 1, 1)) = 1$$

The Hamming distance between **x** and **y** equals the Hamming weight of **x** + **y**:
$$d((1, 1, 0, 1, 1), (1, 0, 0, 1, 1)) = \mathrm{wt}(0, 1, 0, 0, 0).$$

# Minimum distance

- The minimum distance of a linear code $C$ is the smallest Hamming weight of a nonzero code word in $C$.

$$d = \min_{0 \neq \mathbf{c} \in C} \{\mathrm{wt}(\mathbf{c})\} = \min_{\mathbf{b} \neq \mathbf{c} \in C} \{d(\mathbf{b}, \mathbf{c})\}$$

- In code with minimum distance $d = 2t + 1$, any vector $\mathbf{x} = \mathbf{c} + \mathbf{e}$ with $\mathrm{wt}(\mathbf{e}) \leq t$ is uniquely decodable to $\mathbf{c}$;
  i. e. there is no closer code word.

# Decoding problem

Decoding problem: find the closest code word $\mathbf{c} \in C$ to a given $\mathbf{x} \in \mathbb{F}_2^n$, assuming that there is a unique closest code word. Let $\mathbf{x} = \mathbf{c} + \mathbf{e}$. Note that finding $\mathbf{e}$ is an equivalent problem.

- If $\mathbf{c}$ is $t$ errors away from $\mathbf{x}$, i.e., the Hamming weight of $\mathbf{e}$ is $t$, this is called a $t$-error correcting problem.
- There are lots of code families with fast decoding algorithms, e.g., Reed–Solomon codes, Goppa codes/alternant codes, etc.
- However, the general decoding problem is hard: Information-set decoding (see later) takes exponential time.

# The McEliece cryptosystem I

- Due to Robert McEliece 1978.
- Let $C$ be a length-$n$ binary Goppa code $\Gamma$ of dimension $k$ with minimum distance $2t + 1$ where $t \approx (n - k)/\log_2(n)$; original parameters (1978) $n = 1024$, $k = 524$, $t = 50$.
- The McEliece secret key consists of a generator matrix $G$ for $\Gamma$, an efficient $t$-error correcting decoding algorithm for $\Gamma$; an $n \times n$ permutation matrix $P$ and a nonsingular $k \times k$ matrix $S$.
- $n, k, t$ are public; but $\Gamma$, $P$, $S$ are randomly generated secrets.
- The McEliece public key is the $k \times n$ matrix $G' = SGP$.

# The McEliece cryptosystem II

- Encrypt: Compute $\mathbf{m}G'$ and add a random error vector $\mathbf{e}$ of weight $t$ and length $n$. Send $\mathbf{y} = \mathbf{m}G' + \mathbf{e}$.
- Decrypt: Compute $\mathbf{y}P^{-1} = \mathbf{m}G'P^{-1} + \mathbf{e}P^{-1} = (\mathbf{m}S)G + \mathbf{e}P^{-1}$. This works because $\mathbf{e}P^{-1}$ has the same weight as $\mathbf{e}$

# The McEliece cryptosystem II

- Encrypt: Compute $\mathbf{m}G'$ and add a random error vector $\mathbf{e}$ of weight $t$ and length $n$. Send $\mathbf{y} = \mathbf{m}G' + \mathbf{e}$.

- Decrypt: Compute $\mathbf{y}P^{-1} = \mathbf{m}G'P^{-1} + \mathbf{e}P^{-1} = (\mathbf{m}S)G + \mathbf{e}P^{-1}$.
  This works because $\mathbf{e}P^{-1}$ has the same weight as $\mathbf{e}$
  because $P$ is a permutation matrix.
  Use fast decoding to find $\mathbf{m}S$ and $\mathbf{m}$.

- Attacker is faced with decoding $\mathbf{y}$ to nearest code word $\mathbf{m}G'$ in the code generated by $G'$.
  This is general decoding if $G'$ does not expose any structure.