

Private information retrieval with codes

Julien Lavauzelle

LAGA, Université Paris 8 Vincennes – Saint-Denis

Seminario de Álgebra, Geometría algebraica y Singularidades
La Laguna

31/10/2023

1. Private information retrieval

2. PIR schemes with low computation and storage

- Transversal designs and codes

- A PIR scheme with transversal designs

- Collusion-resistant PIR schemes with weighted lifted codes

1. Private information retrieval

2. PIR schemes with low computation and storage

Transversal designs and codes

A PIR scheme with transversal designs

Collusion-resistant PIR schemes with weighted lifted codes

Private information retrieval (PIR):

Consider a set of files F_1, \dots, F_k , stored on a remote system.

One wants to retrieve one file F_i **privately**, that is, by hiding the value of i to the system.

Is it possible?

Private information retrieval (PIR):

Consider a set of files F_1, \dots, F_k , stored on a remote system.

One wants to retrieve one file F_i **privately**, that is, by hiding the value of i to the system.

Is it possible?

Applications: access to medical data, geoprivacy, ...

Private information retrieval (PIR):

Consider a set of files F_1, \dots, F_k , stored on a remote system.


One wants to retrieve one file F_i **privately**, that is, by hiding the value of i to the system.

Is it possible?

Applications: access to medical data, geoprivacy, ...


Trivial solution: download all files.

- ▶ perfect privacy: no information
- ▶ bad download rate...

Introduced in:  *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

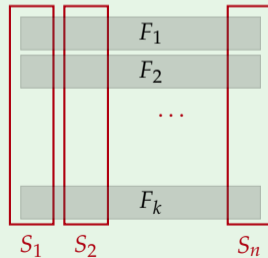
Generally, assume that n servers S_1, \dots, S_n store (in some way) the files F_1, \dots, F_k .

Definition of PIR


Introduced in:  *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

Generally, assume that n servers S_1, \dots, S_n store (in some way) the files F_1, \dots, F_k .

A **Private Information Retrieval protocol** is a set of algorithms $(Q, \mathcal{A}, \mathcal{R})$.



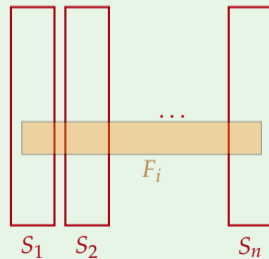
Definition of PIR


Introduced in:  *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

Generally, assume that n servers S_1, \dots, S_n store (in some way) the files F_1, \dots, F_k .

A **Private Information Retrieval protocol** is a set of algorithms $(Q, \mathcal{A}, \mathcal{R})$.

In order to retrieve file F_i :



Introduced in:  *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

Generally, assume that n servers S_1, \dots, S_n store (in some way) the files F_1, \dots, F_k .

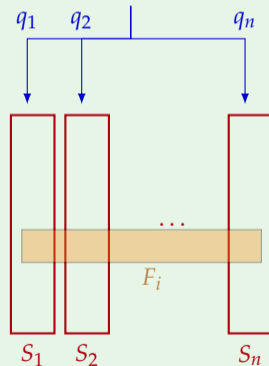
A **Private Information Retrieval protocol** is a set of algorithms $(\mathcal{Q}, \mathcal{A}, \mathcal{R})$.


In order to retrieve file F_i :

1. **Query generation:**

$$q := (q_1, \dots, q_n) \leftarrow \mathcal{Q}(i)$$

Send query q_j to server S_j .



Introduced in:  *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

Generally, assume that n servers S_1, \dots, S_n store (in some way) the files F_1, \dots, F_k .

A **Private Information Retrieval protocol** is a set of algorithms $(\mathcal{Q}, \mathcal{A}, \mathcal{R})$.

In order to retrieve file F_i :

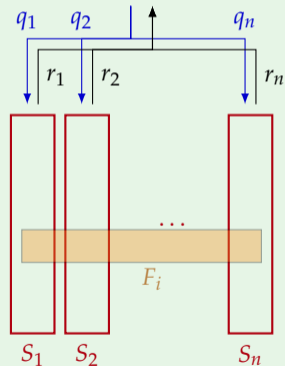
1. **Query generation:**


$$q := (q_1, \dots, q_n) \leftarrow \mathcal{Q}(i)$$

Send query q_j to server S_j .

2. **Response:** server S_j computes and sends back a response

$$r_j := \mathcal{A}(q_j, F_{|S_j})$$



Introduced in:  *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

Generally, assume that n servers S_1, \dots, S_n store (in some way) the files F_1, \dots, F_k .

A **Private Information Retrieval protocol** is a set of algorithms $(\mathcal{Q}, \mathcal{A}, \mathcal{R})$.

In order to retrieve file F_i :

1. **Query generation:**

$$q := (q_1, \dots, q_n) \leftarrow \mathcal{Q}(i)$$

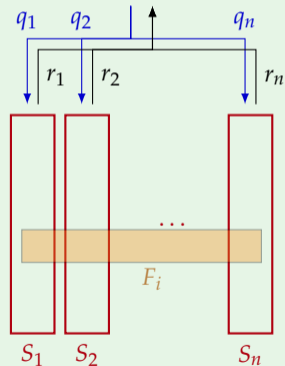
Send query q_j to server S_j .

2. **Response:** server S_j computes and sends back a response

$$r_j := \mathcal{A}(q_j, F_{|S_j})$$

3. **Local reconstruction** of the desired file:

$$F_i = \mathcal{R}(q, r, i).$$

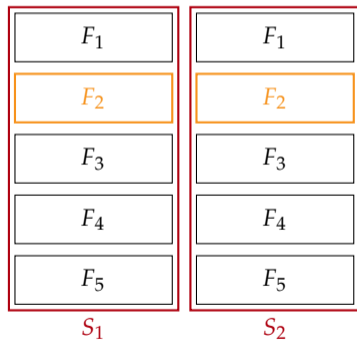


Example

$n = 2$ servers storing a replica of $k = 5$ files

F_1, \dots, F_5

Goal: retrieve file F_2 .

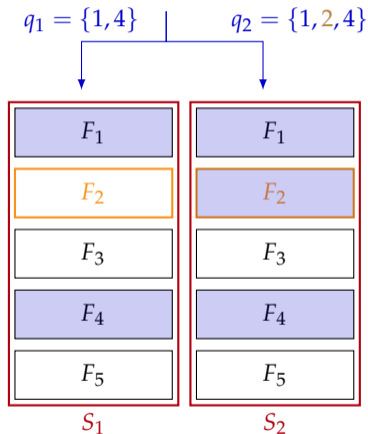


Example

$n = 2$ servers storing a replica of $k = 5$ files
 F_1, \dots, F_5

Goal: retrieve file F_2 .

1. **Query generation.** Pick at random a subset $I \subseteq \{1, \dots, 5\}$, and define:
 - query $q_1 = I$ to server S_1
 - query $q_2 = I \Delta \{2\}$ to server S_2



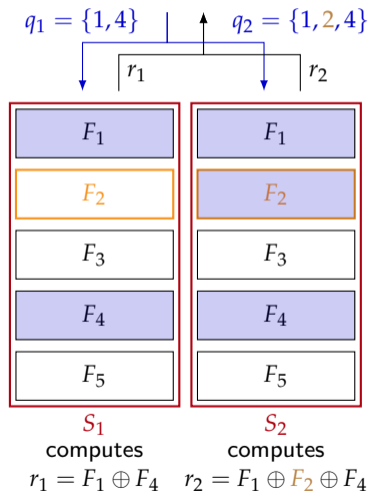
Example

$n = 2$ servers storing a replica of $k = 5$ files

F_1, \dots, F_5

Goal: retrieve file F_2 .

- Query generation.** Pick at random a subset $I \subseteq \{1, \dots, 5\}$, and define:
 - query $q_1 = I$ to server S_1
 - query $q_2 = I \Delta \{2\}$ to server S_2
- Responses.** Each server receives a subset $J \subseteq \{1, \dots, 5\}$, and computes the XOR (=bitwise sum) of files indexed by J .

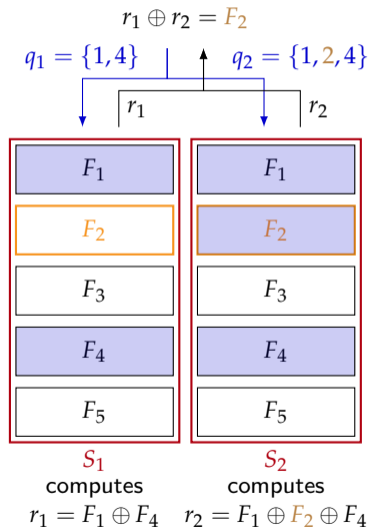


Example

$n = 2$ servers storing a replica of $k = 5$ files
 F_1, \dots, F_5

Goal: retrieve file F_2 .

- Query generation.** Pick at random a subset $I \subseteq \{1, \dots, 5\}$, and define:
 - query $q_1 = I$ to server S_1
 - query $q_2 = I \Delta \{2\}$ to server S_2
- Responses.** Each server receives a subset $J \subseteq \{1, \dots, 5\}$, and computes the XOR (=bitwise sum) of files indexed by J .
- Reconstruction.** One gets file F_2 by XORing the two responses.



Example

$n = 2$ servers storing a replica of $k = 5$ files
 F_1, \dots, F_5

Goal: retrieve file F_2 .

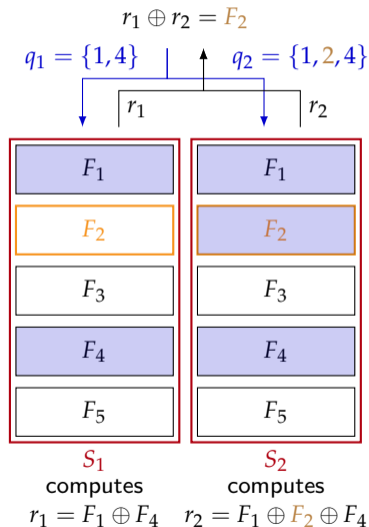
- Query generation.** Pick at random a subset $I \subseteq \{1, \dots, 5\}$, and define:
 - query $q_1 = I$ to server S_1
 - query $q_2 = I \Delta \{2\}$ to server S_2
- Responses.** Each server receives a subset $J \subseteq \{1, \dots, 5\}$, and computes the XOR (=bitwise sum) of files indexed by J .
- Reconstruction.** One gets file F_2 by XORing the two responses.

Upload: $2 \times 5 = 10$ bits to transmit to the servers

Download: $2|F|$ bits to receive from the servers

Server computation: XOR of $\frac{5}{2}|F|$ bits in average

Client computation: XOR of $2|F|$ bits



The adversary: a **collusion of servers** = a subset of servers $\{S_j : j \in T\}$, where $T \subset [1, n]$, which exchange information about queries.

$$t := \max\{|T|, T \subseteq [1, n] \text{ is a collusion}\} \geq 1$$

The adversary: a **collusion of servers** = a subset of servers $\{S_j : j \in T\}$, where $T \subset [1, n]$, which exchange information about queries.

$$t := \max\{|T|, T \subseteq [1, n] \text{ is a collusion}\} \geq 1$$

- **Information-theoretic (IT) privacy:**

$$I(i; \mathbf{q}_{|T}) = 0, \quad \forall T \subseteq [1, n], |T| \leq t.$$

- **Computational privacy:** by varying the index i , distributions of queries $\mathbf{q}_{|T} = \mathcal{Q}(i)_{|T}$ are computationally indistinguishable.

The adversary: a **collusion of servers** = a subset of servers $\{S_j : j \in T\}$, where $T \subset [1, n]$, which exchange information about queries.

$$t := \max\{|T|, T \subseteq [1, n] \text{ is a collusion}\} \geq 1$$

- **Information-theoretic (IT) privacy:**

$$I(i; \mathbf{q}_{|T}) = 0, \quad \forall T \subseteq [1, n], |T| \leq t.$$

- **Computational privacy:** by varying the index i , distributions of queries $\mathbf{q}_{|T} = \mathcal{Q}(i)_{|T}$ are computationally indistinguishable.

Theorem [CGKS95, CG97]. If $t = n$ (in particular if $n = 1$ server), then:

- for IT privacy, **no better solution than full download**,
- **computational** privacy is possible, but remains **expensive** as of now.

We focus on **IT-privacy**
(hence we need $n \geq 2$ servers)

We focus on **IT-privacy**
(hence we need $n \geq 2$ servers)

Parameters to be taken into account:

- **communication** complexity (upload and download)

We focus on **IT-privacy**
(hence we need $n \geq 2$ servers)

Parameters to be taken into account:

- **communication** complexity (upload and download)
- **computation** complexity (client and servers)

We focus on **IT-privacy**
(hence we need $n \geq 2$ servers)

Parameters to be taken into account:

- **communication** complexity (upload and download)
- **computation** complexity (client and servers)
- global **server storage** overhead

We focus on **IT-privacy**
(hence we need $n \geq 2$ servers)

Parameters to be taken into account:

- **communication** complexity (upload and download)
- **computation** complexity (client and servers)
- global **server storage** overhead
- maximum size of collusions (t)

We focus on **IT-privacy**
(hence we need $n \geq 2$ servers)

Parameters to be taken into account:

- **communication** complexity (upload and download)
- **computation** complexity (client and servers)
- global **server storage** overhead
- maximum size of collusions (t)

Several possible **settings**:

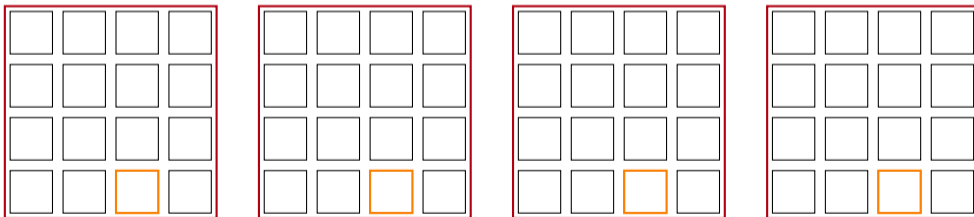
- **replicated** database vs. **coded** database
- unresponsive or **byzantine** servers
- small entries vs. large entries
- bounded vs. unbounded number of entries in the database
- dynamic database vs. static database

📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

16 files are replicated over 4 servers.

Files are indexed by pairs $(i, j) \in \{1, 2, 3, 4\}^2$

Assume one wants to retrieve file $F_{4,3}$.



📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

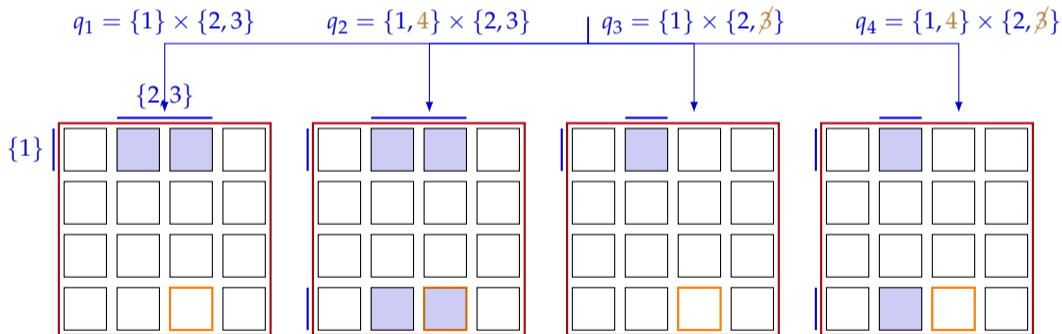
16 files are replicated over 4 servers.

Files are indexed by pairs $(i, j) \in \{1, 2, 3, 4\}^2$

Assume one wants to retrieve file $F_{4,3}$.

Queries are Cartesian products $I \times J$.

We add/remove indices i and j from I and J depending on the server.



📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

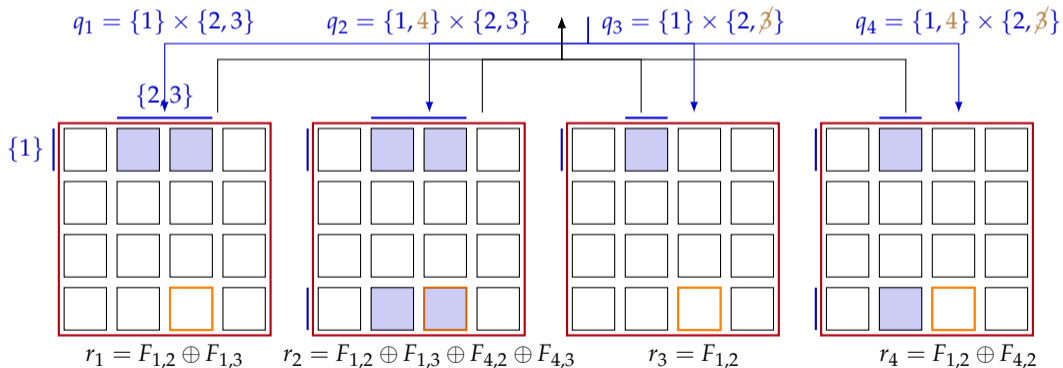
16 files are replicated over 4 servers.

Files are indexed by pairs $(i, j) \in \{1, 2, 3, 4\}^2$

Assume one wants to retrieve file $F_{4,3}$.

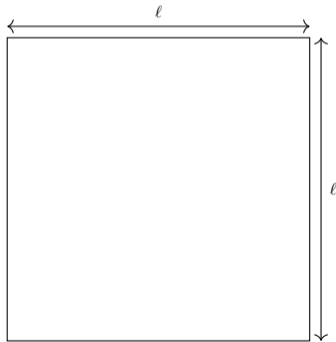
Queries are Cartesian products $I \times J$.

We add/remove indices i and j from I and J depending on the server.



📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

1st generalization: consider k files $F_{i,j}$ where $(i,j) \in [1, \ell]^2$ and $k = \ell^2$.
Files are replicated over $n = 4$ servers $S_{00}, S_{01}, S_{10}, S_{11}$.

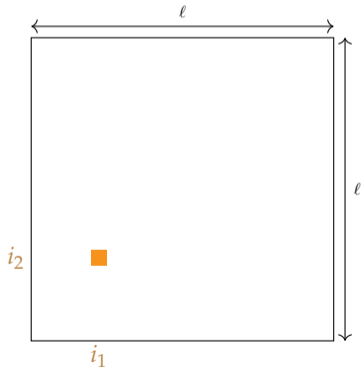


📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

1st generalization: consider k files $F_{i,j}$ where $(i,j) \in [1, \ell]^2$ and $k = \ell^2$.

Files are replicated over $n = 4$ servers $S_{00}, S_{01}, S_{10}, S_{11}$.

Goal: retrieve F_{i_1, i_2} , for $1 \leq i_1, i_2 \leq \ell$.

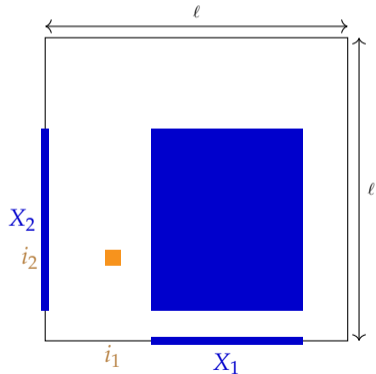


📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

1st generalization: consider k files $F_{i,j}$ where $(i,j) \in [1,\ell]^2$ and $k = \ell^2$.

Files are replicated over $n = 4$ servers $S_{00}, S_{01}, S_{10}, S_{11}$.

Goal: retrieve F_{i_1, i_2} , for $1 \leq i_1, i_2 \leq \ell$.



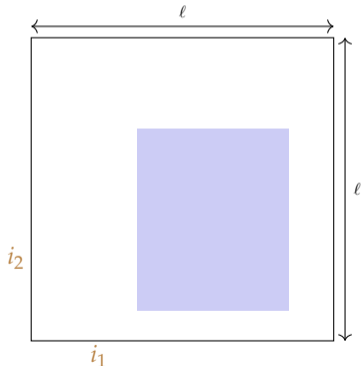
- Query generation:** pick at random two subsets X_1, X_2 of $[1, \ell]$. Then send:

📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

1st generalization: consider k files $F_{i,j}$ where $(i,j) \in [1, \ell]^2$ and $k = \ell^2$.

Files are replicated over $n = 4$ servers $S_{00}, S_{01}, S_{10}, S_{11}$.

Goal: retrieve F_{i_1, i_2} , for $1 \leq i_1, i_2 \leq \ell$.



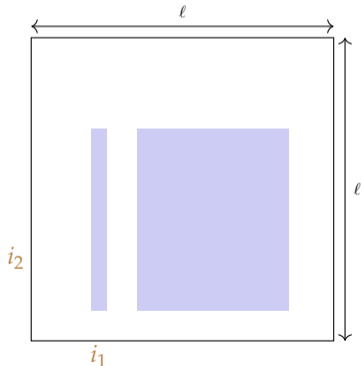
- Query generation:** pick at random two subsets X_1, X_2 of $[1, \ell]$. Then send:
 - (X_1, X_2) to server S_{00} ,

📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

1st generalization: consider k files $F_{i,j}$ where $(i,j) \in [1, \ell]^2$ and $k = \ell^2$.

Files are replicated over $n = 4$ servers $S_{00}, S_{01}, S_{10}, S_{11}$.

Goal: retrieve F_{i_1, i_2} , for $1 \leq i_1, i_2 \leq \ell$.



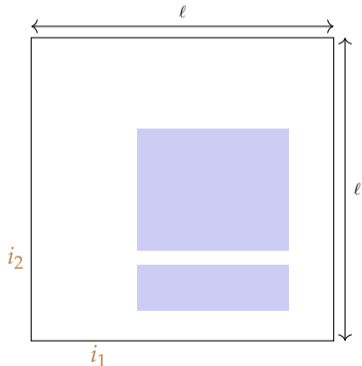
- Query generation:** pick at random two subsets X_1, X_2 of $[1, \ell]$. Then send:
 - (X_1, X_2) to server S_{00} ,
 - $(X_1 \Delta \{i_1\}, X_2)$ to server S_{10} ,

📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

1st generalization: consider k files $F_{i,j}$ where $(i,j) \in [1, \ell]^2$ and $k = \ell^2$.

Files are replicated over $n = 4$ servers $S_{00}, S_{01}, S_{10}, S_{11}$.

Goal: retrieve F_{i_1, i_2} , for $1 \leq i_1, i_2 \leq \ell$.



1. **Query generation:** pick at random two subsets X_1, X_2 of $[1, \ell]$. Then send:

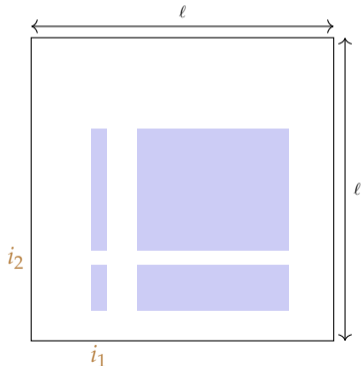
- (X_1, X_2) to server S_{00} ,
- $(X_1 \Delta \{i_1\}, X_2)$ to server S_{10} ,
- $(X_1, X_2 \Delta \{i_2\})$ to server S_{01} ,

📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

1st generalization: consider k files $F_{i,j}$ where $(i,j) \in [1, \ell]^2$ and $k = \ell^2$.

Files are replicated over $n = 4$ servers $S_{00}, S_{01}, S_{10}, S_{11}$.

Goal: retrieve F_{i_1, i_2} , for $1 \leq i_1, i_2 \leq \ell$.



1. **Query generation:** pick at random two subsets X_1, X_2 of $[1, \ell]$. Then send:

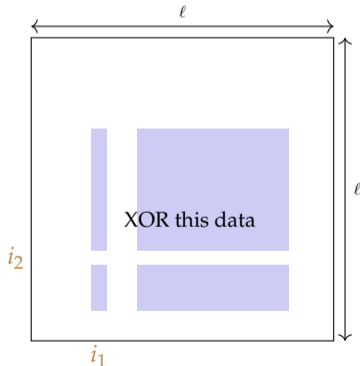
- (X_1, X_2) to server S_{00} ,
- $(X_1 \Delta \{i_1\}, X_2)$ to server S_{10} ,
- $(X_1, X_2 \Delta \{i_2\})$ to server S_{01} ,
- $(X_1 \Delta \{i_1\}, X_2 \Delta \{i_2\})$ to server S_{11} .

📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

1st generalization: consider k files F_{ij} where $(i, j) \in [1, \ell]^2$ and $k = \ell^2$.

Files are replicated over $n = 4$ servers $S_{00}, S_{01}, S_{10}, S_{11}$.

Goal: retrieve F_{i_1, i_2} , for $1 \leq i_1, i_2 \leq \ell$.



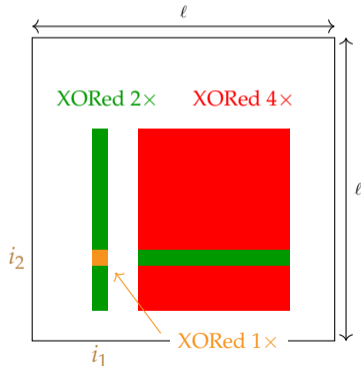
- Query generation:** pick at random two subsets X_1, X_2 of $[1, \ell]$. Then send:
 - (X_1, X_2) to server S_{00} ,
 - $(X_1 \Delta \{i_1\}, X_2)$ to server S_{10} ,
 - $(X_1, X_2 \Delta \{i_2\})$ to server S_{01} ,
 - $(X_1 \Delta \{i_1\}, X_2 \Delta \{i_2\})$ to server S_{11} .
- Answers:** at reception of (Z_1, Z_2) , each server S_j computes $R_j = \bigoplus_{z \in Z_1 \times Z_2} F_z$ and sends back R_j .

📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

1st generalization: consider k files F_{ij} where $(i, j) \in [1, \ell]^2$ and $k = \ell^2$.

Files are replicated over $n = 4$ servers $S_{00}, S_{01}, S_{10}, S_{11}$.

Goal: retrieve F_{i_1, i_2} , for $1 \leq i_1, i_2 \leq \ell$.



- Query generation:** pick at random two subsets X_1, X_2 of $[1, \ell]$. Then send:
 - (X_1, X_2) to server S_{00} ,
 - $(X_1 \Delta \{i_1\}, X_2)$ to server S_{10} ,
 - $(X_1, X_2 \Delta \{i_2\})$ to server S_{01} ,
 - $(X_1 \Delta \{i_1\}, X_2 \Delta \{i_2\})$ to server S_{11} .
- Answers:** at reception of (Z_1, Z_2) , each server S_j computes $R_j = \bigoplus_{z \in Z_1 \times Z_2} F_z$ and sends back R_j .
- Reconstruction:** compute the XOR of the 4 files R_j and retrieves F_{i_1, i_2} .

Correct, and **secure** if no collusion.

Correct, and **secure** if no collusion.

Quantitative results. Assume the k files have same size $|F|$.

With $n = 4$ servers,

- ▶ **Communication:** $8\sqrt{k}$ uploaded bits, $4|F|$ downloaded bits,
- ▶ **Storage:** replication of all files over 4 servers,
- ▶ **Complexity:**
 - for each server: in average, XOR of $(\ell/2)^2 = k/4$ files
 - for the user: XOR of 4 files.

Correct, and **secure** if no collusion.

Quantitative results. Assume the k files have same size $|F|$.

With $n = 4$ servers,

- ▶ **Communication:** $8\sqrt{k}$ uploaded bits, $4|F|$ downloaded bits,
- ▶ **Storage:** replication of all files over 4 servers,
- ▶ **Complexity:**
 - for each server: in average, XOR of $(\ell/2)^2 = k/4$ files
 - for the user: XOR of 4 files.

Generalizable to $n = 2^b$ servers:

- ▶ **Communication:** $b2^b k^{1/b} = n \log(n) k^{1/\log(n)}$ uploaded bits, $n|F|$ downloaded bits,
- ▶ **Storage:** replication of all files over n servers,
- ▶ **Complexity:**
 - ▶ for each server: in average, XOR of k/n files
 - ▶ for the user: XOR of n files.

- 1995: first definition [CGKS95]
- 2000: reduction from smooth locally decodable codes [KT00]
- 2000-10's: many improvements
 - ▶ PIR with 3 servers and subpolynomial communication [Yek08, Efr09]
 - ▶ PIR with 2 servers and subpolynomial communication [DG16]
 - ▶ lower storage overhead with *PIR codes* [FVY15]
- 2016-now: capacity-achieving schemes, schemes dedicated to storage systems
 - ▶ capacity of PIR [SJ17, BU18]
 - ▶ (nearly) capacity-achieving schemes [SRR14, CHY15, TR16, ...]

1. Private information retrieval

2. PIR schemes with low computation and storage

- Transversal designs and codes

- A PIR scheme with transversal designs

- Collusion-resistant PIR schemes with weighted lifted codes

Previous scheme:

- ▶ moderate communication complexity
- ▶ computationally inefficient (linear in $|F|$)
- ▶ huge storage overhead (replicas of $|F|$)

Previous scheme:

- ▶ moderate communication complexity
- ▶ computationally inefficient (linear in $|F|$)
- ▶ huge storage overhead (replicas of $|F|$)

Our goal:

- ▶ moderate communication complexity
- ▶ **optimal computation** (one read for each server)
- ▶ **smaller storage overhead** thanks to a pre-encoding and a distribution of the database

Previous scheme:

- ▶ moderate communication complexity
- ▶ computationally inefficient (linear in $|F|$)
- ▶ huge storage overhead (replicas of $|F|$)

Our goal:

- ▶ moderate communication complexity
- ▶ **optimal computation** (one read for each server)
- ▶ **smaller storage overhead** thanks to a pre-encoding and a distribution of the database

Tools: coding theory and combinatorics

- ▶ transversal designs and associated codes,
- ▶ “lifted” codes.

1. Private information retrieval

2. PIR schemes with low computation and storage

Transversal designs and codes

A PIR scheme with transversal designs

Collusion-resistant PIR schemes with weighted lifted codes

A (linear) **code** \mathcal{C} is a k -dimensional subspace of \mathbb{F}_q^n .

Any code admits **parity-check matrices** $H \in \mathbb{F}_q^{(n-k) \times n}$ such that

$$\mathcal{C} = \{c \in \mathbb{F}_q^n \mid Hc = \mathbf{0}\}.$$

A (linear) **code** \mathcal{C} is a k -dimensional subspace of \mathbb{F}_q^n .

Any code admits **parity-check matrices** $H \in \mathbb{F}_q^{(n-k) \times n}$ such that

$$\mathcal{C} = \{c \in \mathbb{F}_q^n \mid Hc = \mathbf{0}\}.$$

Terminology:

- $h \in \text{RowSpan}(H)$ is a **parity-check equation** for \mathcal{C} .
- **Support:** $\text{supp}(h) := \{i \in \{1, \dots, n\}, h_i \neq 0\}$
- **Weight:** $\text{wt}(h) := |\text{supp}(h)|$.

A (linear) **code** \mathcal{C} is a k -dimensional subspace of \mathbb{F}_q^n .

Any code admits **parity-check matrices** $H \in \mathbb{F}_q^{(n-k) \times n}$ such that

$$\mathcal{C} = \{c \in \mathbb{F}_q^n \mid Hc = \mathbf{0}\}.$$

Terminology:

- $h \in \text{RowSpan}(H)$ is a **parity-check equation** for \mathcal{C} .
- **Support:** $\text{supp}(h) := \{i \in \{1, \dots, n\}, h_i \neq 0\}$
- **Weight:** $\text{wt}(h) := |\text{supp}(h)|$.

Important remark. If $r = \text{wt}(h)$ and $i \in \text{supp}(h)$, then for every $c \in \mathcal{C}$, one can recover $c_i \in \mathbb{F}_q$ by accessing at most $r - 1$ other coordinates c_j of the codeword c :

$$c_i = -\frac{1}{h_i} \sum_{j \in \text{supp}(h) \setminus \{i\}} h_j c_j$$

In that case we call $\text{supp}(h) \setminus \{i\}$ a **helper set** for i .

Our goal. Design a code $\mathcal{C} \subseteq \mathbb{F}_q^n$ such that, for every $i \in \{1, \dots, n\}$, there exists a set of helper sets which uniformly covers $\{1, \dots, n\}$.

Our goal. Design a code $\mathcal{C} \subseteq \mathbb{F}_q^n$ such that, for every $i \in \{1, \dots, n\}$, there exists a set of helper sets which uniformly covers $\{1, \dots, n\}$.

- ▶ This will provide a way to recover c_i by accessing uniformly at random other coordinates of c .
- ▶ Querying these "random coordinates" will leak no information about i to the servers.

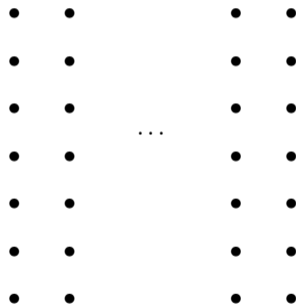
Our goal. Design a code $\mathcal{C} \subseteq \mathbb{F}_q^n$ such that, for every $i \in \{1, \dots, n\}$, there exists a set of helper sets which uniformly covers $\{1, \dots, n\}$.

- ▶ This will provide a way to recover c_i by accessing uniformly at random other coordinates of c .
- ▶ Querying these "random coordinates" will leak no information about i to the servers.

Let's do this with **combinatorics**.

A transversal design $\text{TD}(n, s) = (X, \mathcal{B}, \mathcal{G})$ is given by:

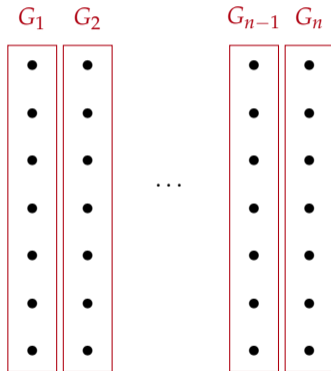
- ▶ X a set of *points*, $|X| = N = ns$,



A transversal design $\text{TD}(n, s) = (X, \mathcal{B}, \mathcal{G})$ is given by:

- ▶ X a set of *points*, $|X| = N = ns$,
- ▶ a partition of X into subsets $\mathcal{G} = \{G_j\}_{1 \leq j \leq n}$ called *groups*:

$$X = \bigsqcup_{j=1}^n G_j \text{ and } |G_j| = s,$$

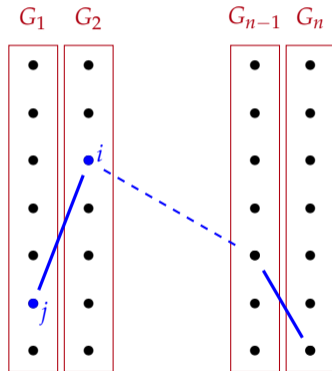


A transversal design $\text{TD}(n, s) = (X, \mathcal{B}, \mathcal{G})$ is given by:

- ▶ X a set of *points*, $|X| = N = ns$,
- ▶ a partition of X into subsets $\mathcal{G} = \{G_j\}_{1 \leq j \leq n}$ called *groups*:

$$X = \bigsqcup_{j=1}^n G_j \text{ and } |G_j| = s,$$

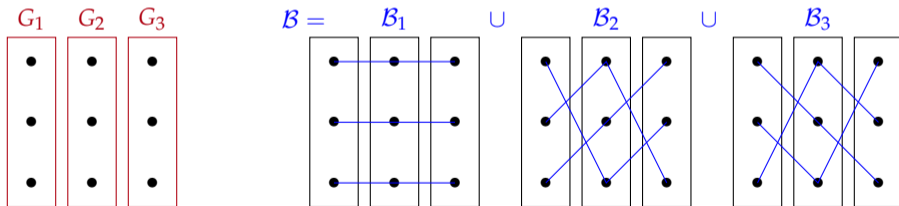
- ▶ a set of subsets of X , “incident” to \mathcal{G} , called *blocks* $B \in \mathcal{B}$:
 - $B \subset X$ and $|B| = n$
 - for all $\{i, j\} \subset X$, the pair $\{i, j\}$ lie either in a single group $G \in \mathcal{G}$, or in a unique block $B \in \mathcal{B}$



Example: a TD(3,3)

An example for a TD(3,3) :

- $ns = 9$ points
- $s = 3$ groups G_1, G_2, G_3 of size 3
- $ns = 9$ blocks of $n = 3$ points, partitionned into 3 parallel classes B_1, B_2, B_3



Let \mathcal{T} be a transversal design $\text{TD}(n, s) = (X, \mathcal{B}, \mathcal{G})$ with points $X = \{x_1, \dots, x_{ns}\}$, blocks $B = \{B_1, \dots, B_{ns}\}$ and groups $\mathcal{G} = \{G_1, \dots, G_n\}$.

Its **incidence matrix** M , of size $|\mathcal{B}| \times |X| = ns \times ns$, is defined by:

$$M_{i,j} = \begin{cases} 1 & \text{if } x_j \in B_i \\ 0 & \text{otherwise.} \end{cases}$$

Let \mathcal{T} be a transversal design $\text{TD}(n, s) = (X, \mathcal{B}, \mathcal{G})$ with points $X = \{x_1, \dots, x_{ns}\}$, blocks $B = \{B_1, \dots, B_{ns}\}$ and groups $\mathcal{G} = \{G_1, \dots, G_n\}$.

Its **incidence matrix** M , of size $|\mathcal{B}| \times |X| = ns \times ns$, is defined by:

$$M_{i,j} = \begin{cases} 1 & \text{if } x_j \in B_i \\ 0 & \text{otherwise.} \end{cases}$$

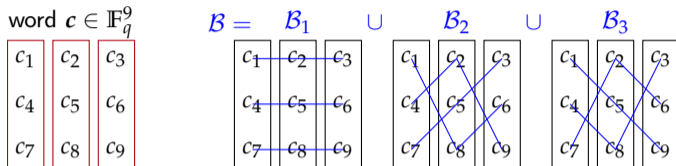
Definition. The **linear code** \mathcal{C} based on \mathcal{T} over \mathbb{F}_q is the \mathbb{F}_q -linear code having M as a parity-check matrix (i.e. \mathcal{C}^\perp is generated by M).

- $\text{length}(\mathcal{C}) = |X| = ns$,
- $\dim(\mathcal{C}) = \dim(\ker M)$,
- every block $B \in \mathcal{B}$ gives a parity-check equation $\mathbf{h} \in \mathcal{C}^\perp$, such that

$$\text{supp}(\mathbf{h}) = B \quad \text{and} \quad \text{wt}(\mathbf{h}|_{G_j}) = 1, \quad \forall j = 1, \dots, n$$

Example

The transversal design $TD(3, 3)$ represented by:

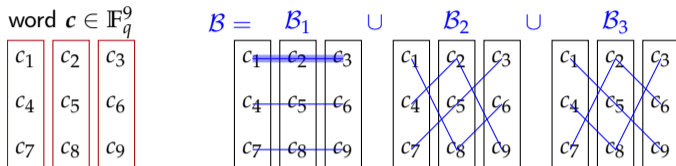


gives a code with the following parity-check matrix:

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

Example

The transversal design $TD(3, 3)$ represented by:

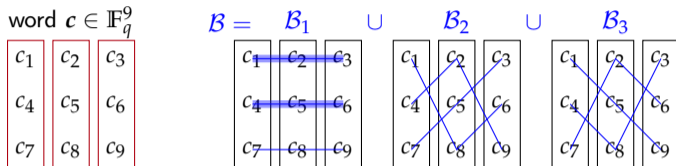


gives a code with the following parity-check matrix:

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

Example

The transversal design $TD(3, 3)$ represented by:

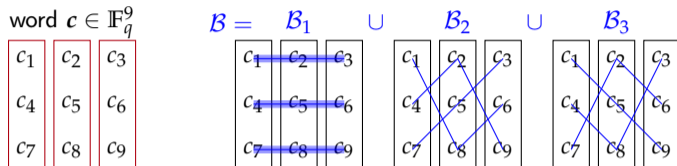


gives a code with the following parity-check matrix:

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

Example

The transversal design $TD(3,3)$ represented by:

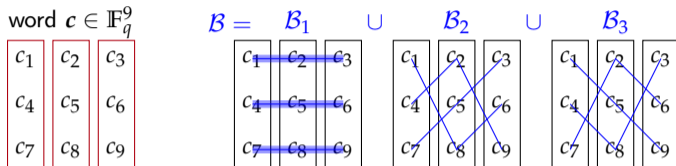


gives a code with the following parity-check matrix:

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

Example

The transversal design $TD(3,3)$ represented by:



gives a code with the following parity-check matrix:

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

Dimension of the code?

- ▶ depends on the characteristic,
- ▶ for instance, over \mathbb{F}_3 , we have $\text{rk}(H) = 6 \implies \dim(\mathcal{C}) = 3$.

1. Private information retrieval

2. PIR schemes with low computation and storage

Transversal designs and codes

A PIR scheme with transversal designs

Collusion-resistant PIR schemes with weighted lifted codes

 *Private Information Retrieval from Transversal Designs*. L.. IEEE-TIT. 2019.

Let $\mathcal{C} \subseteq \mathbb{F}_q^N$ be a code based on a TD(n, s), with $N = ns$.

📄 *Private Information Retrieval from Transversal Designs*. L.. IEEE-TIT. 2019.

Let $\mathcal{C} \subseteq \mathbb{F}_q^N$ be a code based on a TD(n, s), with $N = ns$.

- **Initialisation.** Encode files $(F_1, \dots, F_k) \mapsto c \in \mathcal{C}$, and upload $c|_{G_j}$ on server S_j .

G_1	G_2	G_3
c_1	c_4	c_7
c_2	c_5	c_8
c_3	c_6	c_9

 *Private Information Retrieval from Transversal Designs*. L.. IEEE-TIT. 2019.

Let $\mathcal{C} \subseteq \mathbb{F}_q^N$ be a code based on a TD(n, s), with $N = ns$.

- **Initialisation.** Encode files $(F_1, \dots, F_k) \mapsto c \in \mathcal{C}$, and upload $c|_{G_j}$ on server S_j .
- **To recover** $F_i = c_i$, with $i \in X$:

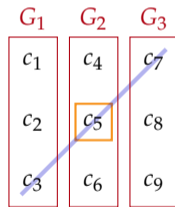
G_1	G_2	G_3
c_1	c_4	c_7
c_2	c_5	c_8
c_3	c_6	c_9

📄 *Private Information Retrieval from Transversal Designs*. L.. IEEE-TIT. 2019.

Let $\mathcal{C} \subseteq \mathbb{F}_q^N$ be a code based on a $\text{TD}(n, s)$, with $N = ns$.

- **Initialisation.** Encode files $(F_1, \dots, F_k) \mapsto c \in \mathcal{C}$, and upload $c|_{G_j}$ on server S_j .
- **To recover** $F_i = c_i$, with $i \in X$:
 1. User randomly picks a block $B \in \mathcal{B}$ containing i .
Then, user defines queries:

$$q_j = Q(i)_j := \begin{cases} \text{unique point in } B \cap G_j & \text{if } i \notin G_j \\ \text{a random point in } G_j & \text{otherwise.} \end{cases}$$

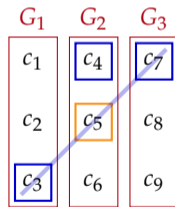


 *Private Information Retrieval from Transversal Designs*. L.. IEEE-TIT. 2019.

Let $\mathcal{C} \subseteq \mathbb{F}_q^N$ be a code based on a $\text{TD}(n, s)$, with $N = ns$.

- **Initialisation.** Encode files $(F_1, \dots, F_k) \mapsto c \in \mathcal{C}$, and upload $c|_{G_j}$ on server S_j .
- **To recover** $F_i = c_i$, with $i \in X$:
 1. User randomly picks a block $B \in \mathcal{B}$ containing i .
Then, user defines queries:

$$q_j = Q(i)_j := \begin{cases} \text{unique point in } B \cap G_j & \text{if } i \notin G_j \\ \text{a random point in } G_j & \text{otherwise.} \end{cases}$$



📄 *Private Information Retrieval from Transversal Designs*. L.. IEEE-TIT. 2019.

Let $\mathcal{C} \subseteq \mathbb{F}_q^N$ be a code based on a $\text{TD}(n, s)$, with $N = ns$.

- **Initialisation.** Encode files $(F_1, \dots, F_k) \mapsto c \in \mathcal{C}$, and upload $c|_{G_j}$ on server S_j .

- **To recover** $F_i = c_i$, with $i \in X$:

1. User randomly picks a block $B \in \mathcal{B}$ containing i .

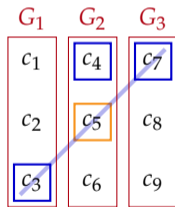
Then, user defines queries:

$$q_j = \mathcal{Q}(i)_j := \begin{cases} \text{unique point in } B \cap G_j & \text{if } i \notin G_j \\ \text{a random point in } G_j & \text{otherwise.} \end{cases}$$

2. Each server S_j sends back c_{q_j}

3. User recovers

$$c_i = - \sum_{j: i \notin G_j} c_{q_j} = - \sum_{b \in B \setminus \{i\}} c_b$$



Theorem. This PIR protocol is information-theoretically private.

Proof:

- the only server which holds F_i received a random query;
- for each other server S_j , query q_j gives no information on the block B which has been picked \Rightarrow no information leaks on i .

Theorem. This PIR protocol is information-theoretically private.

Proof:

- the only server which holds F_i received a random query;
- for each other server S_j , query q_j gives no information on the block B which has been picked \Rightarrow no information leaks on i .

Features.

- ▶ communication complexity: $n \log s$ uploaded bits, $n \log q$ downloaded bits

Theorem. This PIR protocol is information-theoretically private.

Proof:

- the only server which holds F_i received a random query;
- for each other server S_j , query q_j gives no information on the block B which has been picked \Rightarrow no information leaks on i .

Features.

- ▶ communication complexity: $n \log s$ uploaded bits, $n \log q$ downloaded bits
- ▶ computational complexity:
 - ▶ **only 1 read for each server (optimal)**
 - ▶ $\leq n$ additions over \mathbb{F}_q for the user

Theorem. This PIR protocol is information-theoretically private.

Proof:

- the only server which holds F_i received a random query;
- for each other server S_j , query q_j gives no information on the block B which has been picked \Rightarrow no information leaks on i .

Features.

- ▶ communication complexity: $n \log s$ uploaded bits, $n \log q$ downloaded bits
- ▶ computational complexity:
 - ▶ **only 1 read for each server (optimal)**
 - ▶ $\leq n$ additions over \mathbb{F}_q for the user
- ▶ storage overhead: $(ns - k) \log q$ bits, where $k = \dim(\mathcal{C})$

Theorem. This PIR protocol is information-theoretically private.

Proof:

- the only server which holds F_i received a random query;
- for each other server S_j , query q_j gives no information on the block B which has been picked \Rightarrow no information leaks on i .

Features.

- ▶ communication complexity: $n \log s$ uploaded bits, $n \log q$ downloaded bits
- ▶ computational complexity:
 - ▶ **only 1 read for each server (optimal)**
 - ▶ $\leq n$ additions over \mathbb{F}_q for the user
- ▶ storage overhead: $(ns - k) \log q$ bits, where $k = \dim(\mathcal{C})$

Question: transversal designs leading to large dimension codes?

An example: the **classical affine transversal design**:

- ▶ $X = \mathbb{F}_q^m$ for $m \geq 2$,
- ▶ \mathcal{G} a partition of X into q hyperplanes G_1, \dots, G_q ,
- ▶ $\mathcal{B} = \{\text{affine lines } L \text{ secant to each } G_j\}$.

The code has:

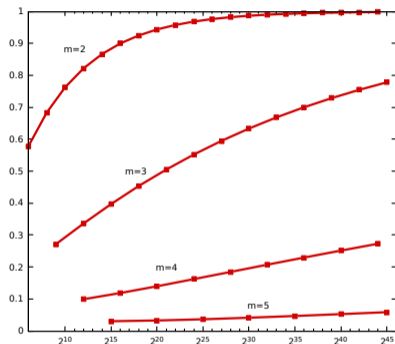
- length $ns = q^m$,
- "locality" $n = q$.

An example: the **classical affine transversal design**:

- ▶ $X = \mathbb{F}_q^m$ for $m \geq 2$,
- ▶ \mathcal{G} a partition of X into q hyperplanes G_1, \dots, G_q ,
- ▶ $\mathcal{B} = \{\text{affine lines } L \text{ secant to each } G_j\}$.

The code has:

- length $ns = q^m$,
- "locality" $n = q$.



Information rate of affine geometry TD-based code depending on the length.

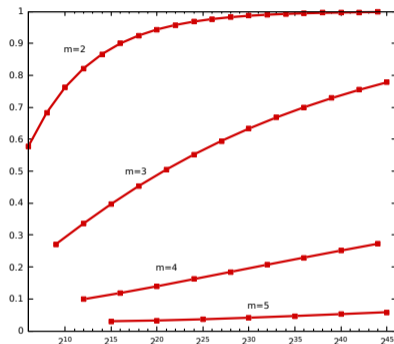
Each curve corresponds to a value of $m \in \{2, 3, 4, 5\}$.

An example: the **classical affine transversal design**:

- ▶ $X = \mathbb{F}_q^m$ for $m \geq 2$,
- ▶ \mathcal{G} a partition of X into q hyperplanes G_1, \dots, G_q ,
- ▶ $\mathcal{B} = \{\text{affine lines } L \text{ secant to each } G_j\}$.

The code has:

- length $ns = q^m$,
- "locality" $n = q$.



Information rate of affine geometry TD-based code depending on the length.

Each curve corresponds to a value of $m \in \{2, 3, 4, 5\}$.

Question: how to deal with collusions and errors?

1. Private information retrieval

2. PIR schemes with low computation and storage

Transversal designs and codes

A PIR scheme with transversal designs

Collusion-resistant PIR schemes with weighted lifted codes

We have seen a **combinatorial** construction of codes for PIR, using **transversal designs**.
Let's now see what we can do **algebraically**.

Definition. The (full-length) **Reed–Solomon code** of dimension k over \mathbb{F}_q is:

$$\text{RS}_q(k) := \{\text{ev}_{\mathbb{A}^1}(f) := (f(x_1), \dots, f(x_q)) \mid \deg(f) \leq k - 1\}.$$

Definition. The (full-length) **Reed–Solomon code** of dimension k over \mathbb{F}_q is:

$$\text{RS}_q(k) := \{\text{ev}_{\mathbb{A}^1}(f) := (f(x_1), \dots, f(x_q)) \mid \deg(f) \leq k - 1\}.$$

The code $\mathcal{C} = \text{RS}_q(k)$ is **MDS**: every codeword $c \in \mathcal{C}$ can be reconstructed from any k -subset of coordinates of c .

Definition. The (full-length) **Reed–Solomon code** of dimension k over \mathbb{F}_q is:

$$\text{RS}_q(k) := \{\text{ev}_{\mathbb{A}^1}(f) := (f(x_1), \dots, f(x_q)) \mid \deg(f) \leq k - 1\}.$$

The code $\mathcal{C} = \text{RS}_q(k)$ is **MDS**: every codeword $c \in \mathcal{C}$ can be reconstructed from any k -subset of coordinates of c .

Definition. The **Reed–Muller code** of order m and degree r over \mathbb{F}_q is:

$$\text{RM}_q(m, r) := \{\text{ev}_{\mathbb{A}^m}(f) \mid f \in \mathbb{F}_q[X] \text{ and } \deg(f) \leq r\}.$$

Definition. The (full-length) **Reed–Solomon code** of dimension k over \mathbb{F}_q is:

$$\text{RS}_q(k) := \{\text{ev}_{\mathbb{A}^1}(f) := (f(x_1), \dots, f(x_q)) \mid \deg(f) \leq k - 1\}.$$

The code $\mathcal{C} = \text{RS}_q(k)$ is **MDS**: every codeword $\mathbf{c} \in \mathcal{C}$ can be reconstructed from any k -subset of coordinates of \mathbf{c} .

Definition. The **Reed–Muller code** of order m and degree r over \mathbb{F}_q is:

$$\text{RM}_q(m, r) := \{\text{ev}_{\mathbb{A}^m}(f) \mid f \in \mathbb{F}_q[X] \text{ and } \deg(f) \leq r\}.$$

Reed–Muller codes have the following property:

$$\forall \mathbf{c} = \text{ev}_{\mathbb{A}^m}(f) \in \text{RM}_q(m, r) \quad \text{and} \quad \forall \text{ affine line } L \subset \mathbb{A}^m, \\ \text{ev}_{\mathbb{A}^1}(f|_L) \in \text{RS}_q(r + 1).$$

(where $f|_L$ is the lowest-degree univariate polynomial interpolating f over L)

Definition. The (full-length) **Reed–Solomon code** of dimension k over \mathbb{F}_q is:

$$\text{RS}_q(k) := \{\text{ev}_{\mathbb{A}^1}(f) := (f(x_1), \dots, f(x_q)) \mid \deg(f) \leq k - 1\}.$$

The code $\mathcal{C} = \text{RS}_q(k)$ is **MDS**: every codeword $c \in \mathcal{C}$ can be reconstructed from any k -subset of coordinates of c .

Definition. The **Reed–Muller code** of order m and degree r over \mathbb{F}_q is:

$$\text{RM}_q(m, r) := \{\text{ev}_{\mathbb{A}^m}(f) \mid f \in \mathbb{F}_q[X] \text{ and } \deg(f) \leq r\}.$$

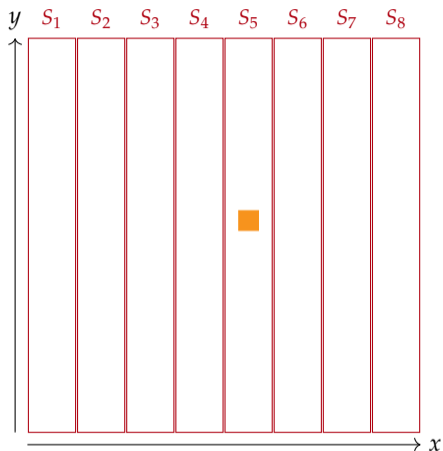
Reed–Muller codes have the following property:

$$\forall c = \text{ev}_{\mathbb{A}^m}(f) \in \text{RM}_q(m, r) \quad \text{and} \quad \forall \text{ affine line } L \subset \mathbb{A}^m, \\ \text{ev}_{\mathbb{A}^1}(f|_L) \in \text{RS}_q(r + 1).$$

(where $f|_L$ is the lowest-degree univariate polynomial interpolating f over L)

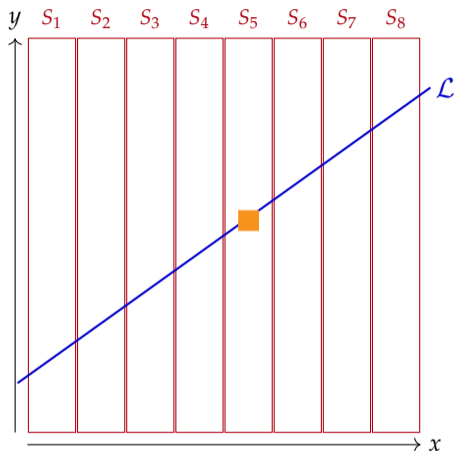
In particular, if $r \leq q - 2$, then $c_i = f(P_i)$ can be reconstructed by interpolating a polynomial of degree r on $q - 1$ other points of **any line** passing through P_i .

A PIR scheme based on Reed–Muller codes



Database F is encoded with $\text{RM}_q(m, r)$, then distributed across the servers

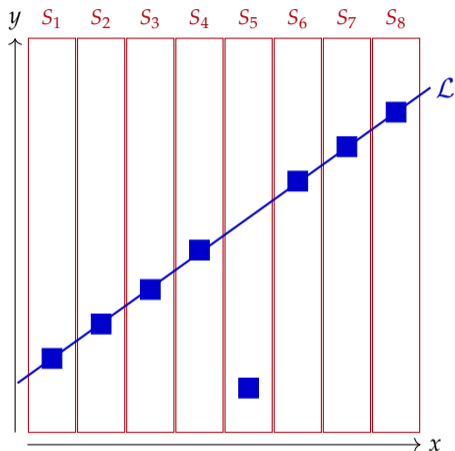
A PIR scheme based on Reed–Muller codes



Database F is encoded with $\text{RM}_q(m, r)$, then distributed across the servers

Assume one wants to extract file F_i .

A PIR scheme based on Reed–Muller codes

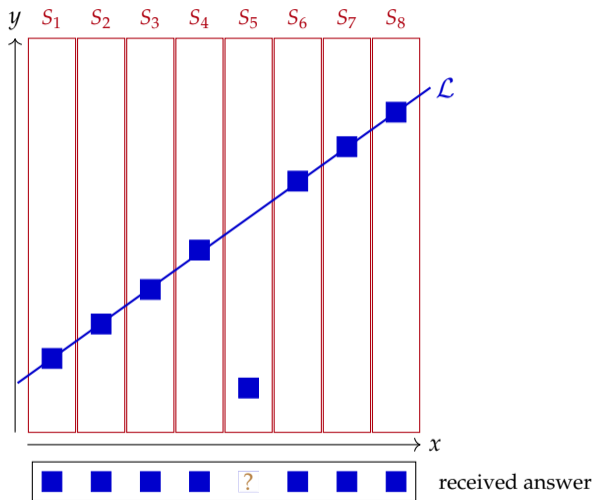


Database F is encoded with $\text{RM}_q(m, r)$, then distributed across the servers

Assume one wants to extract file F_i .

Pick a line L through i ,

A PIR scheme based on Reed–Muller codes

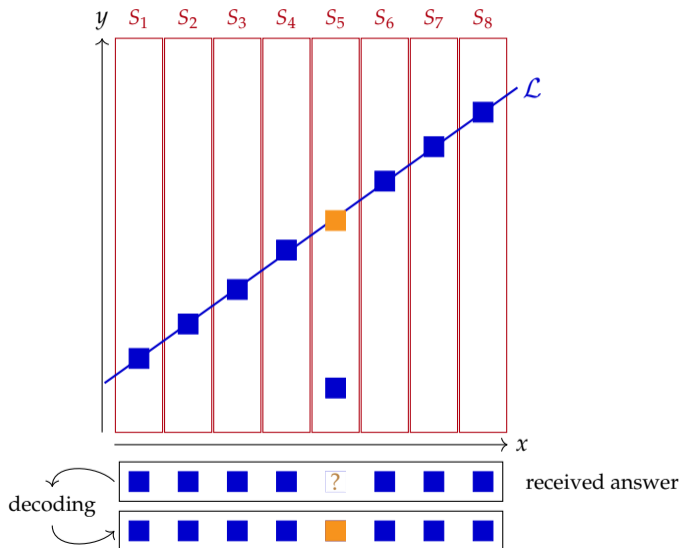


Database F is encoded with $\text{RM}_q(m, r)$, then distributed across the servers

Assume one wants to extract file F_i .

Pick a line L through i , query all files F_j except F_i ,

A PIR scheme based on Reed–Muller codes

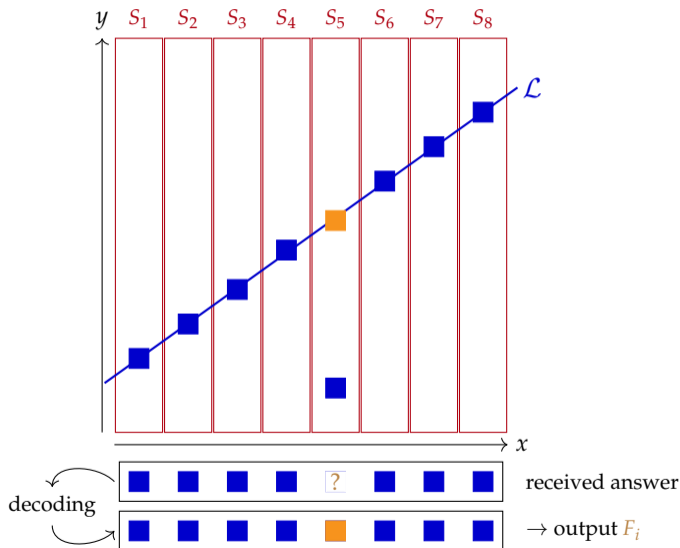


Database F is encoded with $\text{RM}_q(m, r)$, then distributed across the servers

Assume one wants to extract file F_i .

Pick a line L through i , query all files F_j except F_i , and interpolate the corresponding univariate polynomial.

A PIR scheme based on Reed–Muller codes

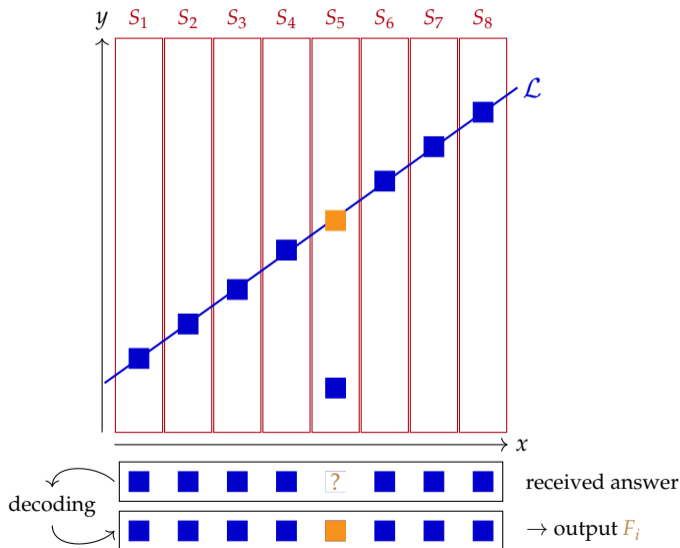


Database F is encoded with $\text{RM}_q(m, r)$, then distributed across the servers

Assume one wants to extract file F_i .

Pick a line L through i , query all files F_j except F_i , and interpolate the corresponding univariate polynomial.

A PIR scheme based on Reed–Muller codes



Database F is encoded with $\text{RM}_q(m, r)$, then distributed across the servers

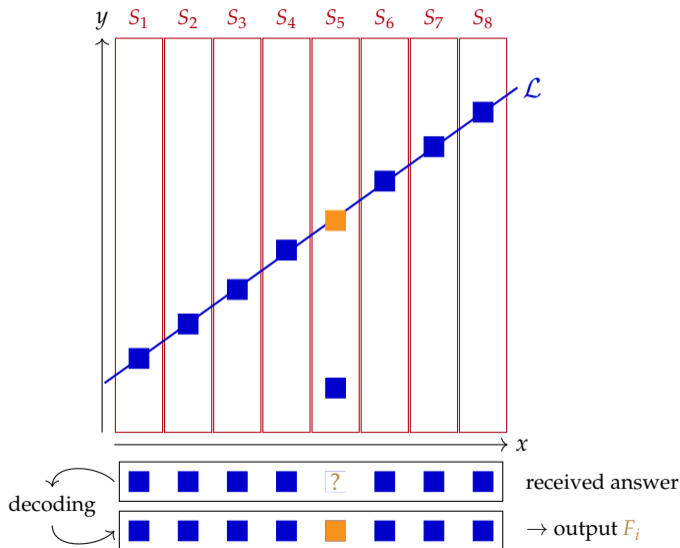
Assume one wants to extract file F_i .

Pick a line L through i , query all files F_j except F_i , and interpolate the corresponding univariate polynomial.

\mathcal{L} is an affine line

\implies no individual server can find where is ■

A PIR scheme based on Reed–Muller codes



Database F is encoded with $\text{RM}_q(m, r)$, then distributed across the servers

Assume one wants to extract file F_i .

Pick a line L through i , query all files F_j except F_i , and interpolate the corresponding univariate polynomial.

\mathcal{L} is an affine line

⇒ no individual server can find where is ■

⇒ the PIR scheme is private


Features with $\text{RM}_q(m, r)$ of length q^m .

- ▶ communication complexity: $(m - 1)q \log q$ uploaded bits, $q \log q$ downloaded bits
- ▶ computational complexity:
 - ▶ **only 1 read for each server (optimal)**
 - ▶ a decoding procedure for $\text{RS}_q(r)$ for the user
- ▶ storage overhead: the **information rate** of $\text{RM}_q(m, r)$ with $r \leq q - 1$ is

$$\simeq \frac{(r/q)^m}{m!} \leq \frac{1}{m!} \dots$$

\implies We need codes with the same “recovery properties”, and with larger dimension.

“Lifted” codes are the largest codes having the same property as Reed–Muller codes.

 *New affine-invariant codes from lifting*. Guo, Kopparty, Sudan. ITCS. 2013.

“Lifted” codes are the largest codes having the same property as Reed–Muller codes.

📄 *New affine-invariant codes from lifting.* Guo, Kopparty, Sudan. ITCS. 2013.

Definition. The *m*-th lifted Reed-Solomon code of degree *r* over \mathbb{F}_q is:

$$\text{Lift}_q(m, r) := \{\text{ev}_{\mathbb{A}^m}(f) \mid f \in \mathbb{F}_q[\mathbf{X}] \text{ and } \forall \text{ affine line } L \subset \mathbb{A}^m, \deg(f|_L) \leq r\}.$$

“Lifted” codes are the largest codes having the same property as Reed–Muller codes.


📄 *New affine-invariant codes from lifting*. Guo, Kopparty, Sudan. ITCS. 2013.

Definition. The *m*-th lifted Reed-Solomon code of degree *r* over \mathbb{F}_q is:

$$\text{Lift}_q(m, r) := \{\text{ev}_{\mathbb{A}^m}(f) \mid f \in \mathbb{F}_q[\mathbf{X}] \text{ and } \forall \text{ affine line } L \subset \mathbb{A}^m, \deg(f|_L) \leq r\}.$$

Lifted codes contain Reed-Muller codes, **sometimes properly**.

“Lifted” codes are the largest codes having the same property as Reed–Muller codes.

 *New affine-invariant codes from lifting.* Guo, Kopparty, Sudan. ITCS. 2013.

Definition. The *m*-th lifted Reed–Solomon code of degree *r* over \mathbb{F}_q is:

$$\text{Lift}_q(m, r) := \{\text{ev}_{\mathbb{A}^m}(f) \mid f \in \mathbb{F}_q[\mathbf{X}] \text{ and } \forall \text{ affine line } L \subset \mathbb{A}^m, \deg(f|_L) \leq r\}.$$

Lifted codes contain Reed–Muller codes, **sometimes properly**.

Example. For $q = 4$, $m = 2$, $r = 2$, consider $f(X, Y) = X^2Y^2$ and an affine line L with equation $(aT + b, cT + d)$. Then,

$$f(aT + b, cT + d) \equiv (a^2d^2 + b^2c^2)T^2 + a^2c^2T + b^2d^2 \pmod{(T^4 - T)}$$

corresponds to a degree–2 polynomial in T .

Hence,

$$\text{ev}(X^2Y^2) \in \text{Lift}_4(2, 2) \quad \text{but} \quad \text{ev}(X^2Y^2) \notin \text{RM}_4(2, 2).$$

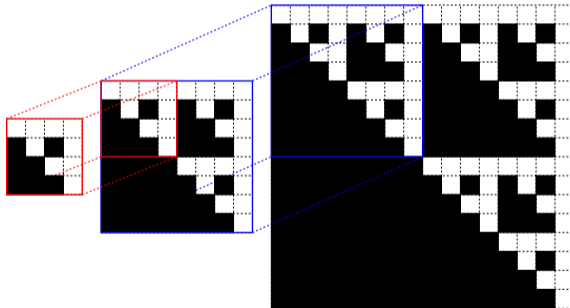
Theorem (Guo, Kopparty, Sudan '13). For every fixed m , and for growing alphabet and length, lifted codes reach arbitrarily large information rates.

Black squares: pairs (i, j) such that $\text{ev}(X^i Y^j) \in \text{Lift}_q(m = 2, r = q - 2)$.

$q = 4$

$q = 8$

$q = 16$



Corollary. For a sufficiently large number of servers, we have PIR with storage overhead $\rightarrow 0$.

Question: how to deal with collusions and byzantine errors?

Question: how to deal with collusions and byzantine errors?

For convenience, here $m = 2$.

Definition. A t -curve is:

$$\mathcal{L} = \{(x, g(x)) \in \mathbb{A}^2 \mid g \in \mathbb{F}_q[X], \deg(g) \leq t\}$$

Question: how to deal with collusions and byzantine errors?

For convenience, here $m = 2$.

Definition. A t -curve is:

$$\mathcal{L} = \{(x, g(x)) \in \mathbb{A}^2 \mid g \in \mathbb{F}_q[X], \deg(g) \leq t\}$$

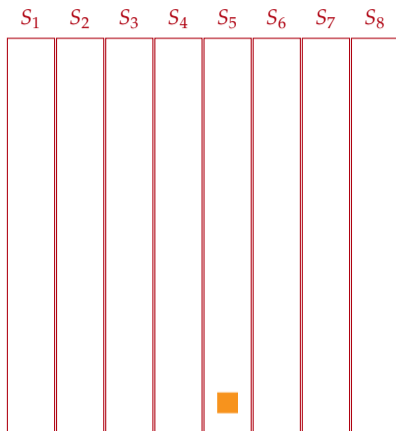
Definition. The **weighted lifted Reed-Solomon code** of degree r and weight t over \mathbb{F}_q is:

$$\text{WLift}_q(t, r) := \{\text{ev}_{\mathbb{A}^2}(f) \mid f \in \mathbb{F}_q[X, Y] \text{ and } \forall t\text{-curve } \mathcal{L} \subset \mathbb{A}^2, \deg(f|_{\mathcal{L}}) \leq r\}$$

Consequence: for every codeword $c \in \text{WLift}_q(t, r)$ and every t -curve \mathcal{L} , we have:

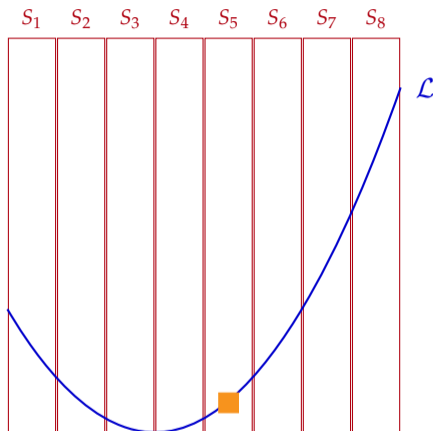
$$c|_{\mathcal{L}} \in \text{RS}_q(r + 1).$$

A PIR scheme based on weighted lifted codes



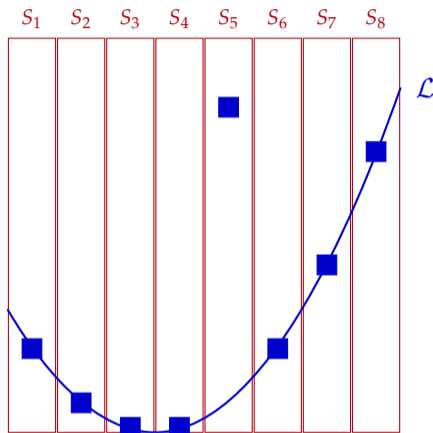
Database F is encoded with $WLift_q(t, r)$, then distributed across q servers

A PIR scheme based on weighted lifted codes



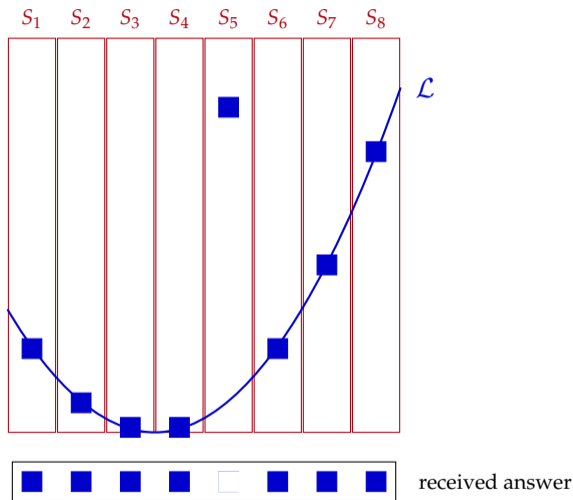
Database F is encoded with $\text{WLift}_q(t, r)$, then distributed across q servers

A PIR scheme based on weighted lifted codes



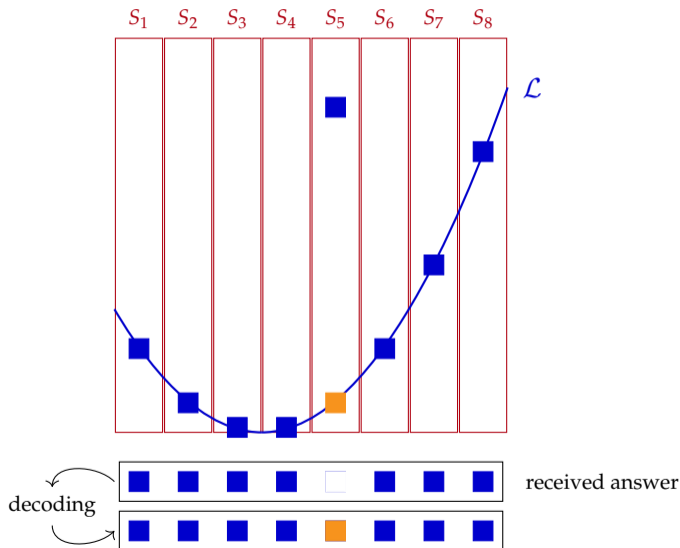
Database F is encoded with $WLift_q(t, r)$, then distributed across q servers

A PIR scheme based on weighted lifted codes



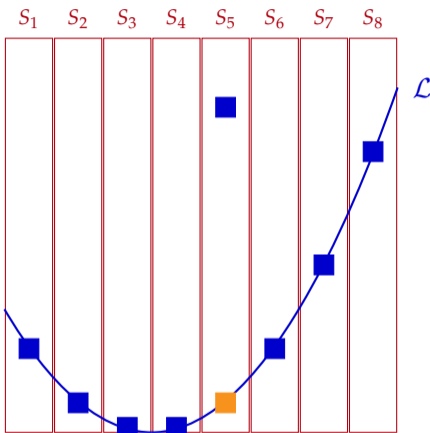
Database F is encoded with $WLift_q(t, r)$, then distributed across q servers

A PIR scheme based on weighted lifted codes



Database F is encoded with $WLift_q(t, r)$, then distributed across q servers

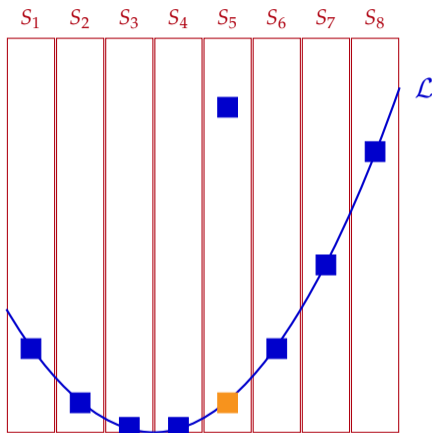
A PIR scheme based on weighted lifted codes



Database F is encoded with $WLift_q(t, r)$, then distributed across q servers



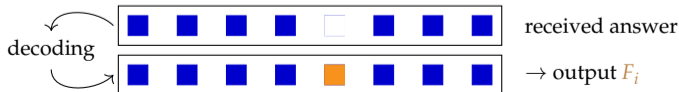
A PIR scheme based on weighted lifted codes



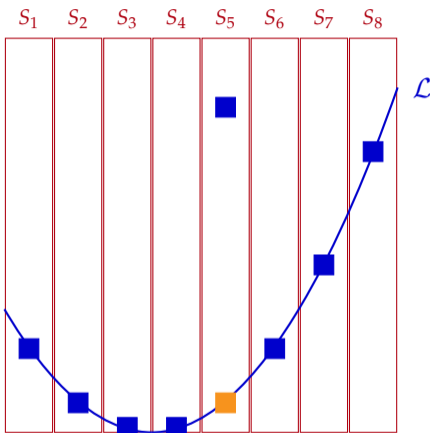
Database F is encoded with $WLift_q(t, r)$, then distributed across q servers

\mathcal{L} is a t -curve

\implies no t -set of servers can find where is ■



A PIR scheme based on weighted lifted codes



Database F is encoded with $WLift_q(t, r)$, then distributed across q servers

\mathcal{L} is a t -curve

\implies no t -set of servers can find where is \blacksquare


\implies the PIR scheme is t -private

decoding




received answer

\rightarrow output F_i

 *Weighted Lifted Codes: Local Correctabilities and Application to Robust Private Information Retrieval.* L., Nardi. IEEE TIT. 2021.

Theorem. Let p be a prime number, $t \geq 1$ and $\alpha \geq 2$ be fixed integers. Then, the information rate $\text{WLift}_{p^e}(t, p^e - \alpha)$ grows to 1 when $e \rightarrow \infty$.

Corollary: we get PIR schemes with relative storage overhead $\rightarrow 0$, for a constant number of adversaries.

 *Weighted Lifted Codes: Local Correctabilities and Application to Robust Private Information Retrieval.* L., Nardi. IEEE TIT. 2021.

Theorem. Let p be a prime number, $t \geq 1$ and $\alpha \geq 2$ be fixed integers. Then, the information rate $\text{WLift}_{p^e}(t, p^e - \alpha)$ grows to 1 when $e \rightarrow \infty$.

Corollary: we get PIR schemes with relative storage overhead $\rightarrow 0$, for a constant number of adversaries.

Theorem. Let p be a prime number, $t \geq 1$ and $c \geq 1$ be fixed integers. Let $\gamma = 1 - p^{-c}$ and $\mathcal{C}_e = \text{WLift}_{p^e}(t, \gamma p^e)$. Then, the information rate R_e of \mathcal{C}_e satisfies:

$$\lim_{e \rightarrow \infty} R_e = K_{t,p,c} > 0$$

Corollary: we get PIR schemes with **constant relative storage overhead**, for a **constant number of collusions** and a **constant fraction of errors**.

Other works. PIR has been a hot topic during for few years.

- ▶ Notion of **PIR capacity**: achievable bounds on the download rate of PIR schemes.
 - requires lot of comutation of the server side
- ▶ **Optimal** constructions over given distributed storage systems:
 - data is already encoded by the storage system
 - we can avoid re-encoding and still do PIR

📄 [SJ17] *The Capacity of Private Information Retrieval*. Sun, Jafar. IEEE-TIT. **2017**.

📄 [TGR18] *PIR from MDS Coded Data in Distributed Storage Systems*. Tajeddine, Gnilke, El Rouayheb. IEEE-TIT. **2018**.

📄 [TGKFH18] *Robust PIR from Coded Systems with Byzantine and Colluding Servers*. Tajeddine, Gnilke, Karpuk, Freij-Hollanti, Hollanti. ISIT. **2018**.

📄 *Private Information Retrieval Schemes With Product-Matrix MBR Codes*. L., Tajeddine, Freij-Hollanti, Hollanti. IEEE IFS. **2021**.

Open questions / future works.

1. Combinatorial bounds on the parameters.
2. Updatable files?
3. Extension to peer-to-peer storage systems (codes on random graphs).

Open questions / future works.

1. Combinatorial bounds on the parameters.
2. Updatable files?
3. Extension to peer-to-peer storage systems (codes on random graphs).

Thank you for your attention!