# Coding Theory:
# From the Past to the Present

## Vitaly Skachek
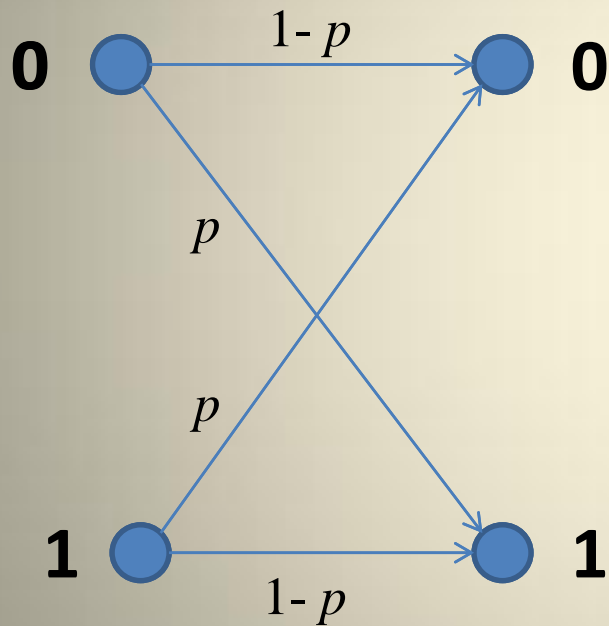
### Institute of Computer Science
### University of Tartu

# Communications Model

Source — $k$ bits — $x$ **0101** → Encoder — $R = k/n$ — $c$ **0101100** → Channel

Destination ← $c$ **0101100** ← Decoder ← $n$ bits — $y$ **0111000** ← Channel
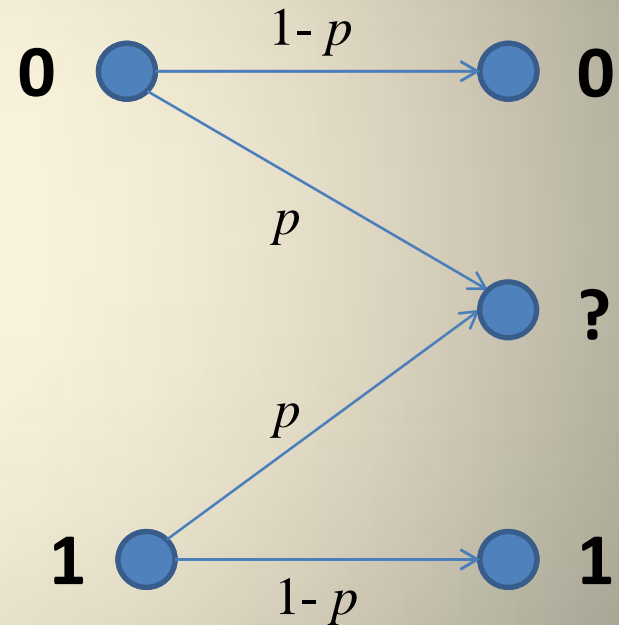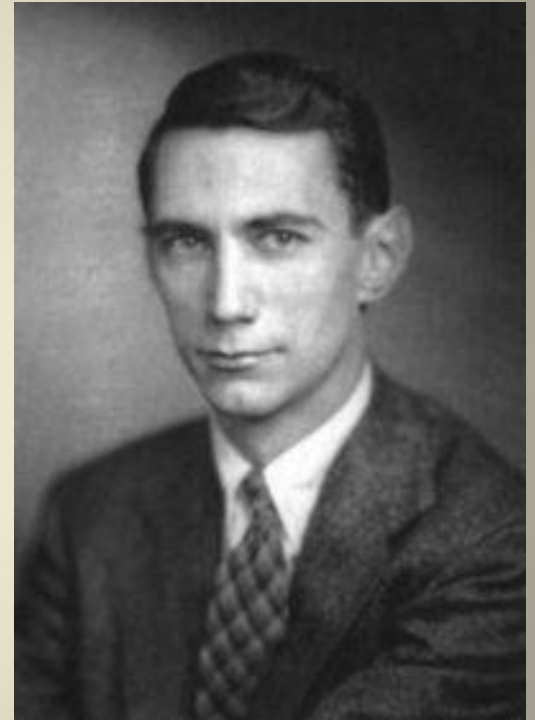
# Communications Channels



Binary Symmetric
Channel

Binary Erasure
Channel

# Shannon's Channel Coding Theorems

- A code is a mapping from the set of all vectors of length $k$ to a set of vectors of length $n$ (over alphabet $\Sigma$)

- Given a channel $S$, there is a quantity $C(S)$ called channel capacity

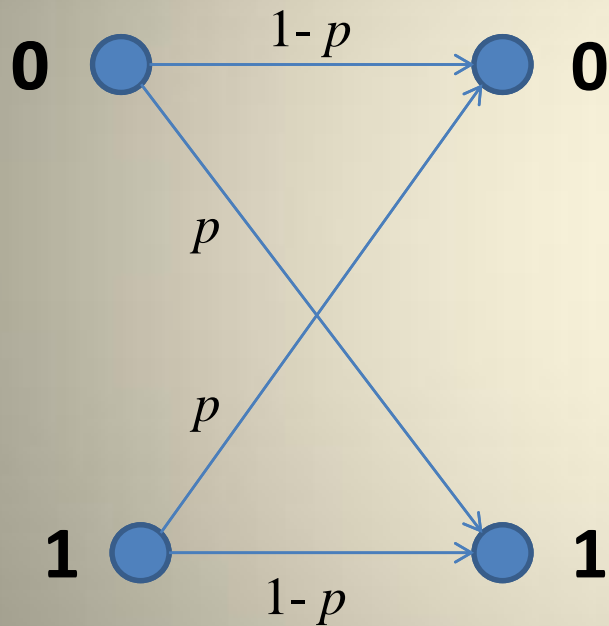Claude Shannon
(1916-2001)

# Shannon's Channel Coding Theorems

For any rate $R < C(S)$, there exists an infinite sequence of block codes $C_i$ of growing lengths $n_i$ such that $\frac{k_i}{n_i} \geq R$ , and there exists a coding scheme for those codes such that the decoding error probability approaches 0 as $i \rightarrow \infty$.
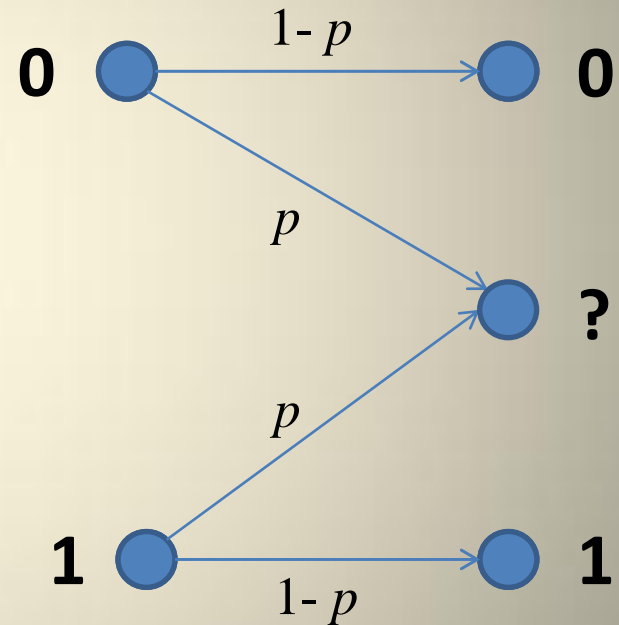
# Shannon's Channel Coding Theorems

For any rate $R < C(S)$, there exists an infinite sequence of block codes $C_i$ of growing lengths $n_i$ such that $\frac{k_i}{n_i} \geq R$ , and there exists a coding scheme for those codes such that the decoding error probability approaches 0 as $i \rightarrow \infty$.

Let $R > C(S)$. For any infinite sequence of block codes $C_i$ of growing lengths $n_i$ such that $\frac{k_i}{n_i} \geq R$ , and for any coding scheme for those codes, the decoding error probability is bounded away from 0 as $i \rightarrow \infty$.
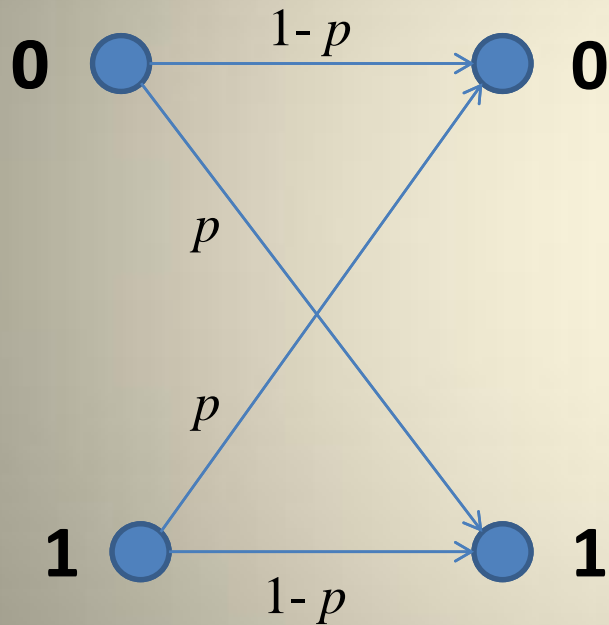
# Communications Channels
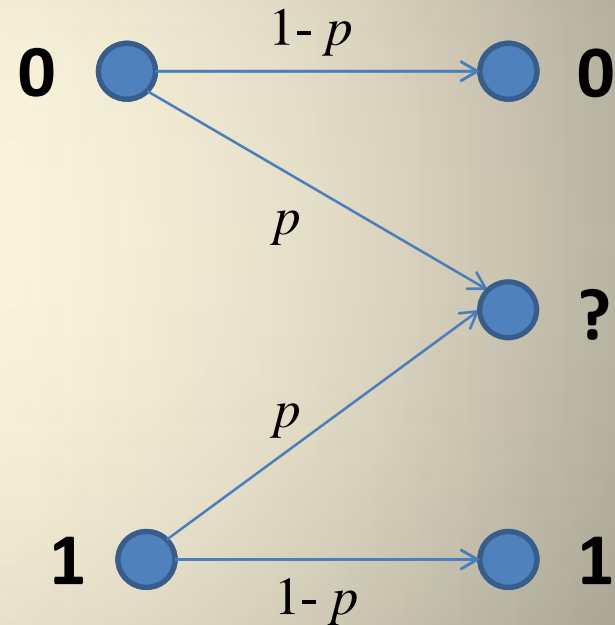


Binary Symmetric Channel

Binary Erasure Channel

# Communications Channels

$C(S)=1\text{-}h_2(p)$             $C(S)=1\text{-}p$
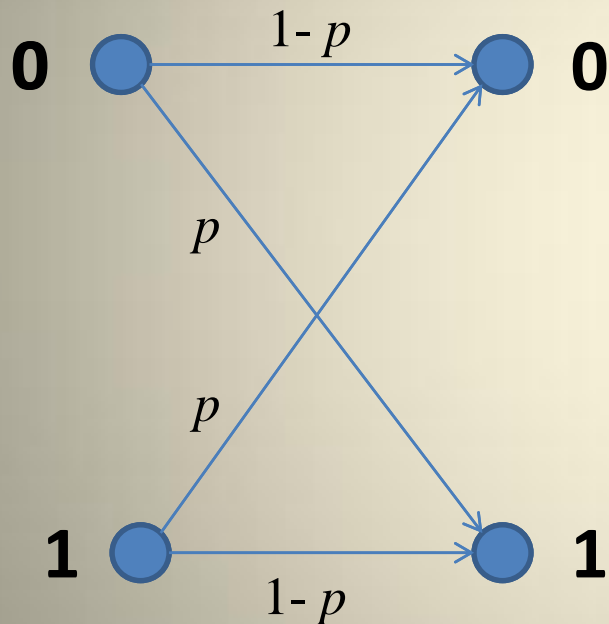


Binary Symmetric
Channel

Binary Erasure
Channel

# Communications Channels

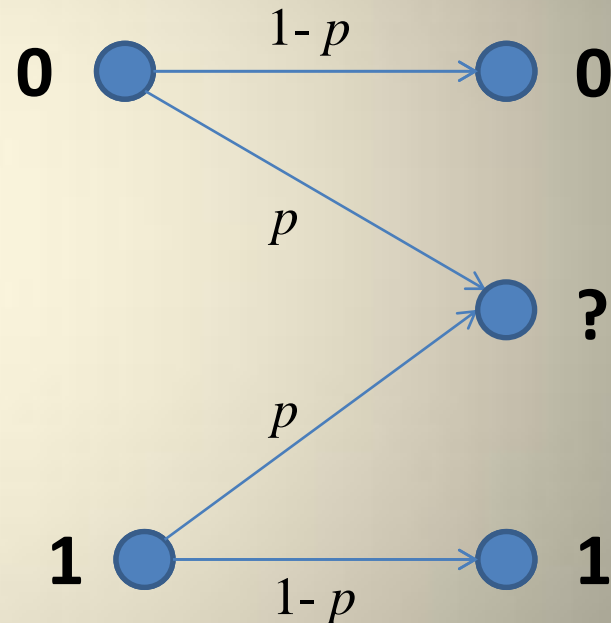$C(S)=1-h_2(p)$                    $C(S)=1-p$

$h_2(x) = -x \log x - (1-x)log(1-x)$
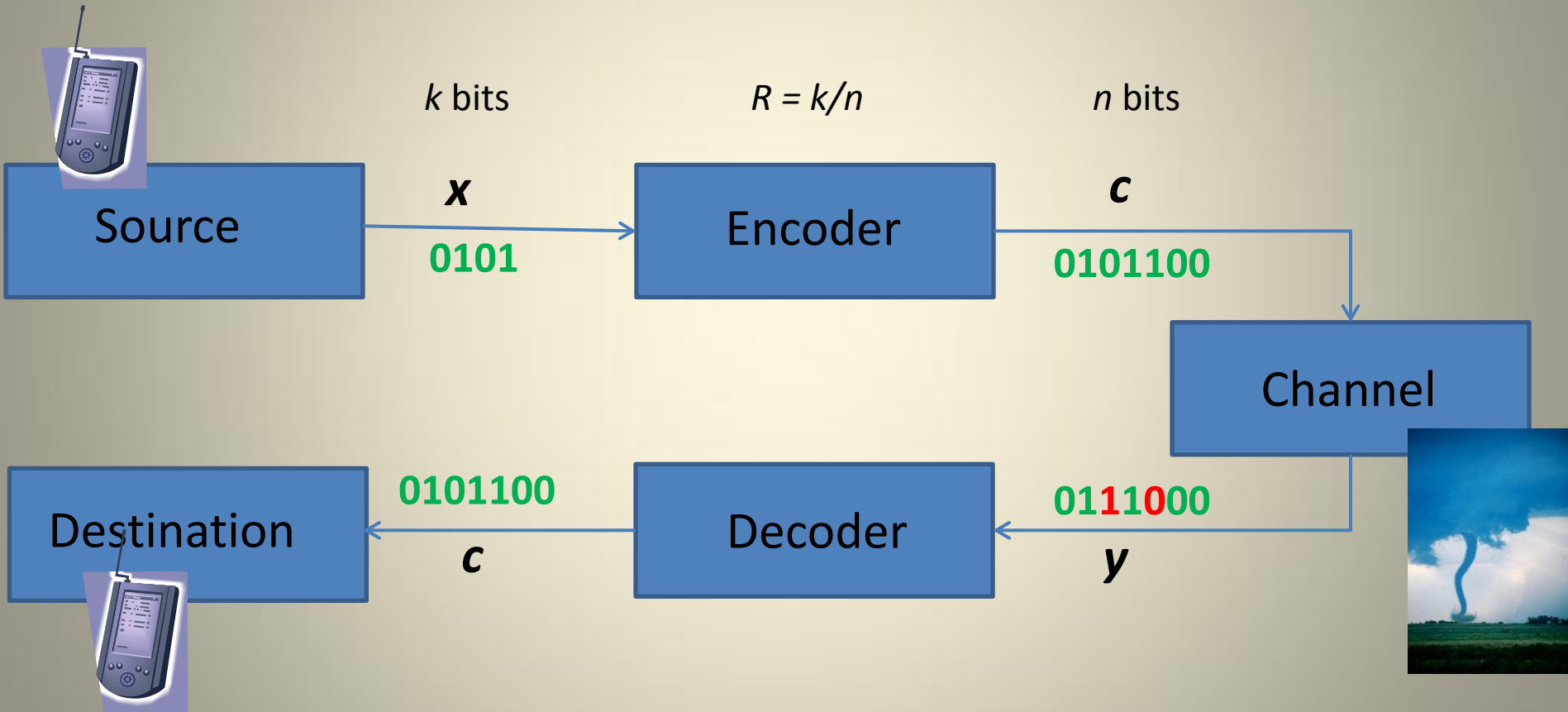


Binary Symmetric
Channel

Binary Erasure
Channel

# Communications Model

# Parameters in Consideration

- Target: optimize the code rate $R = k/n$.

Other parameters in considerations:

- Speed of convergence $\Pr(\text{err}) \to 0$ as $n \to \infty$. Low error probability for short lengths is needed!

- Time complexity of encoding and decoding algorithms. Structured codes are needed!

# Distance

- The Hamming distance between
  $x = (x_1, x_2, \ldots, x_n)$ and $y = (y_1, y_2, \ldots, y_n)$,
  $d(x, y)$, is the number of pairs of symbols
  $(x_i, y_i)$, such that $x_i \neq y_i$.

- The minimum distance of a code $C$ is
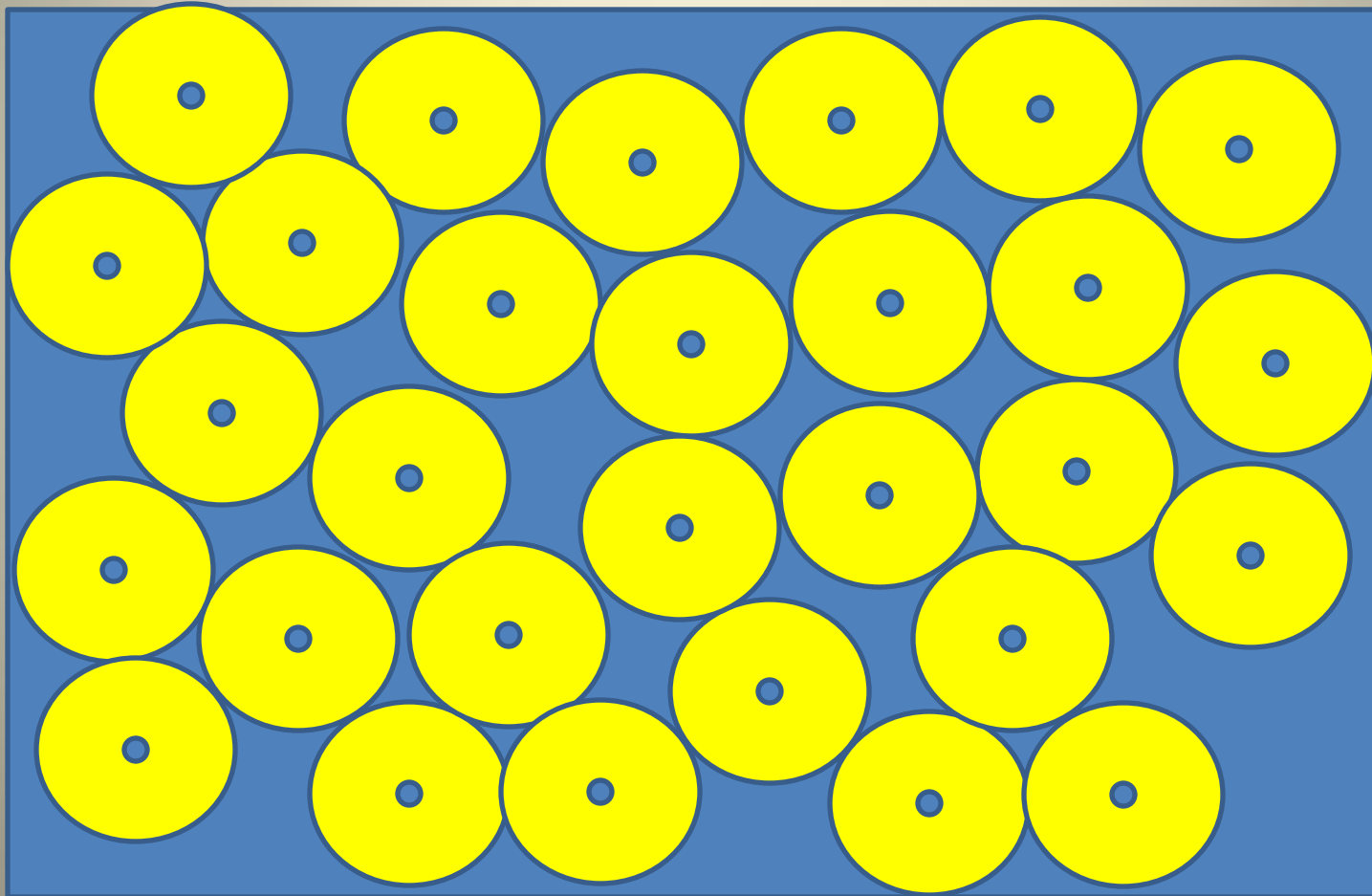  $$d = \min_{\{x, y \in C, x \neq y\}} d(x, y)$$

# Linear Codes

- A code $C$ over field $F$ is a linear [*n, k, d*] code if there exists a matrix $H$ with *n* columns and rank n − k such that
$$H \cdot c^T = 0^T \iff c \in C.$$

- The matrix $H$ is called a parity-check matrix.
- The value $k$ is called the dimension of the code $C$.
- The ratio $R = k/n$ is called the rate of the code $C$.
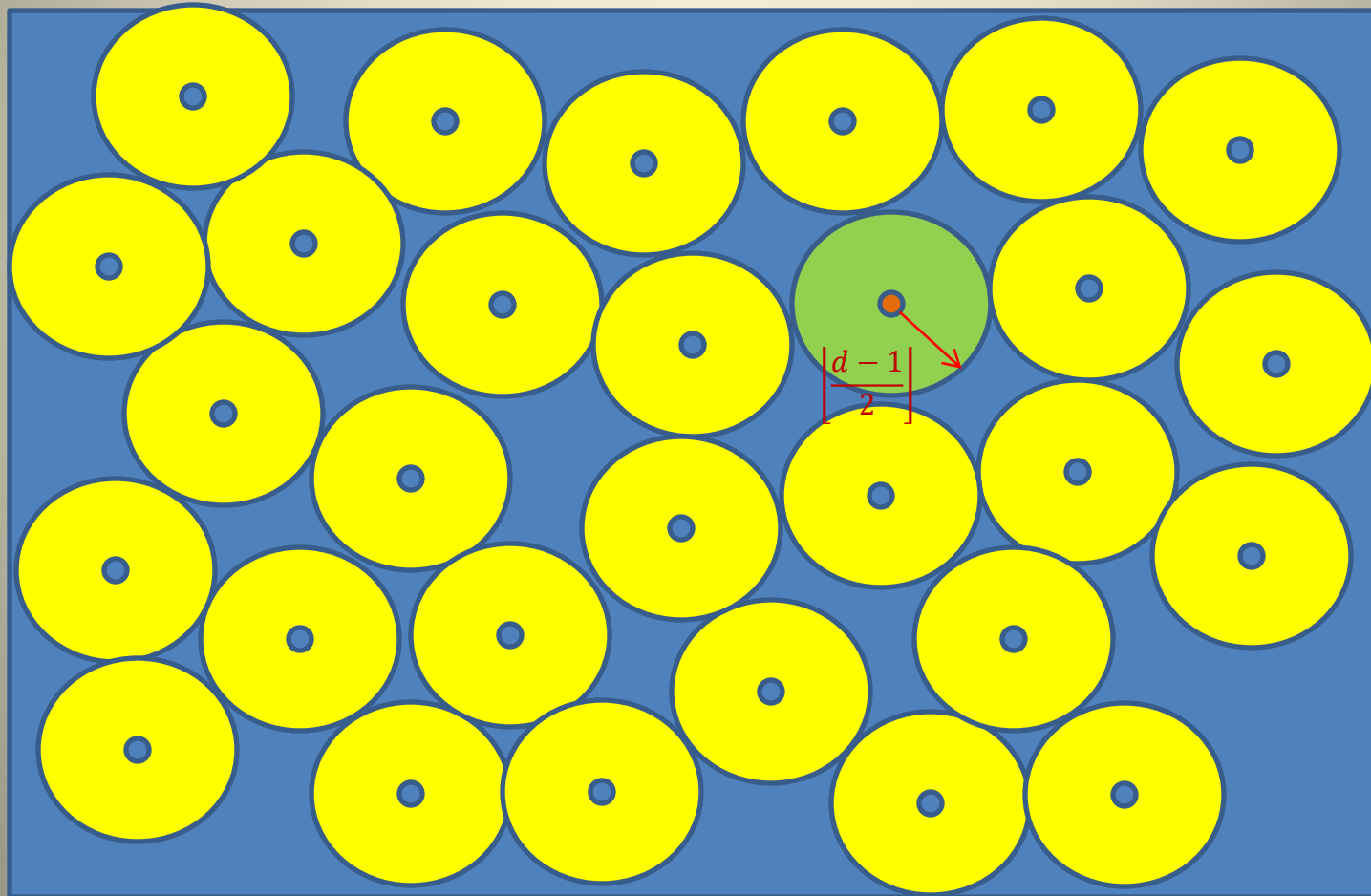- All words of $C$ are exactly all linear combinations of rows of a generating $k \times n$ matrix $G$.

# Sphere-packing idea

# Sphere-packing idea

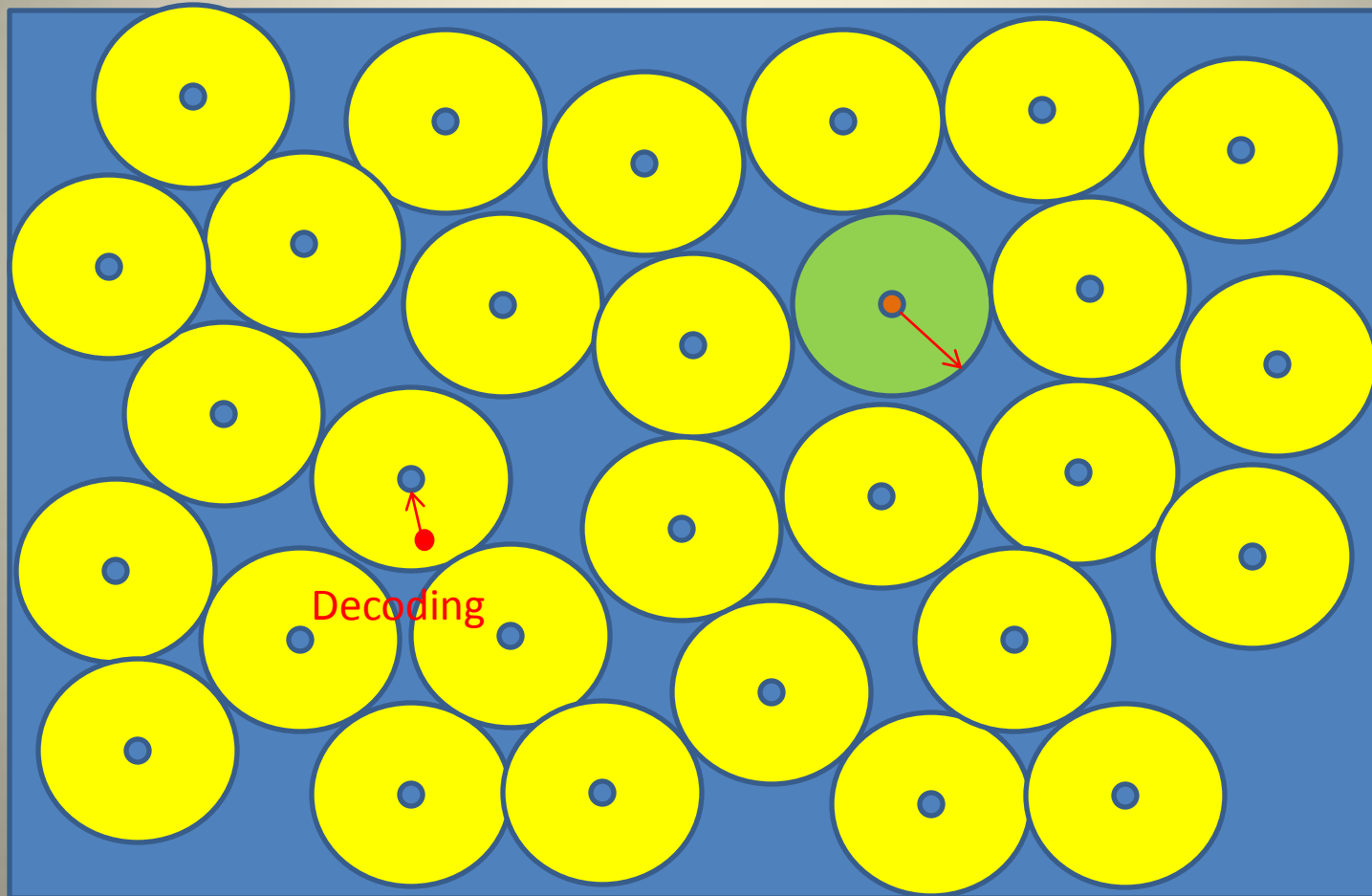# Sphere-packing idea

# Sphere-packing idea



Decoding

# Reed-Solomon Codes

- Let $\alpha_1, \alpha_2, \ldots, \alpha_n \in F$ be *n* distinct elements.
- The generator matrix:

$$G = \begin{bmatrix} 1 & 1 & \ldots & 1 \\ \alpha_1 & \alpha_2 & \ldots & \alpha_n \\ \alpha_1^2 & \alpha_2^2 & \ldots & \alpha_n^2 \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_1^{k-1} & \alpha_2^{k-1} & \ldots & \alpha_n^{k-1} \end{bmatrix}$$

- Satisfies the <span style="color:red">Singleton bound</span>: *n = d + k − 1*
  - Optimal trade-off between the parameters

# Reed-Solomon Codes (cont.)

- Encoding:

$$[x_0 x_1 \dots x_{k-1}] \cdot \begin{bmatrix} 1 & 1 & \dots & 1 \\ \alpha_1 & \alpha_2 & \dots & \alpha_n \\ \alpha_1^2 & \alpha_2^2 & \dots & \alpha_n^2 \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_1^{k-1} & \alpha_2^{k-1} & \dots & \alpha_n^{k-1} \end{bmatrix}$$

# Polynomial Interpolation Viewpoint

- Input vector $[x_0 x_1 \ldots x_{k-1}]$ is associated with polynomial
$$P(z) = x_{k-1}z^{k-1} + x_{k-2}z^{k-2} + x_1 z + x_0$$

- Encoding is a substitution:
$$(P(\alpha_1), P(\alpha_2), \ldots, P(\alpha_n))$$

- Decoding is an interpolation by degree $\leq k-1$ polynomial

# Reed-Solomon Codes are Used in:

- Wired and wireless communications

- Satellite communications
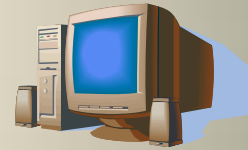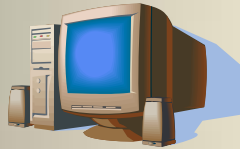
- Hard drives and compact disks

- Flash memory devices

# Application of Reed-Solomon Codes

- Shamir's Secret-Sharing Scheme '79

- $n$ users

- 1 key (number in $F$)

- Any coalition of $< t$ users does not have any information about the key
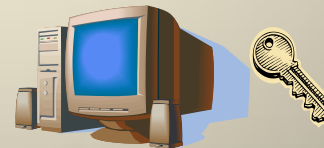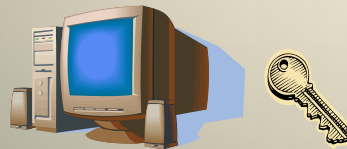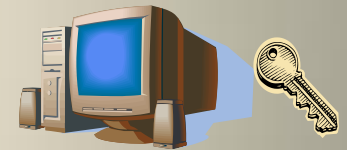
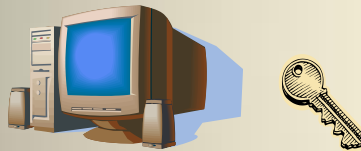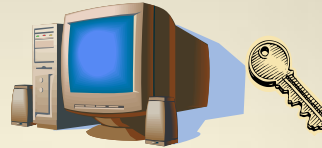- Any coalition of $\geq t$ users can recover the key



Adi Shamir

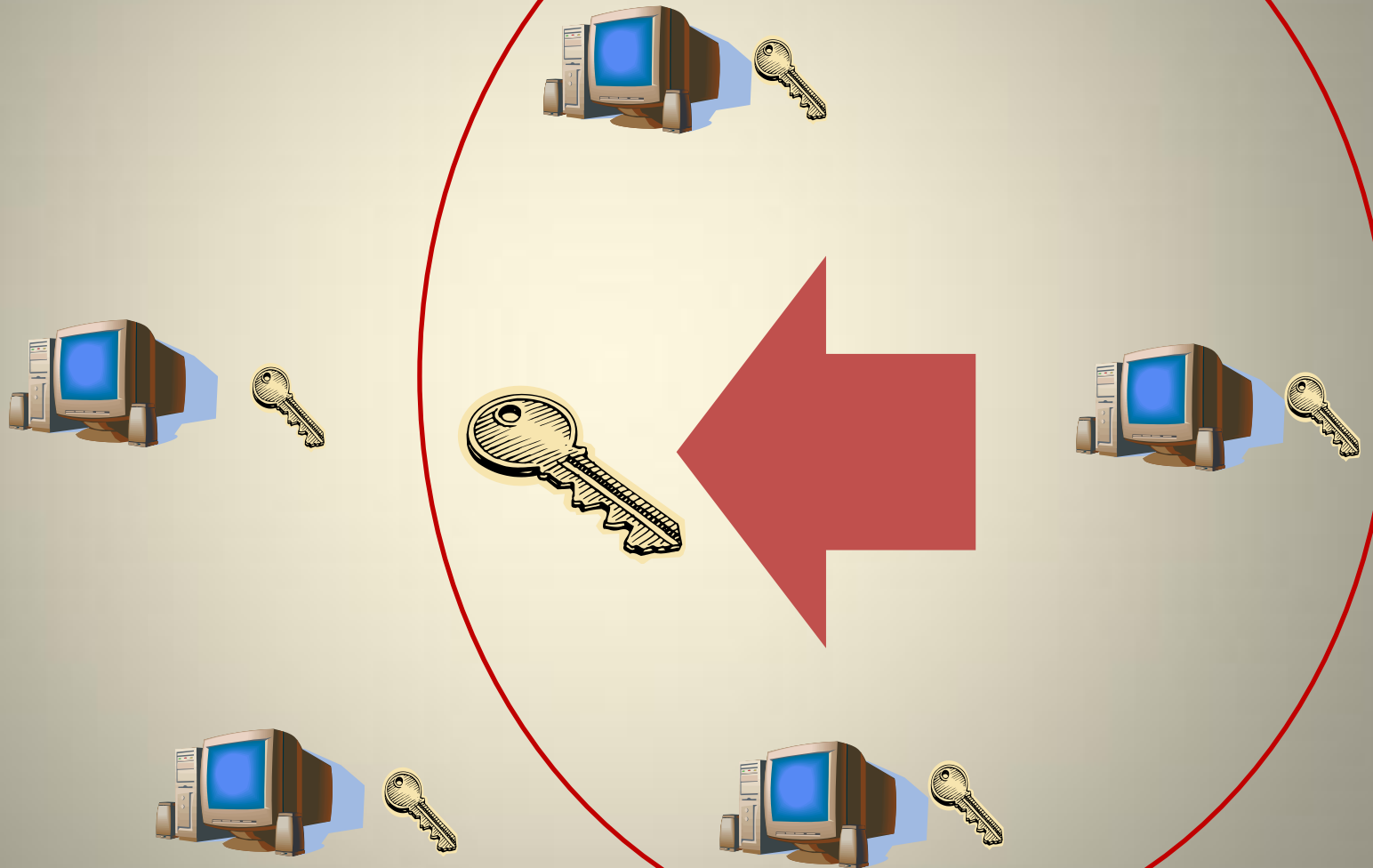# Shamir's Secret Sharing Scheme

# Shamir's Secret Sharing Scheme
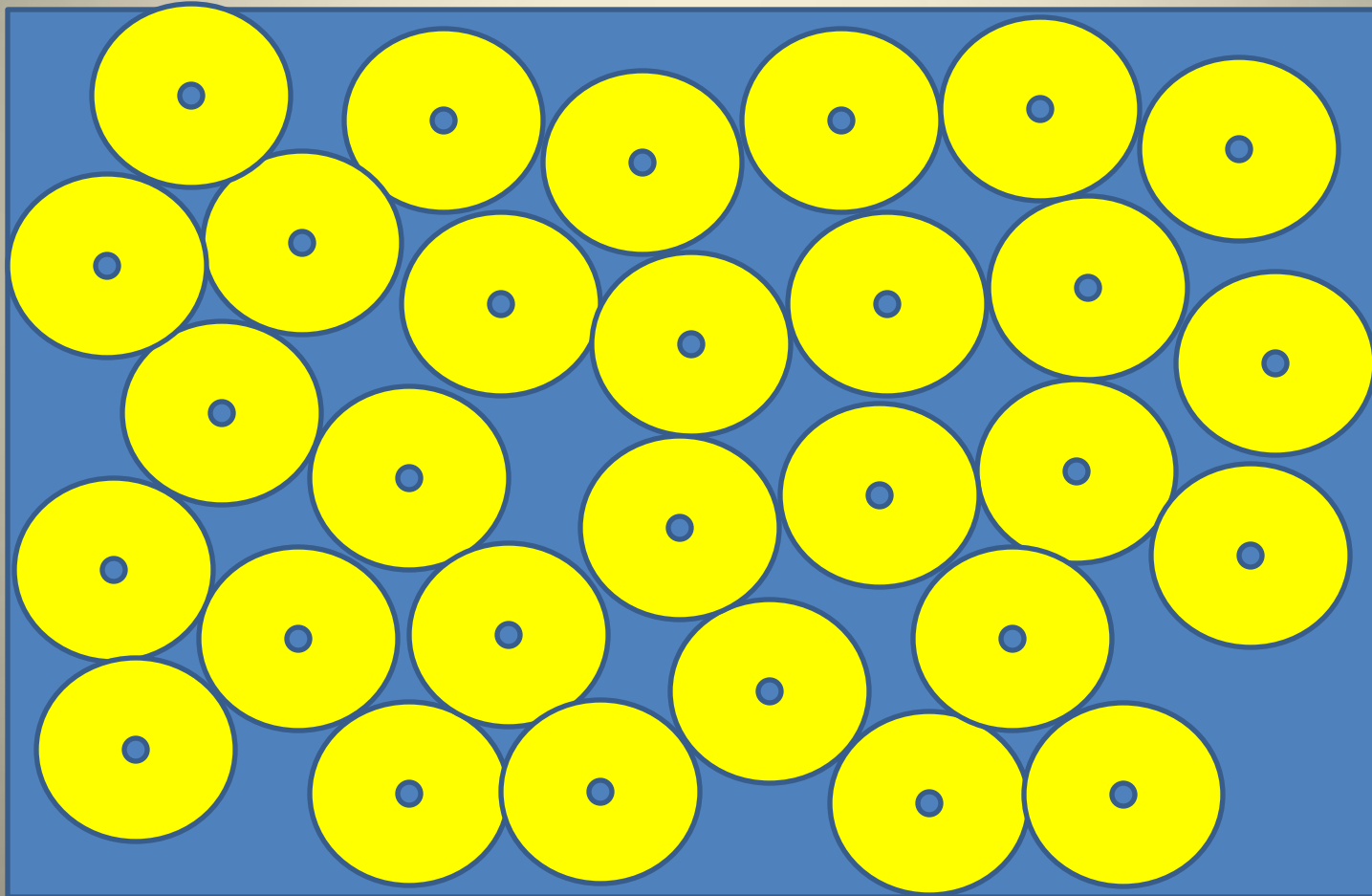
# Shamir's Secret Sharing Scheme

# Shamir's Secret Sharing Scheme (cont.)

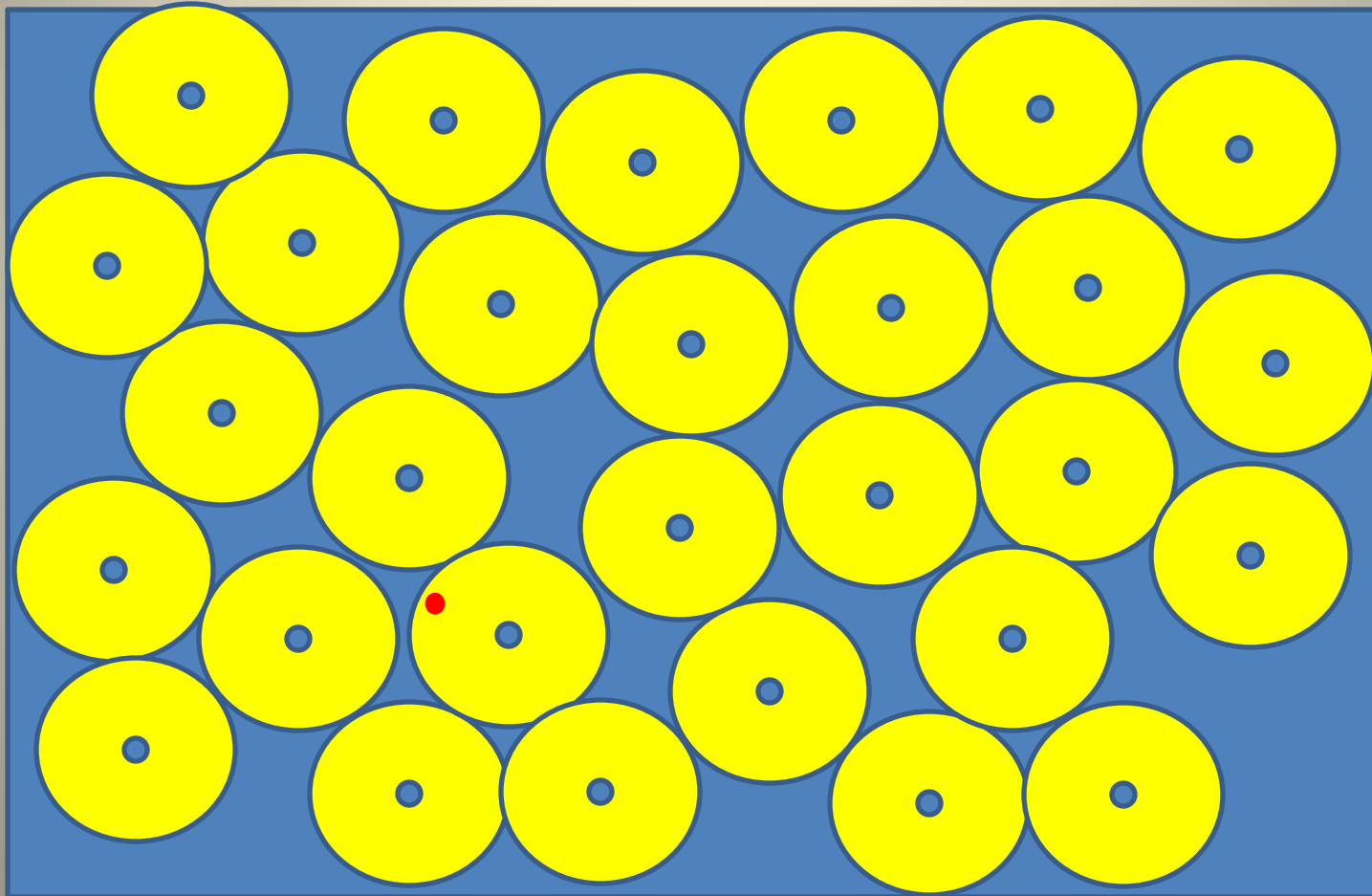- Select randomly $x_1, x_2, \ldots, x_{k-1}$. Let $x_0$ be a secret key. Construct polynomial
$$P(z) = x_{k-1}z^{k-1} + x_{k-2}z^{k-2} + x_1 z + x_0$$
- Give $(\alpha_i, P(\alpha_i))$ to user $i$

- Large coalition has enough points to reconstruct the polynomial

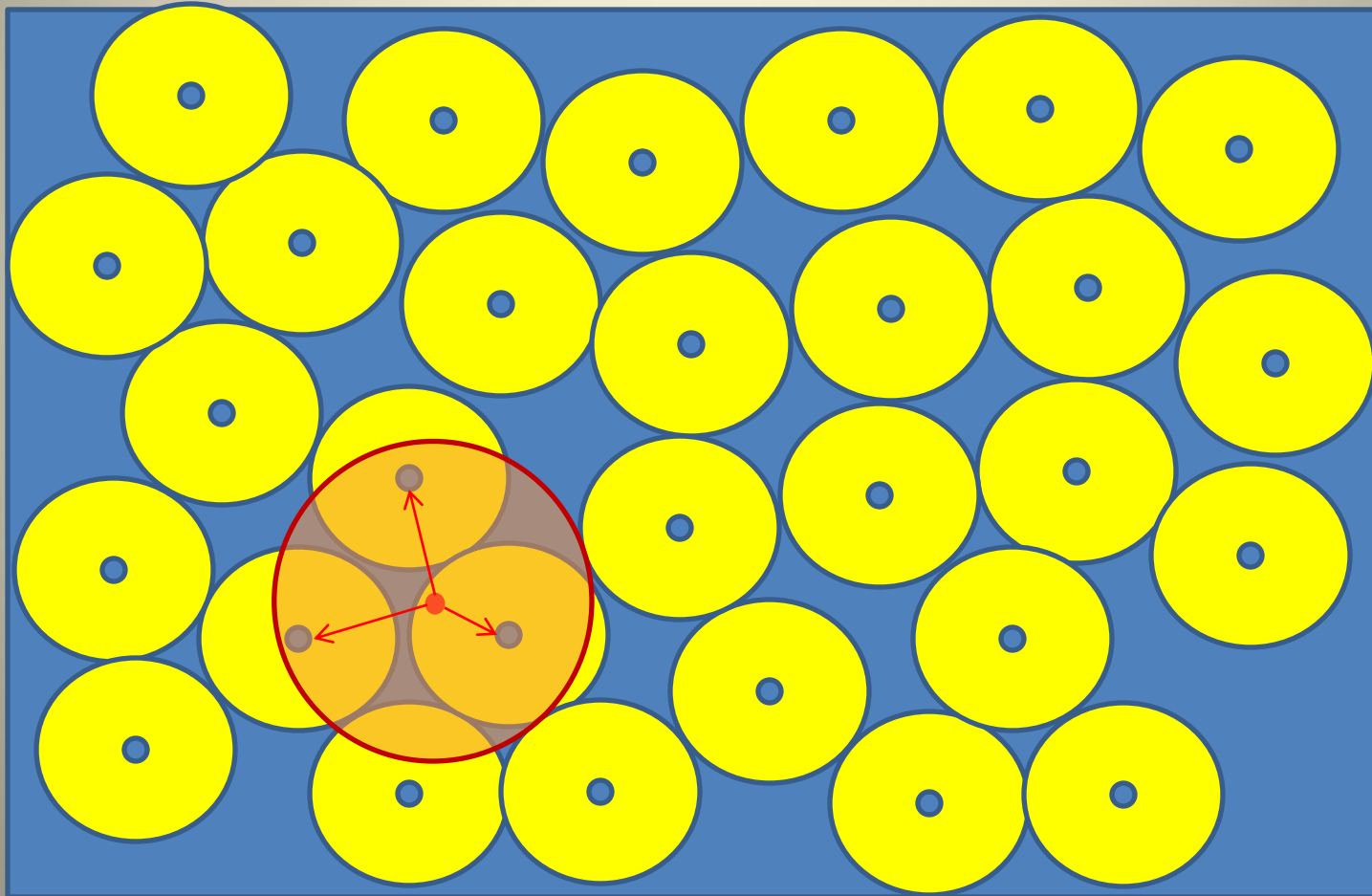- Small coalition has no information about the polynomial

# List-decoding of Reed-Solomon Codes

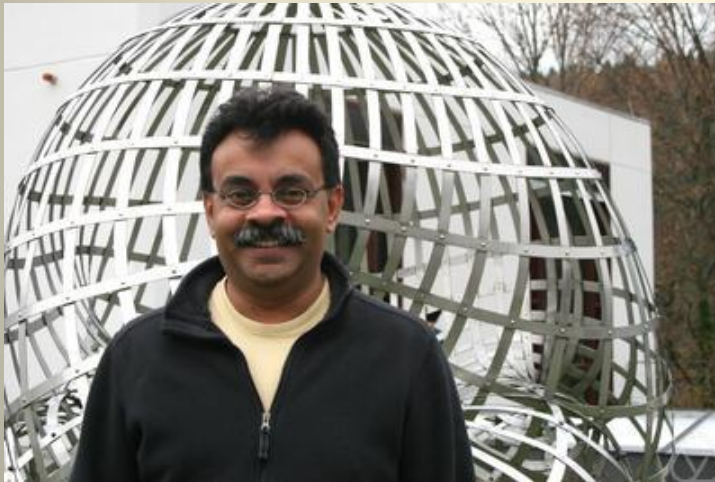# List-decoding of Reed-Solomon Codes

# List-decoding of Reed-Solomon Codes

# List-decoding of Reed-Solomon Codes

- Sudan '97, Guruswami '99, Vardy-Parvaresh '05, Guruswami-Rudra '06
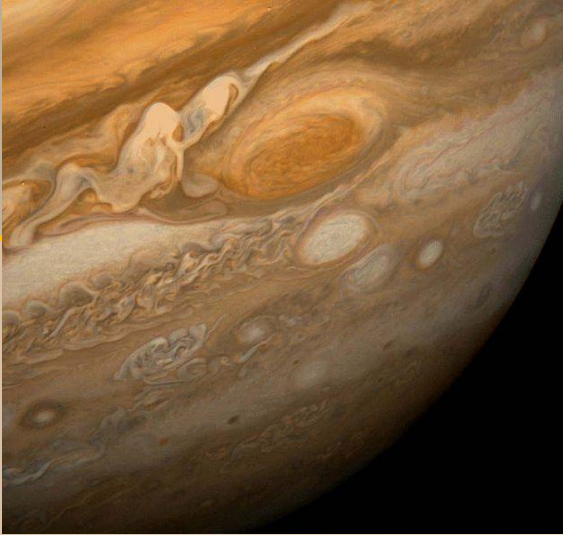


Madhu Sudan

Venkatesan Guruswami
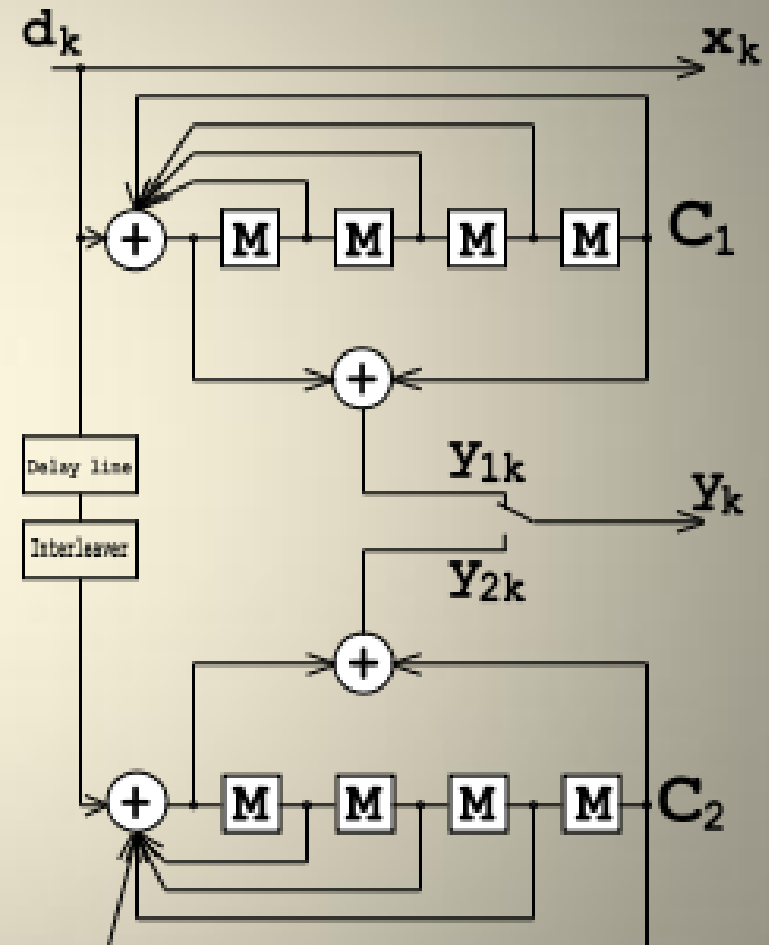
# List Decoding of RS Codes

Voyager 1 – the first manmade object to leave the Solar System. Launched in 1977.

# Turbo Codes

Berrou, Glavieux and Thitimajshima
(Telecom Bretagne) '93

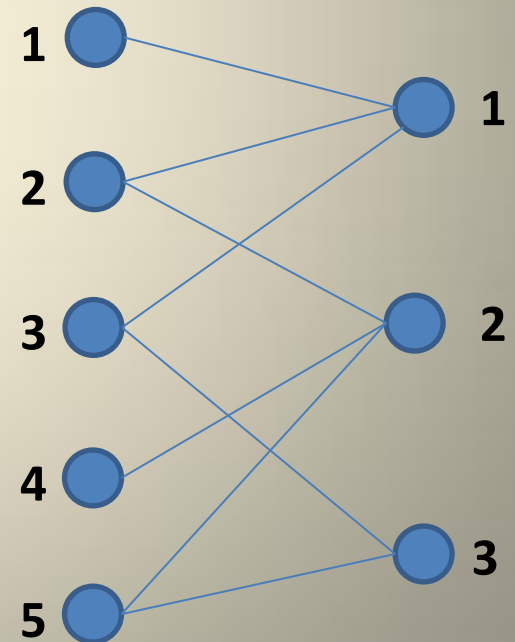- Non-algebraic codes!
- "Killer" of algebraic coding theory

# Low-Density Parity-Check Codes

- Gallager '62
- Urbanke, Richardson and Shokrollahi '01
- Parity-check matrix $H$ is sparse
- Performance extremely close to channel capacity
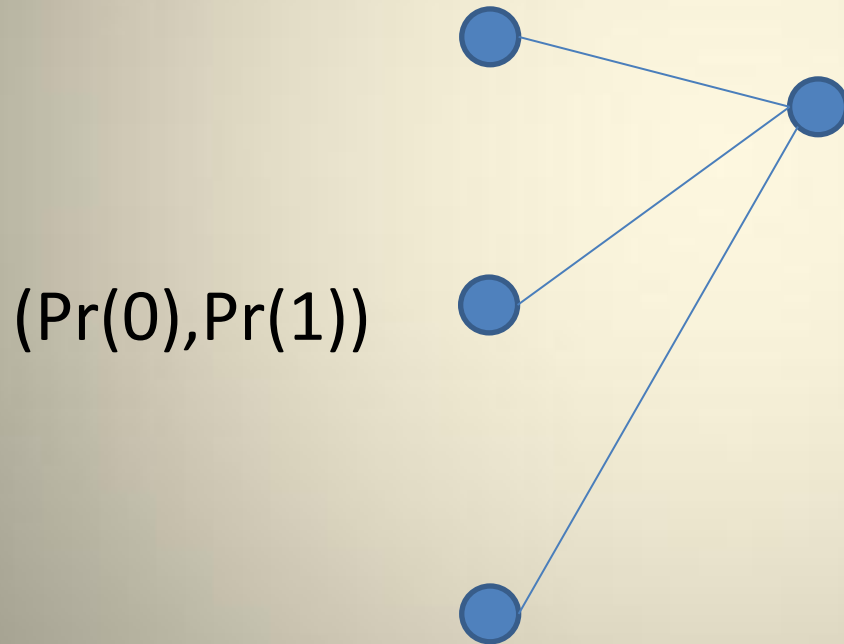- Decoding complexity linear in $n$

Tanner graph:

$$H = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 \end{bmatrix}$$
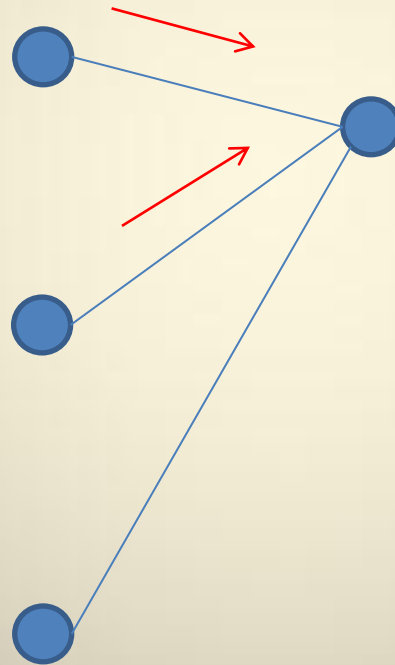
# Low-Density Parity-Check Codes

- Belief-propagation decoding algorithm (message-passing algorithm)

(Pr(0),Pr(1))

# Low-Density Parity-Check Codes
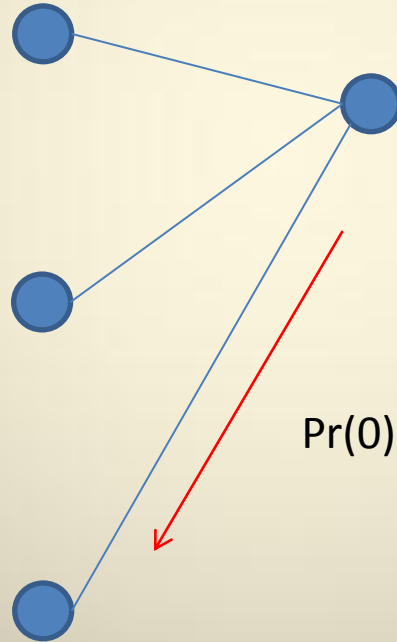
Pr(0) = 0.2, Pr(1) = 0.8

Pr(0) = 0.4, Pr(1) = 0.6

# Low-Density Parity-Check Codes



Pr(0) = 0.2, Pr(1) = 0.8

Pr(0) = 0.4, Pr(1) = 0.6
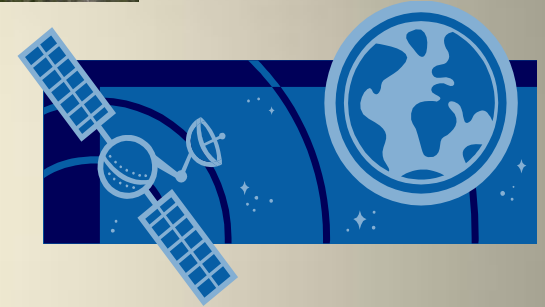
Pr(0) = 0.56, Pr(1) = 0.44

# Reed-Solomon Codes are Used in:

- Wired and wireless communications

- Satellite communications

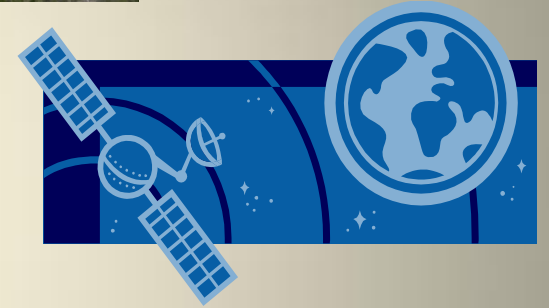- Hard drives and compact disks

- Flash memory devices

# LDPC Codes are Used in:

- Wired and wireless communications

- Satellite communications
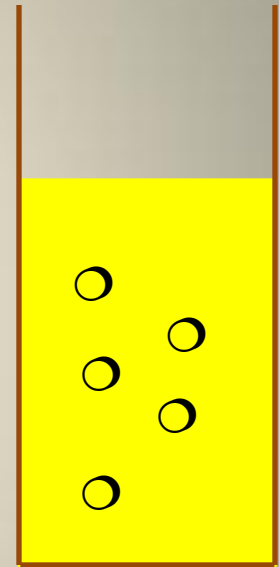
- Hard drives and compact disks

- Flash memory devices

# Emerging Applications of Coding Theory

# Flash memories

- Easy to add electric charge, hard to remove

- The charge "leaks" with the time

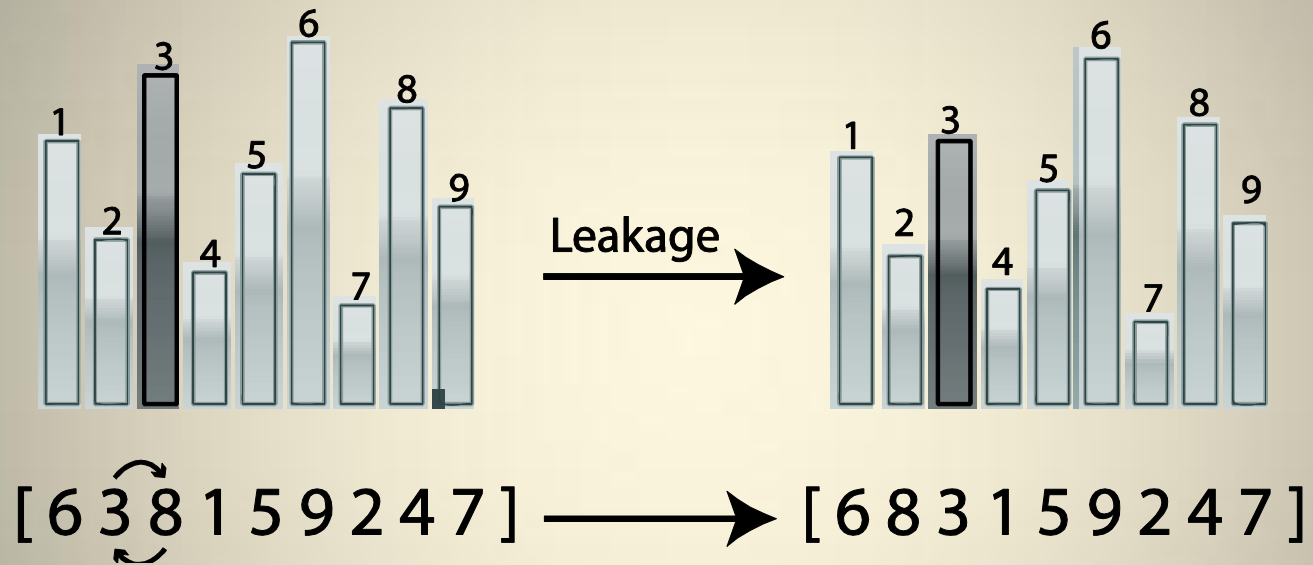- Neighboring cells influence each other

Flash memory cell

# Flash memories

- Rank modulation
- The information is represented using relative levels of charge, invariant to leakage
- Coding over permutations

Jiang, Mateescu, Schwartz, Bruck '2006
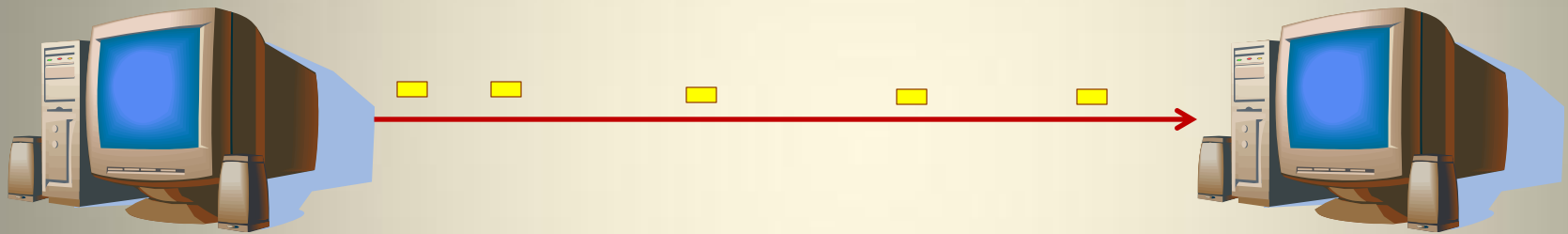
# Flash memories

# Networking

- Raptor Codes



- A. Shokrollahi '2004
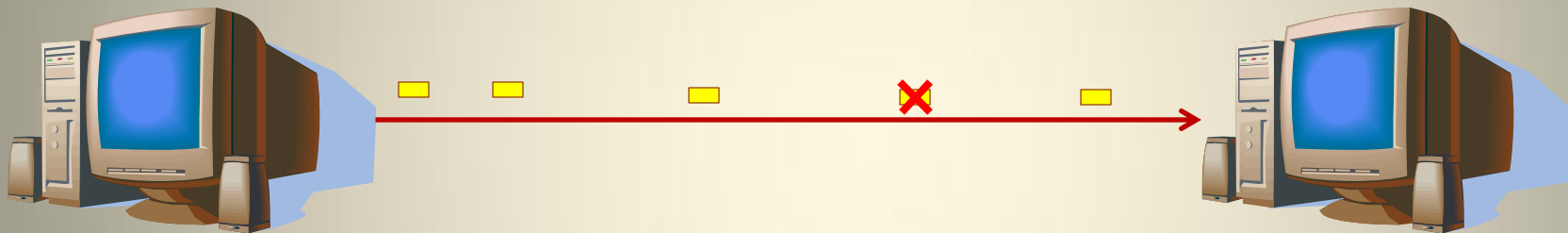- Used in DVB-H standard for IP datacast for handheld devices
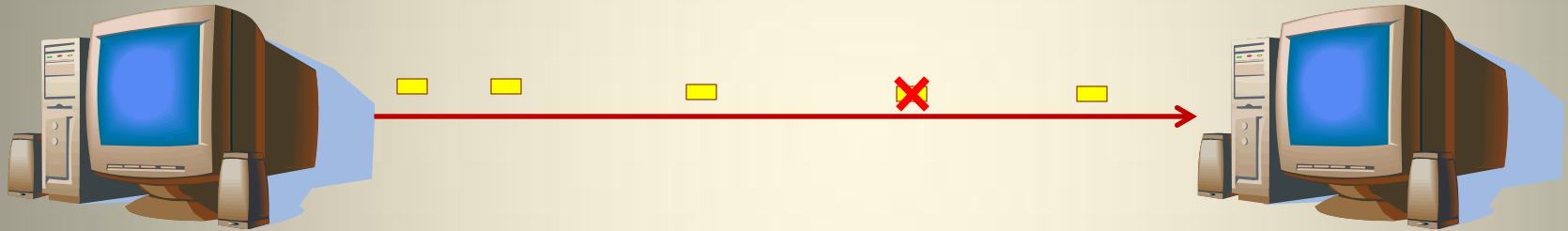
# Networking

- Raptor Codes

# Networking

- Raptor Codes

# Networking

- Raptor Codes



- Possible solution: ARQs (retransmissions) – slow!
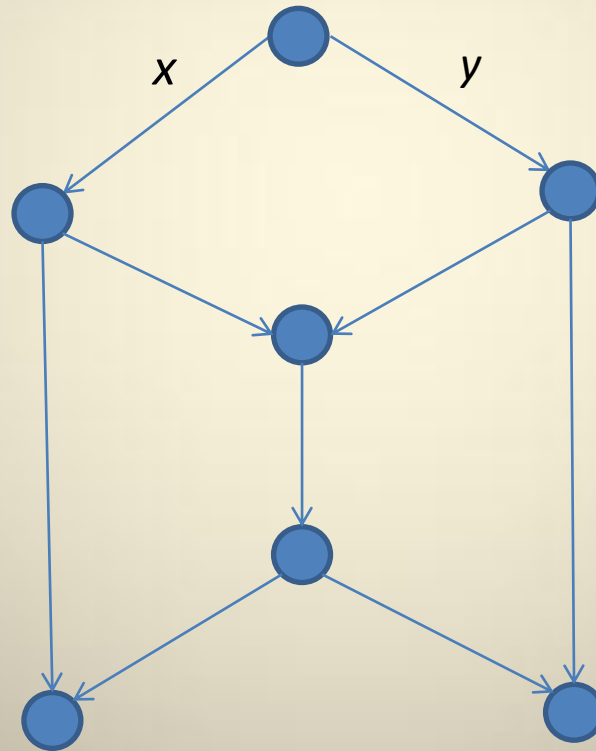- Alternative: large error-correcting code

# Networking

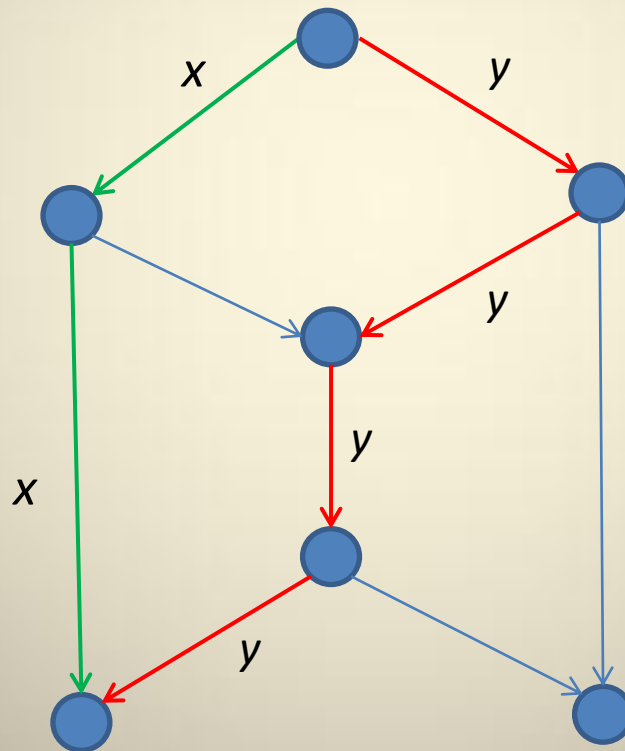- Raptor Codes

# Network coding

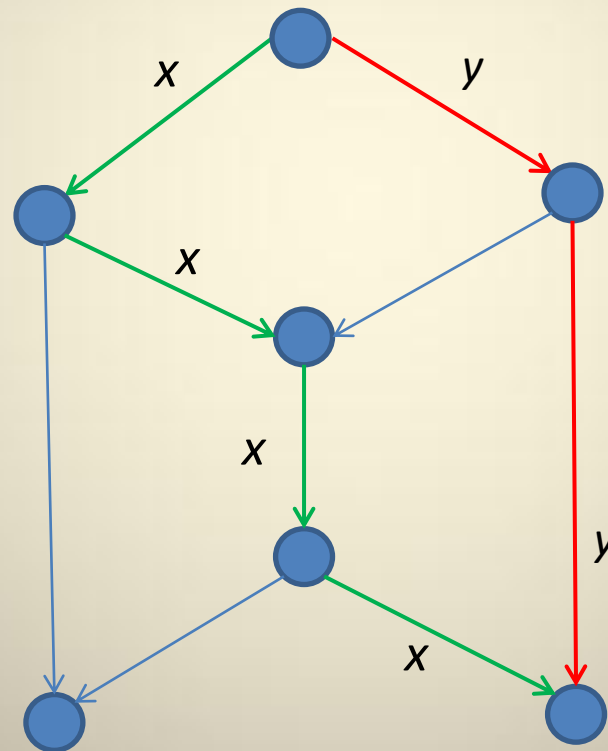- Butterfly network    Ahlswede, Cai, Li and Yeung, 2000
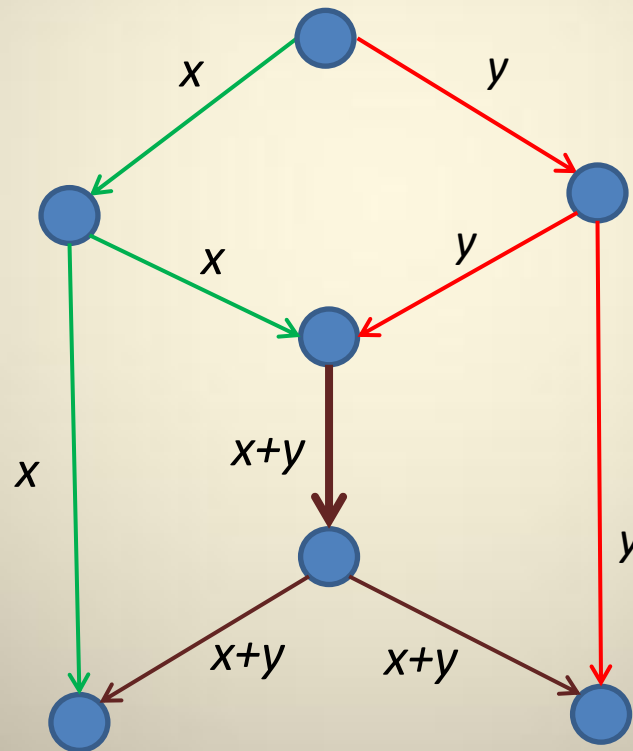
# Network coding

- Butterfly network
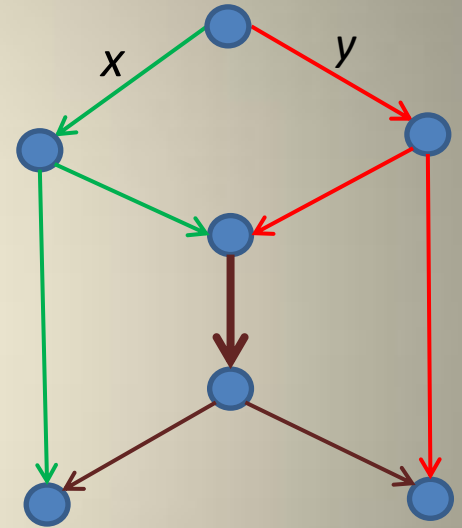
# Network coding

- Butterfly network

# Network coding

- Butterfly network

# Network coding

- The number of bits deliverable to each destination is equal to min-cut between source and each of destinations

- Avalanche P2P Network (Microsoft, 2005)

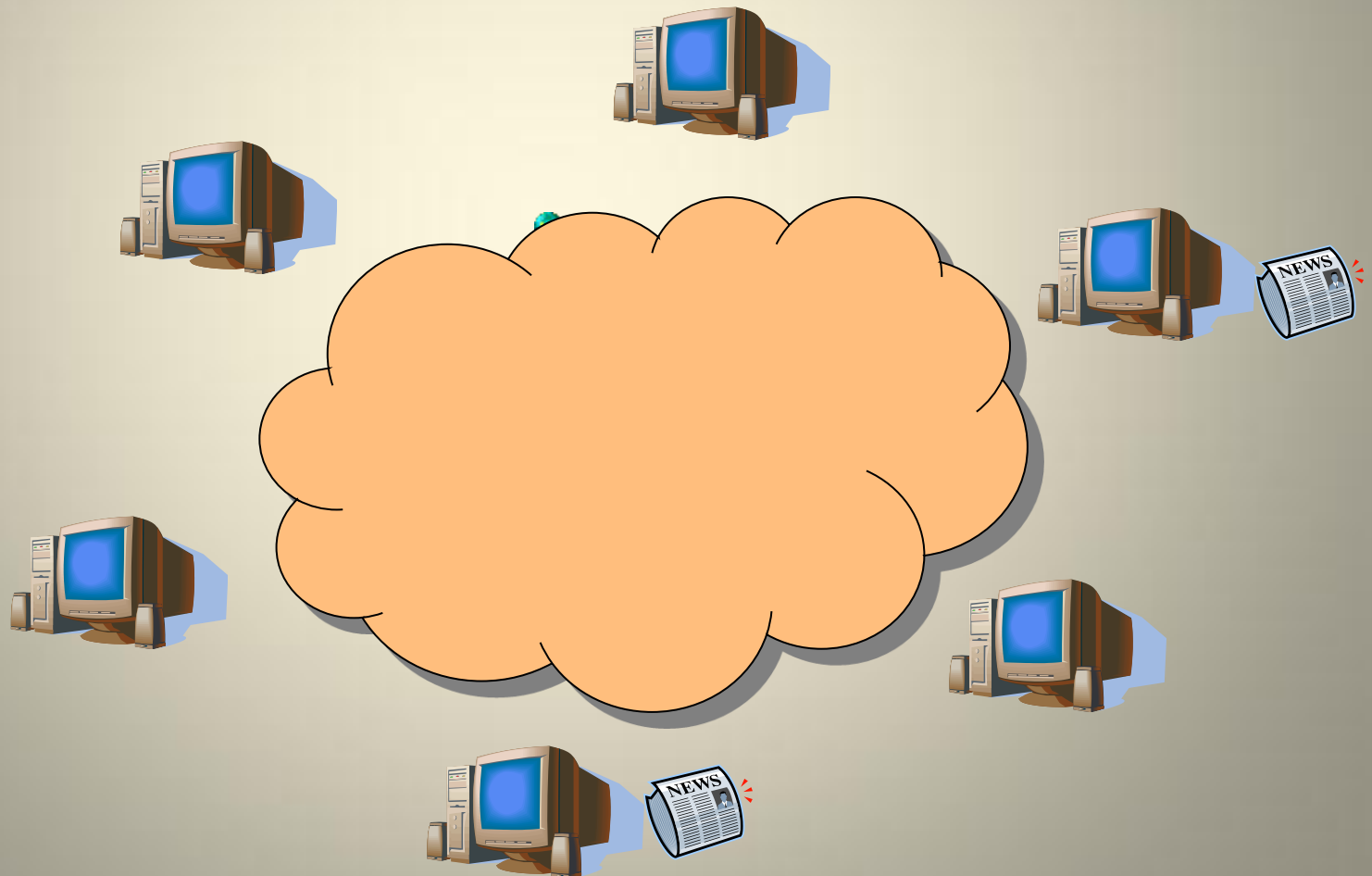- Experiments for use in mobile communications

# Gossip Algorithms

- *n* users in the network
- *k* of them possess a rumor (packet of data) – each rumor is different
- Each users "calls" another user randomly and sends a rumor to him
- Purpose: to distribute all rumors to all users
- Using coding: send a random linear combination of all rumors in your possession
  - Facilitates convergence of the algorithm
    Deb, Medard and Choute 2006

# Gossip Algorithms

- Rumor spreading problem

# Gossip Algorithms

- *n* users in the network
- *k* of them possess a rumor (packet of data) – each rumor is different
- Each users "calls" another user randomly and sends a rumor to him
- Purpose: to distribute all rumors to all users
- Using coding: send a random linear combination of all rumors in your possession
  - Facilitates convergence of the algorithm
    Deb, Medard and Choute 2006

# Distributed Storage

- Huge amounts of data stored by big data companies (Google, Amazon, Facebook, Dropbox)
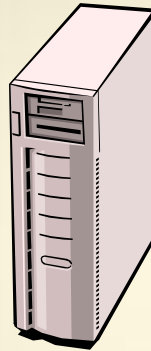


Facebook data center in Oregon
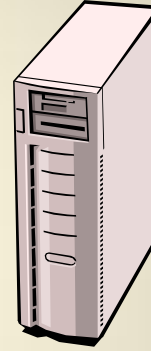


Server room at Wikipedia data center

# Distributed data storage
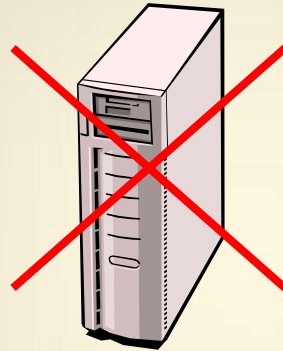
$x$                    $y$                    $x+ y$

Dimakis, Godfrey, Wu, Wainwright, Ramchandran '2008

# Distributed data storage
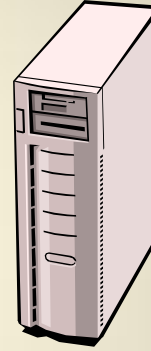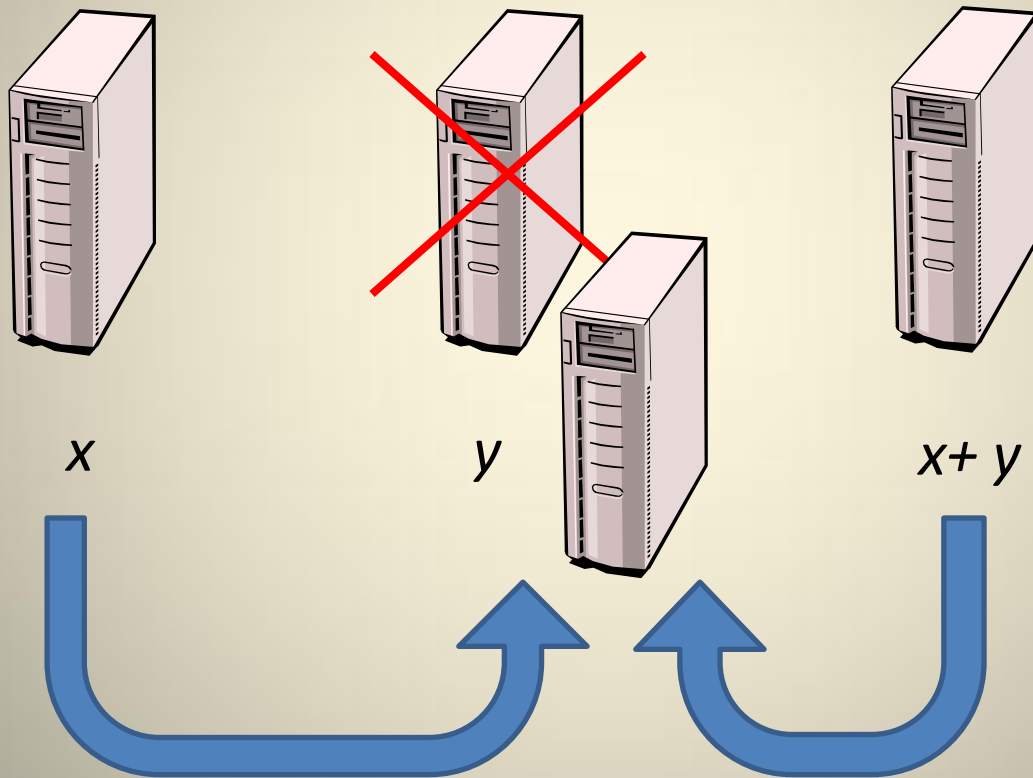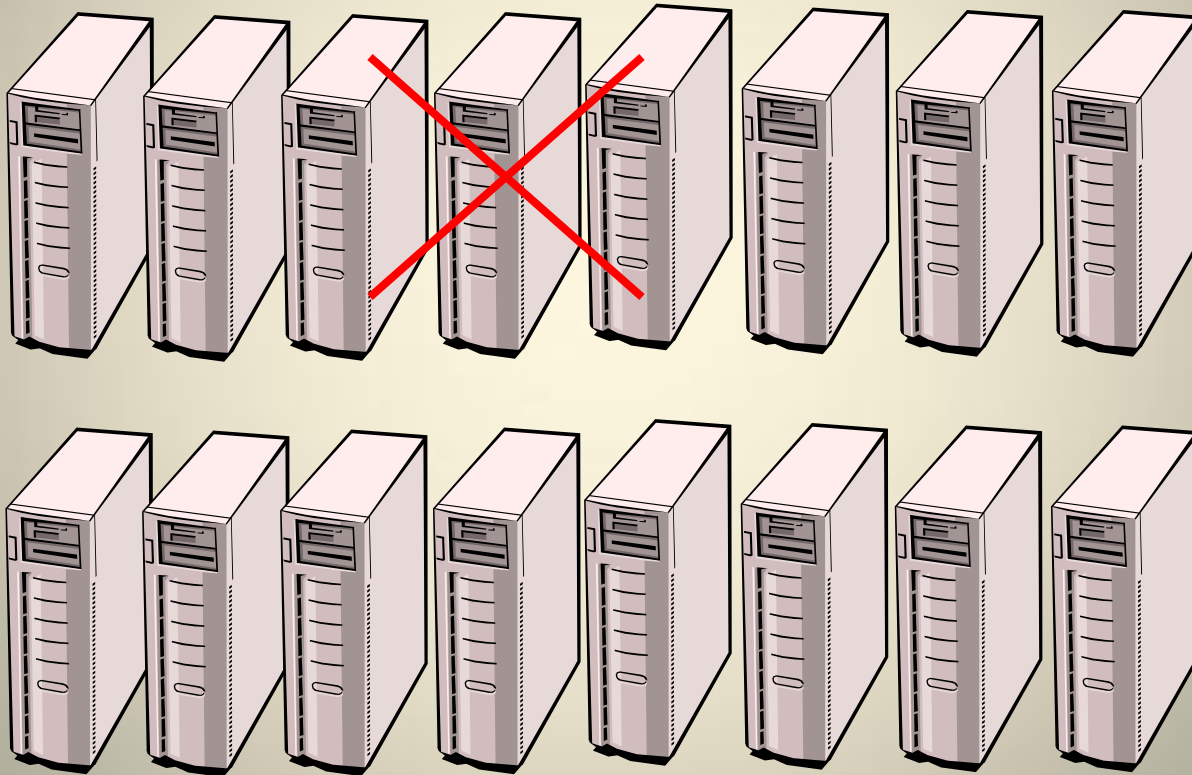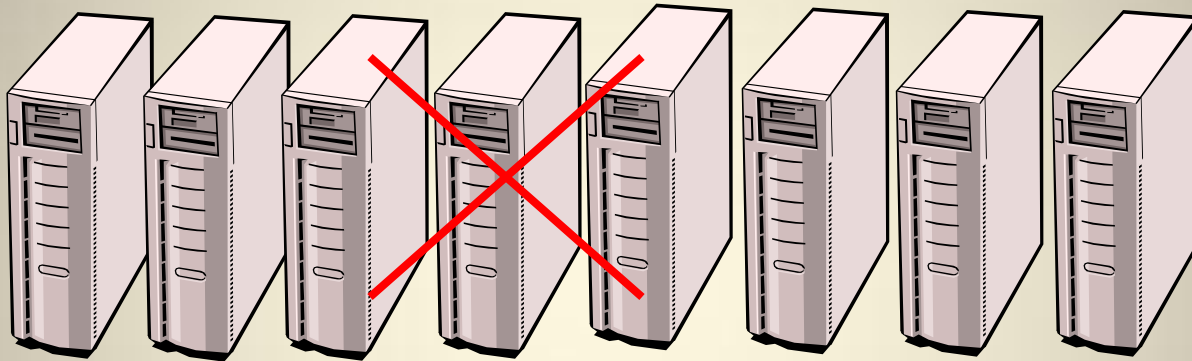


$x$            $y$            $x+ y$

# Distributed data storage



$x$        $y$        $x+y$
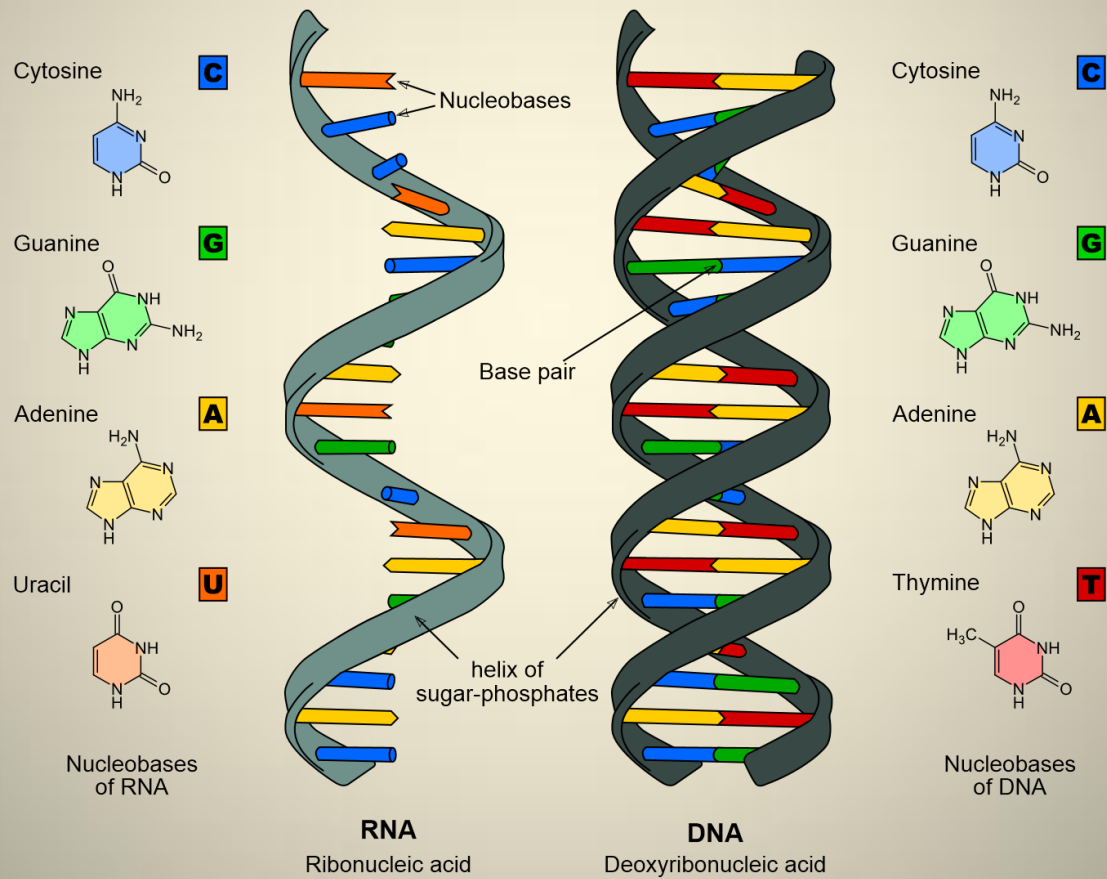
# Distributed data storage

# Distributed data storage



- Classical error-correcting codes can be employed
- Local correction is needed (using few other servers) to facilitate the correction

# DNA Analysis

# String Reconstruction Problem

- Four amino acids: A, F, G, C
- The composition of each protein can be deduced from its weight
- Each protein-sequence bond is cut independently with the same probability

AFGCCGA

CCG

GCCA

AFGC

CGA

Acharya, Das, Milenkovic, Orlitsky, and Pan '2011

# String Reconstruction Problem

- Binary alphabet {0,1}

0010011

0010

100

011

001

Acharya, Das, Milenkovic, Orlitsky, and Pan '2011