

# Des codes correcteurs pour le retrait confidentiel d'information

Julien Lavauzelle

équipe AGC3, LAGA, Université Paris 8 Vincennes – Saint-Denis

Séminaire d'ouverture aux informatiques et de formation

31/05/2022

## 1. Retrait confidentiel d'information

## 2. Codes correcteurs

## 3. Des schémas de PIR efficaces en temps et en mémoire

Avec une base de données répliquée

Avec une base de données distribuée

Résistance aux coalitions

## 1. Retrait confidentiel d'information

## 2. Codes correcteurs

## 3. Des schémas de PIR efficaces en temps et en mémoire

Avec une base de données répliquée

Avec une base de données distribuée

Résistance aux coalitions

Accroissement de l'accès à des données stockées à distance :

- ▶ données publiques (cartographie, informations & évènements)
- ▶ données personnelles (services de stockage en ligne, données médicales)

Accroissement de l'accès à des données stockées à distance :

- ▶ données publiques (cartographie, informations & évènements)
- ▶ données personnelles (services de stockage en ligne, données médicales)

Désir d'un accès **confidentiel** à ces données, par

- ▶ l'**anonymat** : l'identité de l'entité qui accède à l'information est inconnue

Accroissement de l'accès à des données stockées à distance :

- ▶ données publiques (cartographie, informations & événements)
- ▶ données personnelles (services de stockage en ligne, données médicales)

Désir d'un accès **confidentiel** à ces données, par

- ▶ l'**anonymat** : l'identité de l'entité qui accède à l'information est inconnue
- ▶ le **chiffrement** : la valeur de l'information accédée est inconnue

Accroissement de l'accès à des données stockées à distance :

- ▶ données publiques (cartographie, informations & évènements)
- ▶ données personnelles (services de stockage en ligne, données médicales)

Désir d'un accès **confidentiel** à ces données, par

- ▶ l'**anonymat** : l'identité de l'entité qui accède à l'information est inconnue
- ▶ le **chiffrement** : la valeur de l'information accédée est inconnue
- ▶ le **retrait confidentiel** : l'identité de l'information accédée est inconnue

Retrait confidentiel d'information (*private information retrieval*, PIR) :

Soit une base de données distante  $F \in \Sigma^M$  et un indice  $i \in [1, M]$ ,  
peut-on **extraire** l'entrée  $F_i$  de la base,  
**sans révéler aucune information** sur l'indice  $i$  ?



Retrait confidentiel d'information (*private information retrieval*, PIR) :

Soit une base de données distante  $F \in \Sigma^M$  et un indice  $i \in [1, M]$ ,  
peut-on **extraire** l'entrée  $F_i$  de la base,  
**sans révéler aucune information** sur l'indice  $i$  ?

Une **solution triviale** : télécharger complètement la base de données.

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

Un vecteur d'entrées  $F \in \Sigma^M$  est stocké sur  $n$  serveurs  $S_1, \dots, S_n$ .  
L'entité User désire extraire l'entrée  $F_i$  de manière confidentielle.

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

Un vecteur d'entrées  $F \in \Sigma^M$  est stocké sur  $n$  serveurs  $S_1, \dots, S_n$ .

L'entité User désire extraire l'entrée  $F_i$  de manière confidentielle.

Un protocole de **retrait confidentiel d'information** est un ensemble de trois algorithmes (Query, Answer, Reconstruct). Pour extraire  $F_i$  :

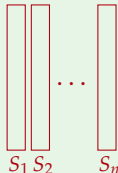
📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

Un vecteur d'entrées  $F \in \Sigma^M$  est stocké sur  $n$  serveurs  $S_1, \dots, S_n$ .

L'entité User désire extraire l'entrée  $F_i$  de manière confidentielle.

Un protocole de **retrait confidentiel d'information** est un ensemble de trois algorithmes (Query, Answer, Reconstruct). Pour extraire  $F_i$  :

User



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

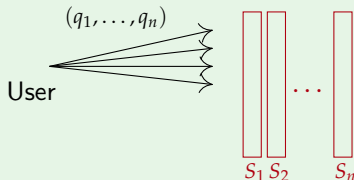
Un vecteur d'entrées  $F \in \Sigma^M$  est stocké sur  $n$  serveurs  $S_1, \dots, S_n$ .  
L'entité User désire extraire l'entrée  $F_i$  de manière confidentielle.

Un protocole de **retrait confidentiel d'information** est un ensemble de trois algorithmes (Query, Answer, Reconstruct). Pour extraire  $F_i$  :

1. User engendre un vecteur de requêtes

$$q := (q_1, \dots, q_n) \leftarrow \text{Query}(i)$$

et envoie  $q_j$  au serveur  $S_j$ .



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

Un vecteur d'entrées  $F \in \Sigma^M$  est stocké sur  $n$  serveurs  $S_1, \dots, S_n$ .  
L'entité User désire extraire l'entrée  $F_i$  de manière confidentielle.

Un protocole de **retrait confidentiel d'information** est un ensemble de trois algorithmes (Query, Answer, Reconstruct). Pour extraire  $F_i$  :

1. User engendre un vecteur de requêtes

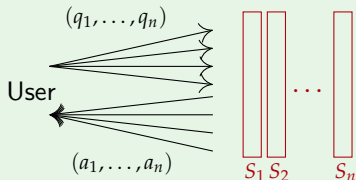
$$q := (q_1, \dots, q_n) \leftarrow \text{Query}(i)$$

et envoie  $q_j$  au serveur  $S_j$ .

2. Chaque serveur  $S_j$  calcule

$$a_j := \text{Answer}(q_j, F|_{S_j})$$

et retourne la valeur  $a_j$  à User.



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

Un vecteur d'entrées  $F \in \Sigma^M$  est stocké sur  $n$  serveurs  $S_1, \dots, S_n$ .  
L'entité User désire extraire l'entrée  $F_i$  de manière confidentielle.

Un protocole de **retrait confidentiel d'information** est un ensemble de trois algorithmes (Query, Answer, Reconstruct). Pour extraire  $F_i$  :

1. User engendre un vecteur de requêtes

$$q := (q_1, \dots, q_n) \leftarrow \text{Query}(i)$$

et envoie  $q_j$  au serveur  $S_j$ .

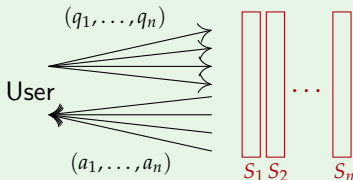
2. Chaque serveur  $S_j$  calcule

$$a_j := \text{Answer}(q_j, F|_{S_j})$$

et retourne la valeur  $a_j$  à User.

3. User reconstitue l'information désirée par :

$$F_i = \text{Reconstruct}(q, a, i).$$



L'**adversaire** : des **coalitions de serveurs** = des sous-ensembles de  $\{S_j : j \in T\}$ , avec  $T \subset [1, n]$ , qui échangent de l'information sur les requêtes reçues.



L'**adversaire** : des **coalitions de serveurs** = des sous-ensembles de  $\{S_j : j \in T\}$ , avec  $T \subset [1, n]$ , qui échangent de l'information sur les requêtes reçues.

On note  $t$  la taille de la plus grande coalition.

L'**adversaire** : des **coalitions de serveurs** = des sous-ensembles de  $\{S_j : j \in T\}$ , avec  $T \subset [1, n]$ , qui échangent de l'information sur les requêtes reçues.

On note  $t$  la taille de la plus grande coalition.

Deux modèles de sécurité :

L'**adversaire** : des **coalitions de serveurs** = des sous-ensembles de  $\{S_j : j \in T\}$ , avec  $T \subset [1, n]$ , qui échangent de l'information sur les requêtes reçues.

On note  $t$  la taille de la plus grande coalition.

Deux modèles de sécurité :

- **Au sens de la théorie de l'information (IT-PIR) :**

Si  $q = \text{Query}(i)$ , on souhaite que :

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

L'**adversaire** : des **coalitions de serveurs** = des sous-ensembles de  $\{S_j : j \in T\}$ , avec  $T \subset [1, n]$ , qui échangent de l'information sur les requêtes reçues.

On note  $t$  la taille de la plus grande coalition.

Deux modèles de sécurité :

- **Au sens de la théorie de l'information (IT-PIR) :**

Si  $q = \text{Query}(i)$ , on souhaite que :

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

- **Au sens calculatoire (C-PIR) :**

À  $T$  fixé et  $i$  variable, les distributions de requêtes  $q_{|T} = \text{Query}(i)_{|T}$  sont calculatoirement indistinguables.

L'**adversaire** : des **coalitions de serveurs** = des sous-ensembles de  $\{S_j : j \in T\}$ , avec  $T \subset [1, n]$ , qui échangent de l'information sur les requêtes reçues.

On note  $t$  la taille de la plus grande coalition.

Deux modèles de sécurité :

- **Au sens de la théorie de l'information (IT-PIR) :**

Si  $q = \text{Query}(i)$ , on souhaite que :

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

- **Au sens calculatoire (C-PIR) :**

À  $T$  fixé et  $i$  variable, les distributions de requêtes  $q_{|T} = \text{Query}(i)_{|T}$  sont calculatoirement indistinguables.

**Théorème [CGKS95, CG97].** Si  $t = n$  (en particulier pour  $n = 1$  serveur), alors IT-PIR implique de télécharger entièrement la base de données.

L'**adversaire** : des **coalitions de serveurs** = des sous-ensembles de  $\{S_j : j \in T\}$ , avec  $T \subset [1, n]$ , qui échangent de l'information sur les requêtes reçues.

On note  $t$  la taille de la plus grande coalition.

Deux modèles de sécurité :

- **Au sens de la théorie de l'information (IT-PIR) :**

Si  $q = \text{Query}(i)$ , on souhaite que :

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

- **Au sens calculatoire (C-PIR) :**

À  $T$  fixé et  $i$  variable, les distributions de requêtes  $q_{|T} = \text{Query}(i)_{|T}$  sont calculatoirement indistinguables.

**Théorème [CGKS95, CG97].** Si  $t = n$  (en particulier pour  $n = 1$  serveur), alors IT-PIR implique de télécharger entièrement la base de données.

**Remarque.** C-PIR est possible avec 1 serveur, mais reste trop coûteux en pratique.

On s'intéresse à **IT-PIR** (donc avec  $n \geq 2$  serveurs)

On s'intéresse à **IT-PIR** (donc avec  $n \geq 2$  serveurs)

**Paramètres** à prendre en compte :

- complexité de **communication** (*upload* et *download*)



On s'intéresse à **IT-PIR** (donc avec  $n \geq 2$  serveurs)

**Paramètres** à prendre en compte :

- complexité de **communication** (*upload* et *download*)
- complexité de **calcul** (client et serveur)

On s'intéresse à **IT-PIR** (donc avec  $n \geq 2$  serveurs)

**Paramètres** à prendre en compte :

- complexité de **communication** (*upload* et *download*)
- complexité de **calcul** (client et serveur)
- coût de **stockage** supplémentaire

On s'intéresse à **IT-PIR** (donc avec  $n \geq 2$  serveurs)

**Paramètres** à prendre en compte :

- complexité de **communication** (*upload* et *download*)
- complexité de **calcul** (client et serveur)
- coût de **stockage** supplémentaire
- taille maximale des coalitions d'adversaires

On s'intéresse à **IT-PIR** (donc avec  $n \geq 2$  serveurs)

**Paramètres** à prendre en compte :

- complexité de **communication** (*upload* et *download*)
- complexité de **calcul** (client et serveur)
- coût de **stockage** supplémentaire
- taille maximale des coalitions d'adversaires

Pour différents **contextes**, suivant :

- si la base de données est **répliquée** ou **ré-encodée**

On s'intéresse à **IT-PIR** (donc avec  $n \geq 2$  serveurs)

**Paramètres** à prendre en compte :

- complexité de **communication** (*upload* et *download*)
- complexité de **calcul** (client et serveur)
- coût de **stockage** supplémentaire
- taille maximale des coalitions d'adversaires

Pour différents **contextes**, suivant :

- si la base de données est **répliquée** ou **ré-encodée**
- la présence de serveurs **sans réponse, malicieux**

On s'intéresse à **IT-PIR** (donc avec  $n \geq 2$  serveurs)

**Paramètres** à prendre en compte :

- complexité de **communication** (*upload* et *download*)
- complexité de **calcul** (client et serveur)
- coût de **stockage** supplémentaire
- taille maximale des coalitions d'adversaires

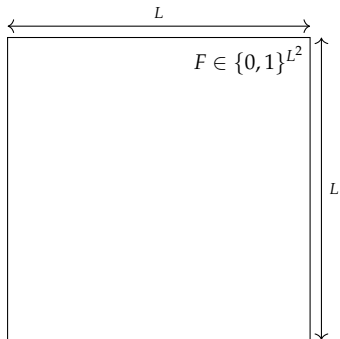
Pour différents **contextes**, suivant :

- si la base de données est **répliquée** ou **ré-encodée**
- la présence de serveurs **sans réponse, malicieux**
- la taille des entrées (clés vs. fichiers vidéos)
- le nombre d'entrées
- la statiticité (mise à jour, ajout d'entrées, etc.)

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

**Contexte :** base de données  $F$  répliquée sur  $n$  serveurs, avec :

- ▶  $|F| = M$  entrées de même taille (e.g. 1 bit), avec  $M = L^2$  et  $[1, M] \simeq [1, L]^2$ .
- ▶  $n = 4$  serveurs  $S_{00}, S_{01}, S_{10}, S_{11}$ , qui stockent chacun une copie de  $F$ .

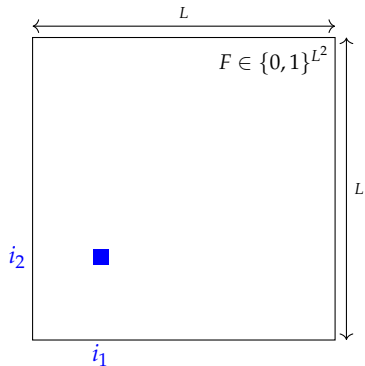


📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

**Contexte :** base de données  $F$  répliquée sur  $n$  serveurs, avec :

- ▶  $|F| = M$  entrées de même taille (e.g. 1 bit), avec  $M = L^2$  et  $[1, M] \simeq [1, L]^2$ .
- ▶  $n = 4$  serveurs  $S_{00}, S_{01}, S_{10}, S_{11}$ , qui stockent chacun une copie de  $F$ .

**Objectif :** extraire  $F_i = F_{(i_1, i_2)}$  pour  $1 \leq i_1, i_2 \leq L$ .



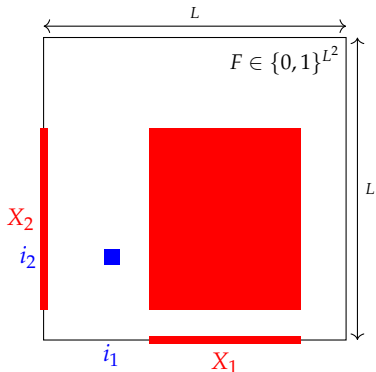


📄 *Private Information Retrieval*. Chor, Goldreich, Kushilevitz, Sudan. FOCS. 1995.

**Contexte :** base de données  $F$  répliquée sur  $n$  serveurs, avec :

- ▶  $|F| = M$  entrées de même taille (e.g. 1 bit), avec  $M = L^2$  et  $[1, M] \simeq [1, L]^2$ .
- ▶  $n = 4$  serveurs  $S_{00}, S_{01}, S_{10}, S_{11}$ , qui stockent chacun une copie de  $F$ .

**Objectif :** extraire  $F_{\mathbf{i}} = F_{(i_1, i_2)}$  pour  $1 \leq i_1, i_2 \leq L$ .



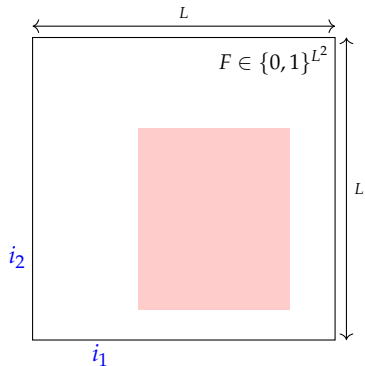
1. User engendre aléatoirement deux sous-ensembles  $X_1, X_2$  de  $[1, L]$ . Ensuite, User envoie :

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

**Contexte :** base de données  $F$  répliquée sur  $n$  serveurs, avec :

- ▶  $|F| = M$  entrées de même taille (e.g. 1 bit), avec  $M = L^2$  et  $[1, M] \simeq [1, L]^2$ .
- ▶  $n = 4$  serveurs  $S_{00}, S_{01}, S_{10}, S_{11}$ , qui stockent chacun une copie de  $F$ .

**Objectif :** extraire  $F_{\mathbf{i}} = F_{(i_1, i_2)}$  pour  $1 \leq i_1, i_2 \leq L$ .



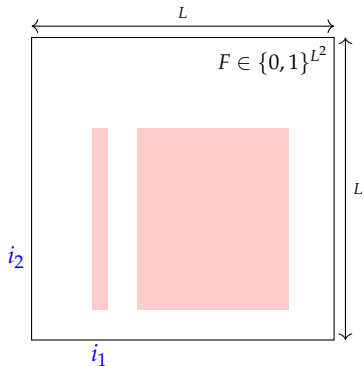
1. User engendre aléatoirement deux sous-ensembles  $X_1, X_2$  de  $[1, L]$ . Ensuite, User envoie :
  - Query( $\mathbf{i}$ )<sub>00</sub> = (  $X_1$  ,  $X_2$  ) à  $S_{00}$ ,

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

**Contexte :** base de données  $F$  répliquée sur  $n$  serveurs, avec :

- ▶  $|F| = M$  entrées de même taille (e.g. 1 bit), avec  $M = L^2$  et  $[1, M] \simeq [1, L]^2$ .
- ▶  $n = 4$  serveurs  $S_{00}, S_{01}, S_{10}, S_{11}$ , qui stockent chacun une copie de  $F$ .

**Objectif :** extraire  $F_{\mathbf{i}} = F_{(i_1, i_2)}$  pour  $1 \leq i_1, i_2 \leq L$ .



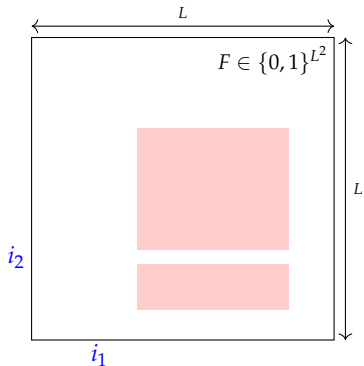
1. User engendre aléatoirement deux sous-ensembles  $X_1, X_2$  de  $[1, L]$ . Ensuite, User envoie :
  - Query( $\mathbf{i}$ )<sub>00</sub> = (  $X_1$  ,  $X_2$  ) à  $S_{00}$ ,
  - Query( $\mathbf{i}$ )<sub>10</sub> = (  $X_1 \Delta \{i_1\}$ ,  $X_2$  ) à  $S_{10}$ ,

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

**Contexte :** base de données  $F$  répliquée sur  $n$  serveurs, avec :

- ▶  $|F| = M$  entrées de même taille (e.g. 1 bit), avec  $M = L^2$  et  $[1, M] \simeq [1, L]^2$ .
- ▶  $n = 4$  serveurs  $S_{00}, S_{01}, S_{10}, S_{11}$ , qui stockent chacun une copie de  $F$ .

**Objectif :** extraire  $F_{\mathbf{i}} = F_{(i_1, i_2)}$  pour  $1 \leq i_1, i_2 \leq L$ .



1. User engendre aléatoirement deux sous-ensembles  $X_1, X_2$  de  $[1, L]$ . Ensuite, User envoie :

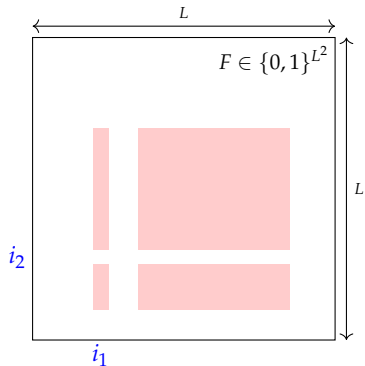
- Query( $\mathbf{i}$ )<sub>00</sub> = (  $X_1$  ,  $X_2$  ) à  $S_{00}$ ,
- Query( $\mathbf{i}$ )<sub>10</sub> = (  $X_1 \Delta \{i_1\}$ ,  $X_2$  ) à  $S_{10}$ ,
- Query( $\mathbf{i}$ )<sub>01</sub> = (  $X_1$  ,  $X_2 \Delta \{i_2\}$  ) à  $S_{01}$ ,

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

**Contexte :** base de données  $F$  répliquée sur  $n$  serveurs, avec :

- ▶  $|F| = M$  entrées de même taille (e.g. 1 bit), avec  $M = L^2$  et  $[1, M] \simeq [1, L]^2$ .
- ▶  $n = 4$  serveurs  $S_{00}, S_{01}, S_{10}, S_{11}$ , qui stockent chacun une copie de  $F$ .

**Objectif :** extraire  $F_{\mathbf{i}} = F_{(i_1, i_2)}$  pour  $1 \leq i_1, i_2 \leq L$ .



1. User engendre aléatoirement deux sous-ensembles  $X_1, X_2$  de  $[1, L]$ . Ensuite, User envoie :

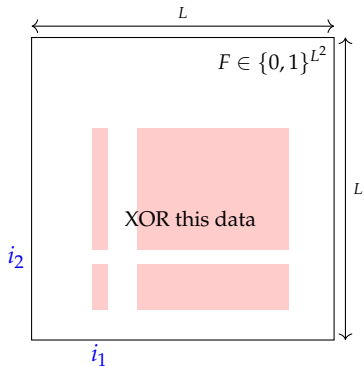
- Query( $\mathbf{i}$ )<sub>00</sub> = (  $X_1$  ,  $X_2$  ) à  $S_{00}$ ,
- Query( $\mathbf{i}$ )<sub>10</sub> = (  $X_1 \Delta \{i_1\}$ ,  $X_2$  ) à  $S_{10}$ ,
- Query( $\mathbf{i}$ )<sub>01</sub> = (  $X_1$  ,  $X_2 \Delta \{i_2\}$  ) à  $S_{01}$ ,
- Query( $\mathbf{i}$ )<sub>11</sub> = (  $X_1 \Delta \{i_1\}$ ,  $X_2 \Delta \{i_2\}$  ) à  $S_{11}$ .

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

**Contexte :** base de données  $F$  répliquée sur  $n$  serveurs, avec :

- ▶  $|F| = M$  entrées de même taille (e.g. 1 bit), avec  $M = L^2$  et  $[1, M] \simeq [1, L]^2$ .
- ▶  $n = 4$  serveurs  $S_{00}, S_{01}, S_{10}, S_{11}$ , qui stockent chacun une copie de  $F$ .

**Objectif :** extraire  $F_{\mathbf{i}} = F_{(i_1, i_2)}$  pour  $1 \leq i_1, i_2 \leq L$ .



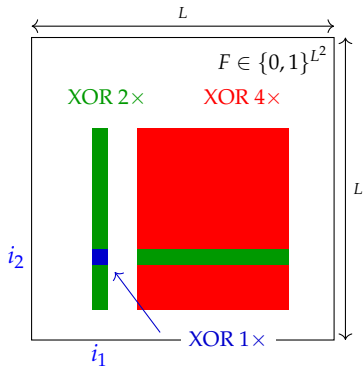
1. User engendre aléatoirement deux sous-ensembles  $X_1, X_2$  de  $[1, L]$ . Ensuite, User envoie :
  - Query( $\mathbf{i}$ )<sub>00</sub> = (  $X_1$  ,  $X_2$  ) à  $S_{00}$ ,
  - Query( $\mathbf{i}$ )<sub>10</sub> = (  $X_1 \Delta \{i_1\}$ ,  $X_2$  ) à  $S_{10}$ ,
  - Query( $\mathbf{i}$ )<sub>01</sub> = (  $X_1$  ,  $X_2 \Delta \{i_2\}$  ) à  $S_{01}$ ,
  - Query( $\mathbf{i}$ )<sub>11</sub> = (  $X_1 \Delta \{i_1\}$ ,  $X_2 \Delta \{i_2\}$  ) à  $S_{11}$ .
2. Sur réception de  $(Z_1, Z_2)$ , chaque serveur calcule et retourne  $a = \bigoplus_{z \in Z_1 \times Z_2} F_z$ .

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

**Contexte :** base de données  $F$  répliquée sur  $n$  serveurs, avec :

- ▶  $|F| = M$  entrées de même taille (e.g. 1 bit), avec  $M = L^2$  et  $[1, M] \simeq [1, L]^2$ .
- ▶  $n = 4$  serveurs  $S_{00}, S_{01}, S_{10}, S_{11}$ , qui stockent chacun une copie de  $F$ .

**Objectif :** extraire  $F_i = F_{(i_1, i_2)}$  pour  $1 \leq i_1, i_2 \leq L$ .



1. User engendre aléatoirement deux sous-ensembles  $X_1, X_2$  de  $[1, L]$ . Ensuite, User envoie :
  - Query( $i$ )<sub>00</sub> = ( $X_1$ ,  $X_2$ ) à  $S_{00}$ ,
  - Query( $i$ )<sub>10</sub> = ( $X_1 \Delta \{i_1\}$ ,  $X_2$ ) à  $S_{10}$ ,
  - Query( $i$ )<sub>01</sub> = ( $X_1$ ,  $X_2 \Delta \{i_2\}$ ) à  $S_{01}$ ,
  - Query( $i$ )<sub>11</sub> = ( $X_1 \Delta \{i_1\}$ ,  $X_2 \Delta \{i_2\}$ ) à  $S_{11}$ .
2. Sur réception de  $(Z_1, Z_2)$ , chaque serveur calcule et retourne  $a = \bigoplus_{z \in Z_1 \times Z_2} F_z$ .
3. User XOR les 4 bits reçus et obtient  $F_i$ .

Le protocole est **correct** et **confidentiel** si  $t = 1$  (pas de coalition).



Le protocole est **correct** et **confidentiel** si  $t = 1$  (pas de coalition).

- ▶ **Communication** :  $8\sqrt{M}$  bits envoyés, 4 bits reçus,
- ▶ **Stockage** : copie de  $F$  sur  $n = 4$  serveurs,
- ▶ **Calcul** :
  - ▶ pour chaque serveur : en moyenne, XOR de  $(L/2)^2 = M/4$  bits
  - ▶ pour le client : XOR de  $n = 4$  bits.

Le protocole est **correct** et **confidentiel** si  $t = 1$  (pas de coalition).

- ▶ **Communication** :  $8\sqrt{M}$  bits envoyés, 4 bits reçus,
- ▶ **Stockage** : copie de  $F$  sur  $n = 4$  serveurs,
- ▶ **Calcul** :
  - ▶ pour chaque serveur : en moyenne, XOR de  $(L/2)^2 = M/4$  bits
  - ▶ pour le client : XOR de  $n = 4$  bits.

Généralisable à  $n = 2^b$  serveurs :

- ▶ **Communication** :  $n \log(n) M^{1/\log(n)}$  bits envoyés,  $n$  bits reçus,
- ▶ **Stockage** : copie de  $F$  sur  $n$  serveurs,
- ▶ **Calcul** :
  - ▶ pour chaque serveur : en moyenne, XOR de  $M/n$  bits
  - ▶ pour l'utilisateur : XOR de  $n$  bits.

### Dans le protocole précédent :

- ▶ complexité de communication sous-linéaire, mais “moyenne”
- ▶ calcul lourd pour les serveurs (linéaire en la taille de  $F$ )
- ▶ nécessite de copier la base de données sur plusieurs serveurs

### Dans le protocole précédent :

- ▶ complexité de communication sous-linéaire, mais “moyenne”
- ▶ calcul lourd pour les serveurs (linéaire en la taille de  $F$ )
- ▶ nécessite de copier la base de données sur plusieurs serveurs

### Objectifs :

- ▶ meilleure complexité de communication
- ▶ **complexité de calcul optimale**
- ▶ plus faible redondance de stockage via un **encodage** et une distribution de la base de données

### Dans le protocole précédent :

- ▶ complexité de communication sous-linéaire, mais “moyenne”
- ▶ calcul lourd pour les serveurs (linéaire en la taille de  $F$ )
- ▶ nécessite de copier la base de données sur plusieurs serveurs

### Objectifs :

- ▶ meilleure complexité de communication
- ▶ **complexité de calcul optimale**
- ▶ plus faible redondance de stockage via un **encodage** et une distribution de la base de données

### Outils : théorie des codes et combinatoire

1. Retrait confidentiel d'information

2. Codes correcteurs

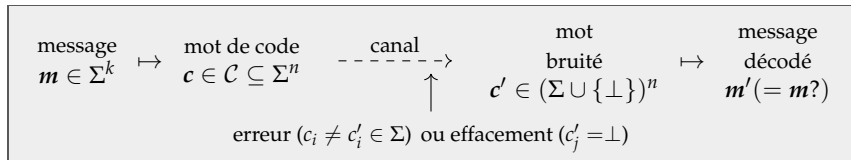
3. Des schémas de PIR efficaces en temps et en mémoire

Avec une base de données répliquée

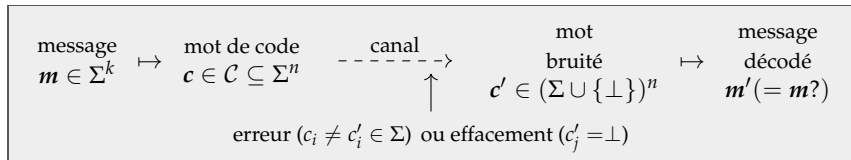
Avec une base de données distribuée

Résistance aux coalitions

Transmission d'information en présence de bruit.



Transmission d'information en présence de bruit.

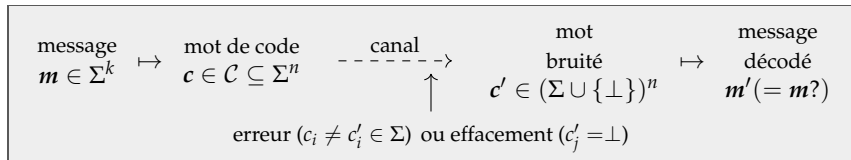


## Remarques :

- ici tous les mots de code  $c$  ont même longueur  $n$
- l'application d'encodage  $\phi : m \mapsto c$  doit être injective



Transmission d'information en présence de bruit.

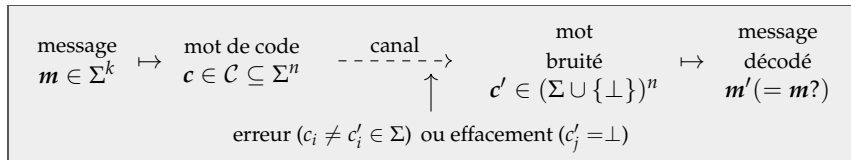


## Remarques :

- ici tous les mots de code  $c$  ont même longueur  $n$
- l'application d'encodage  $\phi : m \mapsto c$  doit être injective

L'ensemble  $\mathcal{C}$  des mots de code (pour  $m \in \Sigma^k$ ) est appelé **code correcteur**.

Transmission d'information en présence de bruit.



## Remarques :

- ici tous les mots de code  $c$  ont même longueur  $n$
- l'application d'encodage  $\phi : m \mapsto c$  doit être injective

L'ensemble  $\mathcal{C}$  des mots de code (pour  $m \in \Sigma^k$ ) est appelé **code correcteur**.

Pour des aspects **pratiques** :

- $\Sigma = \mathbb{F}$  est un corps (e.g.  $\{0, 1\}$ ,  $\mathbb{Z}/p\mathbb{Z}$ ,  $\mathbb{F}_q$ ),
- $\phi$  est  $\mathbb{F}$ -linéaire  $\implies$  le code  $\mathcal{C}$  est un espace vectoriel, on note  $k = \dim_{\mathbb{F}} \mathcal{C}$ .

Soit  $\mathcal{C} \subseteq \mathbb{F}^n$  un code linéaire de dimension  $k$ .

**Définition.** Toute matrice  $G \in \mathbb{F}^{k \times n}$  telle que

$$\mathcal{C} = \{mG \mid m \in \mathbb{F}^k\}$$

est une **matrice génératrice** de  $\mathcal{C}$ .

Soit  $\mathcal{C} \subseteq \mathbb{F}^n$  un code linéaire de dimension  $k$ .

**Définition.** Toute matrice  $G \in \mathbb{F}^{k \times n}$  telle que

$$\mathcal{C} = \{mG \mid m \in \mathbb{F}^k\}$$

est une **matrice génératrice** de  $\mathcal{C}$ .

On appelle **équation de parité** tout vecteur  $h \in \mathbb{F}^n$  tel que

$$\forall c \in \mathcal{C}, \langle h, c \rangle := \sum_{i=1}^n h_i c_i = 0$$

On note  $\mathcal{C}^\perp$  l'ensemble des équations de parité de  $\mathcal{C}$ .

Soit  $\mathcal{C} \subseteq \mathbb{F}^n$  un code linéaire de dimension  $k$ .

**Définition.** Toute matrice  $G \in \mathbb{F}^{k \times n}$  telle que

$$\mathcal{C} = \{mG \mid m \in \mathbb{F}^k\}$$

est une **matrice génératrice** de  $\mathcal{C}$ .

On appelle **équation de parité** tout vecteur  $h \in \mathbb{F}^n$  tel que

$$\forall c \in \mathcal{C}, \langle h, c \rangle := \sum_{i=1}^n h_i c_i = 0$$

On note  $\mathcal{C}^\perp$  l'ensemble des équations de parité de  $\mathcal{C}$ .

**Définition.** Toute matrice  $H \in \mathbb{F}^{m \times n}$  telle que

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

est une **matrice de parité** de  $\mathcal{C}$ .

Soit  $\mathcal{C} \subseteq \mathbb{F}^n$  un code linéaire de dimension  $k$ .

**Définition.** Toute matrice  $G \in \mathbb{F}^{k \times n}$  telle que

$$\mathcal{C} = \{mG \mid m \in \mathbb{F}^k\}$$

est une **matrice génératrice** de  $\mathcal{C}$ .

On appelle **équation de parité** tout vecteur  $h \in \mathbb{F}^n$  tel que

$$\forall c \in \mathcal{C}, \langle h, c \rangle := \sum_{i=1}^n h_i c_i = 0$$

On note  $\mathcal{C}^\perp$  l'ensemble des équations de parité de  $\mathcal{C}$ .

**Définition.** Toute matrice  $H \in \mathbb{F}^{m \times n}$  telle que

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

est une **matrice de parité** de  $\mathcal{C}$ .

On a alors  $\dim_{\mathbb{F}}(\mathcal{C}) = n - \text{rk}_{\mathbb{F}}(H)$ .

Un **effacement**, représenté par  $\perp$ , est une coordonnée d'un mot que l'on **sait** fausse.

Un **effacement**, représenté par  $\perp$ , est une coordonnée d'un mot que l'on **sait** fausse.

Soit  $c \in \mathcal{C}$ , dégradé par un effacement à l'indice  $i$  :  $c_i = \perp$ .

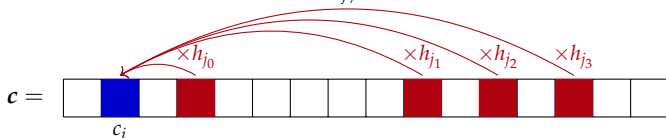


Un **effacement**, représenté par  $\perp$ , est une coordonnée d'un mot que l'on **sait** fausse.

Soit  $c \in \mathcal{C}$ , dégradé par un effacement à l'indice  $i$ :  $c_i = \perp$ .

S'il existe une équation de parité  $h \in \mathcal{C}^\perp$  telle que  $h_i \neq 0$ , alors :

$$c_i = \frac{-1}{h_i} \sum_{j \neq i} h_j c_j.$$

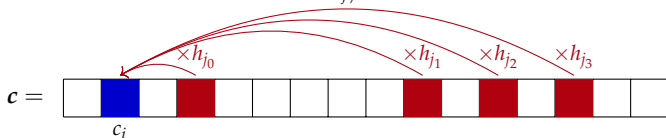


Un **effacement**, représenté par  $\perp$ , est une coordonnée d'un mot que l'on **sait** fausse.

Soit  $c \in \mathcal{C}$ , dégradé par un effacement à l'indice  $i$ :  $c_i = \perp$ .

S'il existe une équation de parité  $h \in \mathcal{C}^\perp$  telle que  $h_i \neq 0$ , alors :

$$c_i = \frac{-1}{h_i} \sum_{j \neq i} h_j c_j.$$



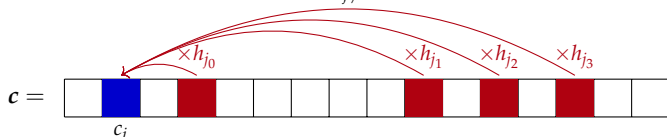
Notons  $\text{supp}(h) = \{j \in [1, n] \mid h_j \neq 0\}$  le **support** de  $h$ .

Un **effacement**, représenté par  $\perp$ , est une coordonnée d'un mot que l'on **sait** fausse.

Soit  $c \in \mathcal{C}$ , dégradé par un effacement à l'indice  $i$  :  $c_i = \perp$ .

S'il existe une équation de parité  $h \in \mathcal{C}^\perp$  telle que  $h_i \neq 0$ , alors :

$$c_i = \frac{-1}{h_i} \sum_{j \neq i} h_j c_j.$$



Notons  $\text{supp}(h) = \{j \in [1, n] \mid h_j \neq 0\}$  le **support** de  $h$ .

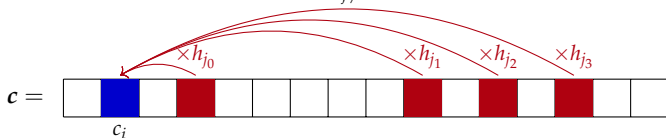
- on peut réparer la coordonnée  $c_i$  s'il n'y a pas d'erreur/effacement sur  $\text{supp}(h) \setminus \{i\}$  ;

Un **effacement**, représenté par  $\perp$ , est une coordonnée d'un mot que l'on **sait** fausse.

Soit  $c \in \mathcal{C}$ , dégradé par un effacement à l'indice  $i$ :  $c_i = \perp$ .

S'il existe une équation de parité  $h \in \mathcal{C}^\perp$  telle que  $h_i \neq 0$ , alors :

$$c_i = \frac{-1}{h_i} \sum_{j \neq i} h_j c_j.$$



Notons  $\text{supp}(h) = \{j \in [1, n] \mid h_j \neq 0\}$  le **support** de  $h$ .

- on peut réparer la coordonnée  $c_i$  s'il n'y a pas d'erreur/effacement sur  $\text{supp}(h) \setminus \{i\}$  ;
- pour une réparation efficace, on souhaite  $|\text{supp}(h)|$  le plus petit possible.

Une **erreur** est un symbole dégradé d'un mot dont on ne connaît pas la position.

Une **erreur** est un symbole dégradé d'un mot dont on ne connaît pas la position.

Certains codes sont munis d'**algorithmes de décodage** efficaces.

Les algorithmes de décodage classiques nécessitent, pour corriger une erreur, de **lire entièrement** le mot erroné  $\implies$  complexité  $\Omega(n)$ .

Une **erreur** est un symbole dégradé d'un mot dont on ne connaît pas la position.

Certains codes sont munis d'**algorithmes de décodage** efficaces.

Les algorithmes de décodage classiques nécessitent, pour corriger une erreur, de **lire entièrement** le mot erroné  $\implies$  complexité  $\Omega(n)$ .

**Question légitime** : décodage (probabiliste) d'erreur en complexité sous-linéaire ?

Une **erreur** est un symbole dégradé d'un mot dont on ne connaît pas la position.

Certains codes sont munis d'**algorithmes de décodage** efficaces.

Les algorithmes de décodage classiques nécessitent, pour corriger une erreur, de **lire entièrement** le mot erroné  $\implies$  complexité  $\Omega(n)$ .

**Question légitime** : décodage (probabiliste) d'erreur en complexité sous-linéaire ?

**Idée (informelle)** : construire un code avec beaucoup d'équations de parité de même petit poids  $\ell \ll n$ , **réparties uniformément** sur  $[1, n]$ . Puis, tirer au hasard l'une d'entre elles et l'utiliser pour réparer  $c_i$ .



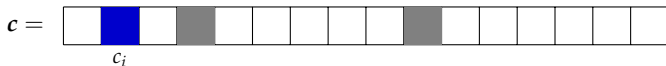
Une **erreur** est un symbole dégradé d'un mot dont on ne connaît pas la position.

Certains codes sont munis d'**algorithmes de décodage** efficaces.

Les algorithmes de décodage classiques nécessitent, pour corriger une erreur, de **lire entièrement** le mot erroné  $\implies$  complexité  $\Omega(n)$ .

**Question légitime** : décodage (probabiliste) d'erreur en complexité sous-linéaire ?

**Idée (informelle)** : construire un code avec beaucoup d'équations de parité de même petit poids  $\ell \ll n$ , **réparties uniformément** sur  $[1, n]$ . Puis, tirer au hasard l'une d'entre elles et l'utiliser pour réparer  $c_i$ .



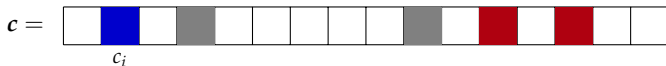
Une **erreur** est un symbole dégradé d'un mot dont on ne connaît pas la position.

Certains codes sont munis d'**algorithmes de décodage** efficaces.

Les algorithmes de décodage classiques nécessitent, pour corriger une erreur, de **lire entièrement** le mot erroné  $\implies$  complexité  $\Omega(n)$ .

**Question légitime** : décodage (probabiliste) d'erreur en complexité sous-linéaire ?

**Idée (informelle)** : construire un code avec beaucoup d'équations de parité de même petit poids  $\ell \ll n$ , **réparties uniformément** sur  $[1, n]$ . Puis, tirer au hasard l'une d'entre elles et l'utiliser pour réparer  $c_i$ .



Une **erreur** est un symbole dégradé d'un mot dont on ne connaît pas la position.

Certains codes sont munis d'**algorithmes de décodage** efficaces.

Les algorithmes de décodage classiques nécessitent, pour corriger une erreur, de **lire entièrement** le mot erroné  $\implies$  complexité  $\Omega(n)$ .

**Question légitime** : décodage (probabiliste) d'erreur en complexité sous-linéaire ?

**Idée (informelle)** : construire un code avec beaucoup d'équations de parité de même petit poids  $\ell \ll n$ , **réparties uniformément** sur  $[1, n]$ . Puis, tirer au hasard l'une d'entre elles et l'utiliser pour réparer  $c_i$ .



Une **erreur** est un symbole dégradé d'un mot dont on ne connaît pas la position.

Certains codes sont munis d'**algorithmes de décodage** efficaces.

Les algorithmes de décodage classiques nécessitent, pour corriger une erreur, de **lire entièrement** le mot erroné  $\implies$  complexité  $\Omega(n)$ .

**Question légitime** : décodage (probabiliste) d'erreur en complexité sous-linéaire ?

**Idée (informelle)** : construire un code avec beaucoup d'équations de parité de même petit poids  $\ell \ll n$ , **réparties uniformément** sur  $[1, n]$ . Puis, tirer au hasard l'une d'entre elles et l'utiliser pour réparer  $c_i$ .



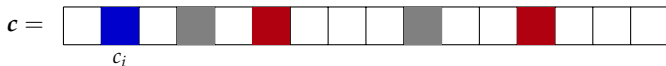
Une **erreur** est un symbole dégradé d'un mot dont on ne connaît pas la position.

Certains codes sont munis d'**algorithmes de décodage** efficaces.

Les algorithmes de décodage classiques nécessitent, pour corriger une erreur, de **lire entièrement** le mot erroné  $\implies$  complexité  $\Omega(n)$ .

**Question légitime** : décodage (probabiliste) d'erreur en complexité sous-linéaire ?

**Idée (informelle)** : construire un code avec beaucoup d'équations de parité de même petit poids  $\ell \ll n$ , **réparties uniformément** sur  $[1, n]$ . Puis, tirer au hasard l'une d'entre elles et l'utiliser pour réparer  $c_i$ .



Une **erreur** est un symbole dégradé d'un mot dont on ne connaît pas la position.

Certains codes sont munis d'**algorithmes de décodage** efficaces.

Les algorithmes de décodage classiques nécessitent, pour corriger une erreur, de **lire entièrement** le mot erroné  $\implies$  complexité  $\Omega(n)$ .

**Question légitime** : décodage (probabiliste) d'erreur en complexité sous-linéaire ?

**Idée (informelle)** : construire un code avec beaucoup d'équations de parité de même petit poids  $\ell \ll n$ , **réparties uniformément** sur  $[1, n]$ . Puis, tirer au hasard l'une d'entre elles et l'utiliser pour réparer  $c_i$ .



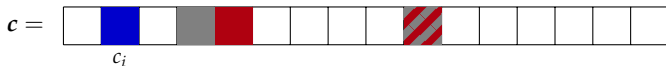
Une **erreur** est un symbole dégradé d'un mot dont on ne connaît pas la position.

Certains codes sont munis d'**algorithmes de décodage** efficaces.

Les algorithmes de décodage classiques nécessitent, pour corriger une erreur, de **lire entièrement** le mot erroné  $\implies$  complexité  $\Omega(n)$ .

**Question légitime** : décodage (probabiliste) d'erreur en complexité sous-linéaire ?

**Idée (informelle)** : construire un code avec beaucoup d'équations de parité de même petit poids  $\ell \ll n$ , **réparties uniformément** sur  $[1, n]$ . Puis, tirer au hasard l'une d'entre elles et l'utiliser pour réparer  $c_i$ .



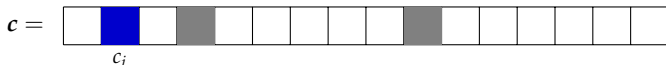
Une **erreur** est un symbole dégradé d'un mot dont on ne connaît pas la position.

Certains codes sont munis d'**algorithmes de décodage** efficaces.

Les algorithmes de décodage classiques nécessitent, pour corriger une erreur, de **lire entièrement** le mot erroné  $\implies$  complexité  $\Omega(n)$ .

**Question légitime** : décodage (probabiliste) d'erreur en complexité sous-linéaire ?

**Idée (informelle)** : construire un code avec beaucoup d'équations de parité de même petit poids  $\ell \ll n$ , **réparties uniformément** sur  $[1, n]$ . Puis, tirer au hasard l'une d'entre elles et l'utiliser pour réparer  $c_i$ .



**Résultats** : avec probabilité de réussite constante,

- ▶ on peut le faire en temps  $O_n(1)$  ( $\simeq$  théorème PCP), mais alors  $n = \exp(k)$ ,
- ▶ on peut le faire en temps  $O_n(2^{\sqrt{\log n \log \log n}})$ , avec  $n = O(k)$



Messages  $\mathbf{m} = (m_1, \dots, m_k) \in \{0, 1\}^k$ .

Messages  $m = (m_1, \dots, m_k) \in \{0, 1\}^k$ .

Soit  $\{0, 1\}^k =: \{a^{(1)}, \dots, a^{(2^k)}\}$

Messages  $m = (m_1, \dots, m_k) \in \{0, 1\}^k$ .

Soit  $\{0, 1\}^k =: \{a^{(1)}, \dots, a^{(2^k)}\}$

**Définition.** Le code de Hadamard est défini par l'encodage

$$m \mapsto c \in \{0, 1\}^{2^k} \quad \text{où} \quad c_i = \sum_{j=1}^k a_j^{(i)} m_j$$

C'est un code de dimension  $k$  et de longueur  $n = 2^k$ .

Messages  $m = (m_1, \dots, m_k) \in \{0, 1\}^k$ .

Soit  $\{0, 1\}^k =: \{a^{(1)}, \dots, a^{(2^k)}\}$

**Définition.** Le code de Hadamard est défini par l'encodage

$$m \mapsto c \in \{0, 1\}^{2^k} \quad \text{où} \quad c_i = \sum_{j=1}^k a_j^{(i)} m_j$$

C'est un code de dimension  $k$  et de longueur  $n = 2^k$ .

**Propriété importante :**

$$\text{si } a^{(i_1)} = a^{(i_2)} + a^{(i_3)}, \quad \text{alors } c_{i_1} = c_{i_2} + c_{i_3}$$

Messages  $m = (m_1, \dots, m_k) \in \{0, 1\}^k$ .

Soit  $\{0, 1\}^k =: \{a^{(1)}, \dots, a^{(2^k)}\}$

**Définition.** Le code de Hadamard est défini par l'encodage

$$m \mapsto c \in \{0, 1\}^{2^k} \quad \text{où} \quad c_i = \sum_{j=1}^k a_j^{(i)} m_j$$

C'est un code de dimension  $k$  et de longueur  $n = 2^k$ .

**Propriété importante :**

$$\text{si } a^{(i_1)} = a^{(i_2)} + a^{(i_3)}, \quad \text{alors } c_{i_1} = c_{i_2} + c_{i_3}$$

**Conséquence :** il existe des équations de parité de poids  $\ell = 3$ , réparties uniformément sur  $[1, n]$ .

1. Retrait confidentiel d'information

2. Codes correcteurs

3. Des schémas de PIR efficaces en temps et en mémoire

Avec une base de données répliquée

Avec une base de données distribuée

Résistance aux coalitions

1. Retrait confidentiel d'information

2. Codes correcteurs

3. Des schémas de PIR efficaces en temps et en mémoire

Avec une base de données répliquée

Avec une base de données distribuée

Résistance aux coalitions

## Paramètres globaux :

- $\mathcal{C} \subseteq \mathbb{F}^n$  un code localement décodable de localité  $\ell$
- $\mathcal{H} = \bigcup_{i=1}^n \mathcal{H}^{(i)} \subseteq \mathcal{C}^\perp$ , où  $\mathcal{H}^{(i)}$  est un sous-ensemble d'équations de parité de poids  $\ell$ , dont le support contient  $i$ , réparties uniformément sur  $[1, n]$ .



## Paramètres globaux :

- $\mathcal{C} \subseteq \mathbb{F}^n$  un code localement décodable de localité  $\ell$
- $\mathcal{H} = \bigcup_{i=1}^n \mathcal{H}^{(i)} \subseteq \mathcal{C}^\perp$ , où  $\mathcal{H}^{(i)}$  est un sous-ensemble d'équations de parité de poids  $\ell$ , dont le support contient  $i$ , réparties uniformément sur  $[1, n]$ .

**Initialisation :** la base de donnée  $F$  est copiée sur  $\ell - 1$  serveurs qui ne coopèrent pas.

## Paramètres globaux :

- $\mathcal{C} \subseteq \mathbb{F}^n$  un code localement décodable de localité  $\ell$
- $\mathcal{H} = \bigcup_{i=1}^n \mathcal{H}^{(i)} \subseteq \mathcal{C}^\perp$ , où  $\mathcal{H}^{(i)}$  est un sous-ensemble d'équations de parité de poids  $\ell$ , dont le support contient  $i$ , réparties uniformément sur  $[1, n]$ .

**Initialisation :** la base de donnée  $F$  est copiée sur  $\ell - 1$  serveurs qui ne coopèrent pas.

**Protocole** pour extraire  $F_i = c_i$ .

## Paramètres globaux :

- $\mathcal{C} \subseteq \mathbb{F}^n$  un code localement décodable de localité  $\ell$
- $\mathcal{H} = \bigcup_{i=1}^n \mathcal{H}^{(i)} \subseteq \mathcal{C}^\perp$ , où  $\mathcal{H}^{(i)}$  est un sous-ensemble d'équations de parité de poids  $\ell$ , dont le support contient  $i$ , réparties uniformément sur  $[1, n]$ .

**Initialisation :** la base de donnée  $F$  est copiée sur  $\ell - 1$  serveurs qui ne coopèrent pas.

**Protocole** pour extraire  $F_i = c_i$ .

1. Pour définir  $\text{Query}(i)$  :

- on tire aléatoirement  $\mathbf{h} \leftarrow \mathcal{H}^{(i)}$ , avec  $\text{supp}(\mathbf{h}) = \{h_i, h_{m_1}, \dots, h_{m_{\ell-1}}\}$ ,
- on définit  $\text{Query}(i)_j = m_j$ .

## Paramètres globaux :

- $\mathcal{C} \subseteq \mathbb{F}^n$  un code localement décodable de localité  $\ell$
- $\mathcal{H} = \bigcup_{i=1}^n \mathcal{H}^{(i)} \subseteq \mathcal{C}^\perp$ , où  $\mathcal{H}^{(i)}$  est un sous-ensemble d'équations de parité de poids  $\ell$ , dont le support contient  $i$ , réparties uniformément sur  $[1, n]$ .

**Initialisation :** la base de donnée  $F$  est copiée sur  $\ell - 1$  serveurs qui ne coopèrent pas.

**Protocole** pour extraire  $F_i = c_i$ .

1. Pour définir  $\text{Query}(i)$  :

- on tire aléatoirement  $\mathbf{h} \leftarrow \mathcal{H}^{(i)}$ , avec  $\text{supp}(\mathbf{h}) = \{h_i, h_{m_1}, \dots, h_{m_{\ell-1}}\}$ ,
- on définit  $\text{Query}(i)_j = m_j$ .

2. Chaque serveur  $S_j$  retourne  $a_j = \text{Enc}_{\mathcal{C}}(F)_{m_j}$ .

## Paramètres globaux :

- $\mathcal{C} \subseteq \mathbb{F}^n$  un code localement décodable de localité  $\ell$
- $\mathcal{H} = \bigcup_{i=1}^n \mathcal{H}^{(i)} \subseteq \mathcal{C}^\perp$ , où  $\mathcal{H}^{(i)}$  est un sous-ensemble d'équations de parité de poids  $\ell$ , dont le support contient  $i$ , réparties uniformément sur  $[1, n]$ .

**Initialisation :** la base de donnée  $F$  est copiée sur  $\ell - 1$  serveurs qui ne coopèrent pas.

**Protocole** pour extraire  $F_i = c_i$ .

1. Pour définir  $\text{Query}(i)$  :

- on tire aléatoirement  $\mathbf{h} \leftarrow \mathcal{H}^{(i)}$ , avec  $\text{supp}(\mathbf{h}) = \{h_i, h_{m_1}, \dots, h_{m_{\ell-1}}\}$ ,
- on définit  $\text{Query}(i)_j = m_j$ .

2. Chaque serveur  $S_j$  retourne  $a_j = \text{Enc}_{\mathcal{C}}(F)_{m_j}$ .

3. User reconstruit  $F_i = -\frac{1}{h_i} \sum_j h_{m_j} a_j$ .

**Théorème (informel, Katz-Trevisan 00).** Tout code localement décodable de localité  $\ell$  induit un protocole de PIR à  $\ell - 1$  serveurs.

**Théorème (informel, Katz-Trevisan 00).** Tout code localement décodable de localité  $\ell$  induit un protocole de PIR à  $\ell - 1$  serveurs.

## Grandeurs.

- ▶ la base de données est copiée sur les  $\ell - 1$  serveurs (coûteux)
- ▶ à chaque exécution, chaque serveur doit recalculer le symbole  $\text{Enc}_{\mathcal{C}}(F)_{m_j}$  (très coûteux)

**Théorème (informel, Katz-Trevisan 00).** Tout code localement décodable de localité  $\ell$  induit un protocole de PIR à  $\ell - 1$  serveurs.

## Grandeurs.

- ▶ la base de données est copiée sur les  $\ell - 1$  serveurs (coûteux)
- ▶ à chaque exécution, chaque serveur doit recalculer le symbole  $\text{Enc}_{\mathcal{C}}(F)_{m_j}$  (très coûteux)

**Solution “classique” :** compromis temps/mémoire par un précalcul

- ▶ le coût supplémentaire de stockage dépend de la dimension du code
- ▶ souvent très coûteux en mémoire



**Théorème (informel, Katz-Trevisan 00).** Tout code localement décodable de localité  $\ell$  induit un protocole de PIR à  $\ell - 1$  serveurs.

## Grandeurs.

- ▶ la base de données est copiée sur les  $\ell - 1$  serveurs (coûteux)
- ▶ à chaque exécution, chaque serveur doit recalculer le symbole  $\text{Enc}_{\mathcal{C}}(F)_{m_j}$  (très coûteux)

**Solution “classique” :** compromis temps/mémoire par un précalcul

- ▶ le coût supplémentaire de stockage dépend de la dimension du code
- ▶ souvent très coûteux en mémoire

**Nouvelle idée** (initiée par Augot, Levy-dit-Vehel, Shikfa, 2014) :

- ▶ partitionner l'encodage de la base de données sur les serveurs,
- ▶ nécessite de concevoir un code  $\mathcal{C}$  dont l'ensemble d'équations de parité de petits poids  $\mathcal{H}$  puisse correspondre à cette partition.

1. Retrait confidentiel d'information

2. Codes correcteurs

3. Des schémas de PIR efficaces en temps et en mémoire

Avec une base de données répliquée

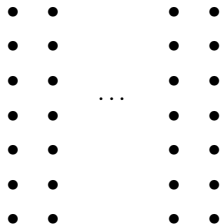
Avec une base de données distribuée

Résistance aux coalitions

Pour concevoir ces codes, un peu de combinatoire...

Un **design transversal**  $\text{TD}(n, s) = (X, \mathcal{B}, \mathcal{G})$  est donné par :

- ▶  $X$  un ensemble de *points*,  $|X| = N = ns$ ,

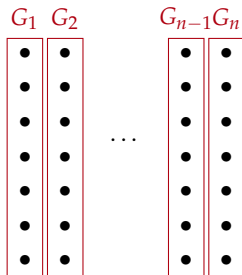


Pour concevoir ces codes, un peu de combinatoire...

Un **design transversal**  $TD(n, s) = (X, \mathcal{B}, \mathcal{G})$  est donné par :

- ▶  $X$  un ensemble de *points*,  $|X| = N = ns$ ,
- ▶ une partition de  $X$  en *groupes*  $\mathcal{G} = \{G_j\}_{1 \leq j \leq n}$

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



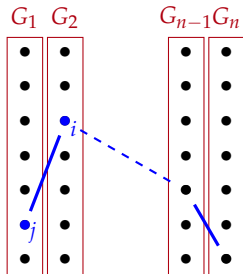
Pour concevoir ces codes, un peu de combinatoire...

Un **design transversal**  $TD(n, s) = (X, \mathcal{B}, \mathcal{G})$  est donné par :

- ▶  $X$  un ensemble de *points*,  $|X| = N = ns$ ,
- ▶ une partition de  $X$  en *groupes*  $\mathcal{G} = \{G_j\}_{1 \leq j \leq n}$

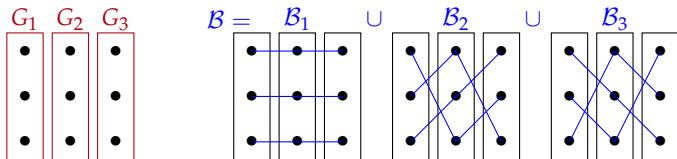
$$X = \bigsqcup_{j=1}^n G_j \text{ et } |G_j| = s,$$

- ▶ un ensemble de *blocs*  $B \in \mathcal{B}$  tels que
  - $B \subset X$  et  $|B| = n$ ;
  - toute paire  $\{i, j\} \subset X$  est incluse :  
 ou bien dans un unique groupe  $G \in \mathcal{G}$ ,  
 ou bien dans un unique bloc  $B \in \mathcal{B}$



Exemple pour TD(3,3) :

- $ns = 9$  points
- $s = 3$  groupes  $G_1, G_2, G_3$  de cardinal 3
- $ns = 9$  blocs de  $n = 3$  points (ici présentées en 3 classes  $B_1, B_2, B_3$ )



Soit  $\mathcal{T}$  un design transversal  $\text{TD}(n, s) = (X, \mathcal{B}, \mathcal{G})$ .

Sa **matrice d'incidence**  $M$ , de taille  $|\mathcal{B}| \times |X| = ns \times ns$ , est définie comme :

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

Soit  $\mathcal{T}$  un design transversal  $\text{TD}(n, s) = (X, \mathcal{B}, \mathcal{G})$ .

Sa **matrice d'incidence**  $M$ , de taille  $|\mathcal{B}| \times |X| = ns \times ns$ , est définie comme :

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

**Définition.** Le **code  $\mathcal{C}$  construit sur le design  $\mathcal{T}$  sur le corps  $\mathbb{F}$**  est le code  $\mathbb{F}$ -linéaire qui admet  $M$  comme matrice de parité.

- la longueur du code est  $\text{length}(\mathcal{C}) = |X| = ns$ ,
- la dimension du code est  $\dim(\mathcal{C}) = \dim(\ker M)$ ,
- chaque bloc  $B \in \mathcal{B}$  produit une équation de parité  $\mathbf{h} \in \mathbb{F}^{ns}$ , telle que

$$\text{wt}(\mathbf{h}|_{\mathcal{G}_j}) = 1, \quad \forall j = 1, \dots, n$$



# Example

On reprend TD(3,3) :

mot de code  $c \in \mathbb{F}^9$

$c_1$	$c_2$	$c_3$
$c_4$	$c_5$	$c_6$
$c_7$	$c_8$	$c_9$

$\mathcal{B} = \mathcal{B}_1 \cup$

$c_1$	$c_2$	$c_3$
$c_4$	$c_5$	$c_6$
$c_7$	$c_8$	$c_9$

$\cup \mathcal{B}_2 \cup$

$c_1$	$c_2$	$c_3$
$c_4$	$c_5$	$c_6$
$c_7$	$c_8$	$c_9$

$\cup \mathcal{B}_3$

$c_1$	$c_2$	$c_3$
$c_4$	$c_5$	$c_6$
$c_7$	$c_8$	$c_9$

La matrice d'incidence du design (= matrice de parité du code) est :

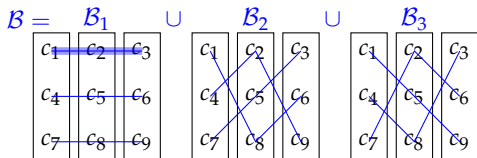
$$H = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 &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

# Example

On reprend TD(3,3) :

mot de code  $c \in \mathbb{F}^9$

$c_1$	$c_2$	$c_3$
$c_4$	$c_5$	$c_6$
$c_7$	$c_8$	$c_9$



La matrice d'incidence du design (= matrice de parité du code) est :

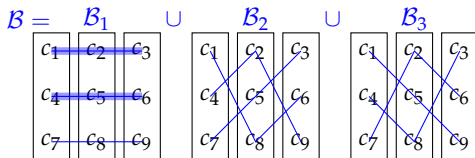
$$H = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 &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

# Example

On reprend TD(3,3) :

mot de code  $c \in \mathbb{F}^9$

$c_1$	$c_2$	$c_3$
$c_4$	$c_5$	$c_6$
$c_7$	$c_8$	$c_9$



La matrice d'incidence du design (= matrice de parité du code) est :

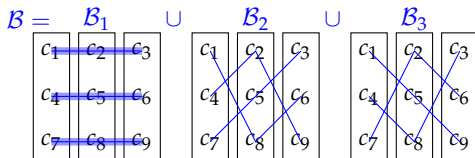
$$H = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 &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

# Example

On reprend TD(3,3) :

mot de code  $c \in \mathbb{F}^9$

$c_1$	$c_2$	$c_3$
$c_4$	$c_5$	$c_6$
$c_7$	$c_8$	$c_9$



La matrice d'incidence du design (= matrice de parité du code) est :

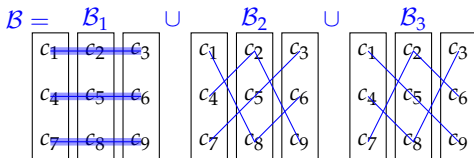
$$H = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 &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

# Example

On reprend TD(3,3) :

mot de code  $c \in \mathbb{F}^9$

$c_1$	$c_2$	$c_3$
$c_4$	$c_5$	$c_6$
$c_7$	$c_8$	$c_9$



La matrice d'incidence du design (= matrice de parité du code) est :

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

Le noyau de la matrice dépend du corps  $\mathbb{F}$  considéré.

Cas le plus favorable : sur  $\mathbb{F}_3 := \mathbb{Z}/3\mathbb{Z}$ , on a  $\text{rk}(H) = 6 \implies \dim_{\mathbb{F}_3}(\mathcal{C}) = 3$ .

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

Soit  $\mathcal{C} \subseteq \mathbb{F}_q^{ns}$  un code construit sur un TD( $n, s$ ).

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

Soit  $\mathcal{C} \subseteq \mathbb{F}_q^{ns}$  un code construit sur un TD( $n, s$ ).

- **Initialisation.** User encode le message  $m = F$  en un mot du code  $c \in \mathcal{C}$ .  
Puis, User distribue  $c|_{G_j}$  au serveur  $S_j$ .

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

Soit  $\mathcal{C} \subseteq \mathbb{F}_q^{ns}$  un code construit sur un TD( $n, s$ ).

• **Initialisation.** User encode le message  $m = F$  en un mot du code  $c \in \mathcal{C}$ .  
Puis, User distribue  $c|_{G_j}$  au serveur  $S_j$ .

• **Pour extraire**  $F_i = c_i$  où  $i \in X$ :

1. User tire aléatoirement un bloc  $B \in \mathcal{B}$  qui contient  $i$ .

Puis, User définit la requête :

$$q_j = \text{Query}(i)_j := \begin{cases} \text{unique élément de } B \cap G_j & \text{si } i \notin G_j \\ \text{un point aléatoire de } G_j & \text{sinon.} \end{cases}$$



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

Soit  $\mathcal{C} \subseteq \mathbb{F}_q^{ns}$  un code construit sur un TD( $n, s$ ).

• **Initialisation.** User encode le message  $m = F$  en un mot du code  $c \in \mathcal{C}$ .  
Puis, User distribue  $c|_{G_j}$  au serveur  $S_j$ .

• **Pour extraire**  $F_i = c_i$  où  $i \in X$ :

1. User tire aléatoirement un bloc  $B \in \mathcal{B}$  qui contient  $i$ .

Puis, User définit la requête :

$$q_j = \text{Query}(i)_j := \begin{cases} \text{unique élément de } B \cap G_j & \text{si } i \notin G_j \\ \text{un point aléatoire de } G_j & \text{sinon.} \end{cases}$$

2. Chaque serveur  $S_j$  retourne  $c_{q_j}$

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

Soit  $\mathcal{C} \subseteq \mathbb{F}_q^{ns}$  un code construit sur un TD( $n, s$ ).

• **Initialisation.** User encode le message  $m = F$  en un mot du code  $c \in \mathcal{C}$ .  
Puis, User distribue  $c|_{G_j}$  au serveur  $S_j$ .

• **Pour extraire**  $F_i = c_i$  où  $i \in X$ :

1. User tire aléatoirement un bloc  $B \in \mathcal{B}$  qui contient  $i$ .  
Puis, User définit la requête :

$$q_j = \text{Query}(i)_j := \begin{cases} \text{unique élément de } B \cap G_j & \text{si } i \notin G_j \\ \text{un point aléatoire de } G_j & \text{sinon.} \end{cases}$$

2. Chaque serveur  $S_j$  retourne  $c_{q_j}$
3. User reconstruit

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

**Théorème.** Ce protocole PIR est correct et sûr du point de vue IT.

Preuve :

- le seul serveur qui détient  $F_i$  reçoit une requête aléatoire ;
- pour tout autre serveur  $S_j$ , la requête  $q_j$  ne fournit aucune information sur le bloc  $B$  qui a été choisi par User  $\Rightarrow$  aucune information sur  $i$  n'est révélée.

**Théorème.** Ce protocole PIR est correct et sûr du point de vue IT.

Preuve :

- le seul serveur qui détient  $F_i$  reçoit une requête aléatoire ;
- pour tout autre serveur  $S_j$ , la requête  $q_j$  ne fournit aucune information sur le bloc  $B$  qui a été choisi par User  $\Rightarrow$  aucune information sur  $i$  n'est révélée.

**Caractéristiques.**

- ▶ complexité de communication :  $n \log s$  bits envoyés,  $n \log q$  bits reçus

**Théorème.** Ce protocole PIR est correct et sûr du point de vue IT.

Preuve :

- le seul serveur qui détient  $F_i$  reçoit une requête aléatoire ;
- pour tout autre serveur  $S_j$ , la requête  $q_j$  ne fournit aucune information sur le bloc  $B$  qui a été choisi par User  $\Rightarrow$  aucune information sur  $i$  n'est révélée.

**Caractéristiques.**

- ▶ complexité de communication :  $n \log s$  bits envoyés,  $n \log q$  bits reçus
- ▶ complexité de calcul :
  - ▶ **seulement 1 lecture pour chaque serveur (optimal)**
  - ▶  $\leq n$  additions sur  $\mathbb{F}$  pour User

**Théorème.** Ce protocole PIR est correct et sûr du point de vue IT.

Preuve :

- le seul serveur qui détient  $F_i$  reçoit une requête aléatoire ;
- pour tout autre serveur  $S_j$ , la requête  $q_j$  ne fournit aucune information sur le bloc  $B$  qui a été choisi par User  $\Rightarrow$  aucune information sur  $i$  n'est révélée.

**Caractéristiques.**

- ▶ complexité de communication :  $n \log s$  bits envoyés,  $n \log q$  bits reçus
- ▶ complexité de calcul :
  - ▶ **seulement 1 lecture pour chaque serveur (optimal)**
  - ▶  $\leq n$  additions sur  $\mathbb{F}$  pour User
- ▶ stockage : redondance de  $(ns - k) \log q$  bits, où  $k = \dim_{\mathbb{F}}(\mathcal{C})$

**Théorème.** Ce protocole PIR est correct et sûr du point de vue IT.

Preuve :

- le seul serveur qui détient  $F_i$  reçoit une requête aléatoire ;
- pour tout autre serveur  $S_j$ , la requête  $q_j$  ne fournit aucune information sur le bloc  $B$  qui a été choisi par User  $\Rightarrow$  aucune information sur  $i$  n'est révélée.

**Caractéristiques.**

- ▶ complexité de communication :  $n \log s$  bits envoyés,  $n \log q$  bits reçus
- ▶ complexité de calcul :
  - ▶ **seulement 1 lecture pour chaque serveur (optimal)**
  - ▶  $\leq n$  additions sur  $\mathbb{F}$  pour User
- ▶ stockage : redondance de  $(ns - k) \log q$  bits, où  $k = \dim_{\mathbb{F}}(\mathcal{C})$

**Question :** quels designs transversaux mènent à des codes de grande dimension ?

Design issu de la **géométrie affine** sur le corps  $\mathbb{F} = \mathbb{F}_q$ .

- ▶  $X = \mathbb{F}_q^m$  pour  $m \geq 2$  fixé,
- ▶  $\mathcal{G}$  une partition de  $X$  en  $q$  hyperplans  $G_1, \dots, G_q$ ,
- ▶  $\mathcal{B} = \{\text{droites affines } L \text{ sécantes à chaque hyperplan } G_j\}$ .



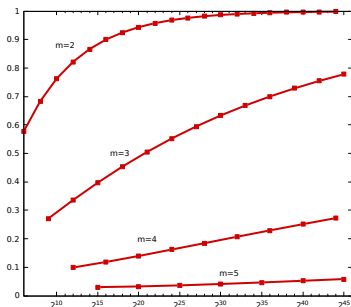
Design issu de la **géométrie affine** sur le corps  $\mathbb{F} = \mathbb{F}_q$ .

- ▶  $X = \mathbb{F}_q^m$  pour  $m \geq 2$  fixé,
- ▶  $\mathcal{G}$  une partition de  $X$  en  $q$  hyperplans  $G_1, \dots, G_q$ ,
- ▶  $\mathcal{B} = \{\text{droites affines } L \text{ sécantes à chaque hyperplan } G_j\}$ .

Ce code a :

- longueur  $ns = q^m$ ,
- "localité"  $n = q$ .

taux d'information  $k/N$



longueur  $N = ns = 2^{em}$

1. Retrait confidentiel d'information

2. Codes correcteurs

3. Des schémas de PIR efficaces en temps et en mémoire

Avec une base de données répliquée

Avec une base de données distribuée

Résistance aux coalitions

**Définition.** Le code de **Reed–Solomon** de dimension  $k$  sur le corps fini  $\mathbb{F}_q$  est

$$\text{RS}_q(k) := \{ \text{ev}_1(f) := (f(x_1), \dots, f(x_q)) \mid f \in \mathbb{F}_q[X] \text{ et } \deg(f) \leq k - 1 \}.$$

**Définition.** Le code de **Reed–Solomon** de dimension  $k$  sur le corps fini  $\mathbb{F}_q$  est

$$\text{RS}_q(k) := \{ \text{ev}_1(f) := (f(x_1), \dots, f(x_q)) \mid f \in \mathbb{F}_q[X] \text{ et } \deg(f) \leq k - 1 \}.$$

**Propriété remarquable :** tout symbole  $c_i$  d'un mot de code  $c \in \text{RS}_q(k)$  peut être reconstruit à partir de n'importe quel sous-ensemble de  $k$  coordonnées de  $c$ .

**Définition.** Le code de **Reed–Solomon** de dimension  $k$  sur le corps fini  $\mathbb{F}_q$  est

$$\text{RS}_q(k) := \{ \text{ev}_1(f) := (f(x_1), \dots, f(x_q)) \mid f \in \mathbb{F}_q[X] \text{ et } \deg(f) \leq k - 1 \}.$$

**Propriété remarquable :** tout symbole  $c_i$  d'un mot de code  $c \in \text{RS}_q(k)$  peut être reconstruit à partir de n'importe quel sous-ensemble de  $k$  coordonnées de  $c$ .

**Définition.** Le code de **Reed–Muller** d'ordre  $m$  et de degré  $r$  sur  $\mathbb{F}_q$  est

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

**Définition.** Le code de **Reed–Solomon** de dimension  $k$  sur le corps fini  $\mathbb{F}_q$  est

$$\text{RS}_q(k) := \{ \text{ev}_1(f) := (f(x_1), \dots, f(x_q)) \mid f \in \mathbb{F}_q[X] \text{ et } \deg(f) \leq k-1 \}.$$

**Propriété remarquable :** tout symbole  $c_i$  d'un mot de code  $c \in \text{RS}_q(k)$  peut être reconstruit à partir de n'importe quel sous-ensemble de  $k$  coordonnées de  $c$ .

**Définition.** Le code de **Reed–Muller** d'ordre  $m$  et de degré  $r$  sur  $\mathbb{F}_q$  est

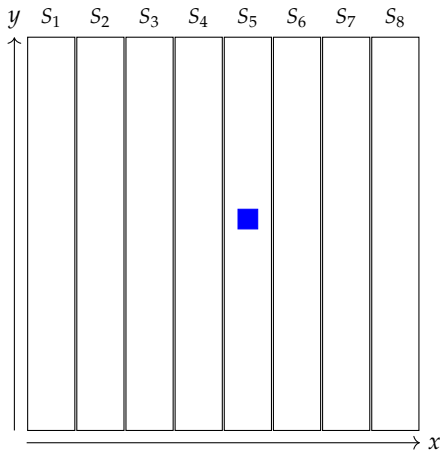
$$\text{RM}_q(m, r) := \{ \text{ev}_m(f) \mid f \in \mathbb{F}_q[X] \text{ et } \deg(f) \leq r \}.$$

**Propriété remarquable :**

$$\forall c = \text{ev}_m(f) \in \text{RM}_q(m, r) \quad \text{et} \quad \forall \text{ droite affine } L \subset \mathbb{F}_q^m, \\ \text{ev}_1(f|_L) \in \text{RS}_q(r+1).$$

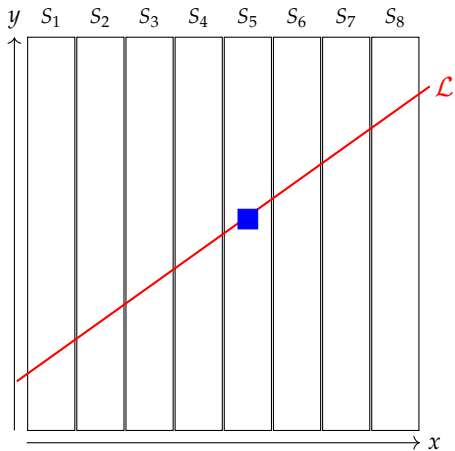
(où  $f|_L$  est le polynôme interpolateur de  $f$  sur  $L$ )

# Un protocole de PIR construit sur les codes de Reed–Muller



La base de données  $F$  est encodée avec  $\text{RM}_q(m, r)$ , puis distribuée sur les serveurs selon une partition en hyperplans.

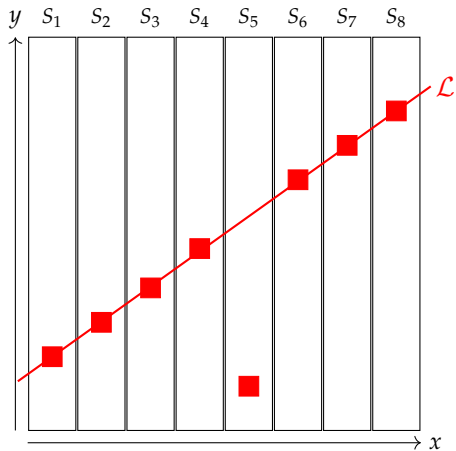
# Un protocole de PIR construit sur les codes de Reed–Muller



La base de données  $F$  est encodée avec  $\text{RM}_q(m, r)$ , puis distribuée sur les serveurs selon une partition en hyperplans.

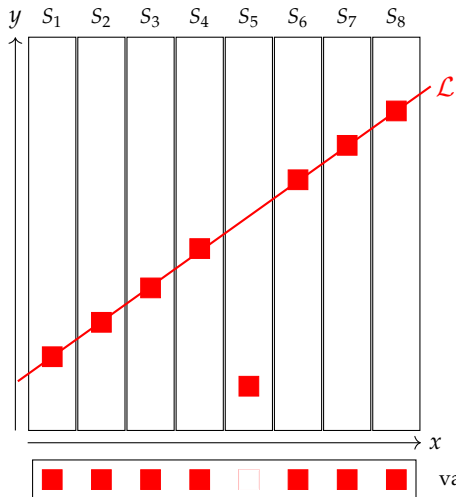


# Un protocole de PIR construit sur les codes de Reed–Muller



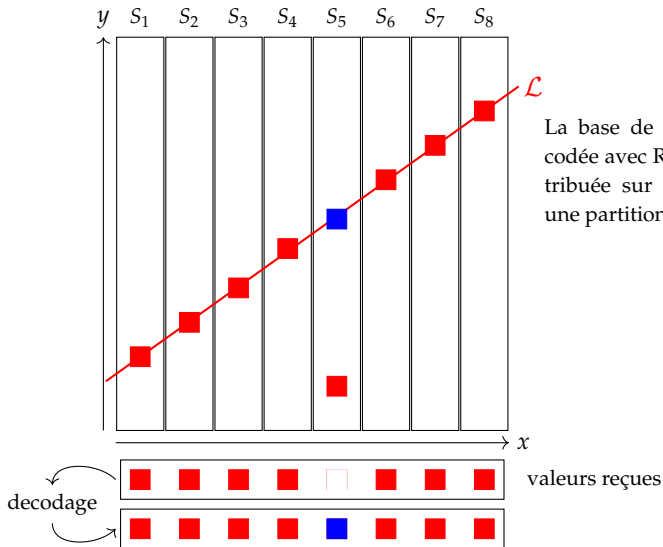
La base de données  $F$  est encodée avec  $\text{RM}_q(m, r)$ , puis distribuée sur les serveurs selon une partition en hyperplans.

# Un protocole de PIR construit sur les codes de Reed–Muller



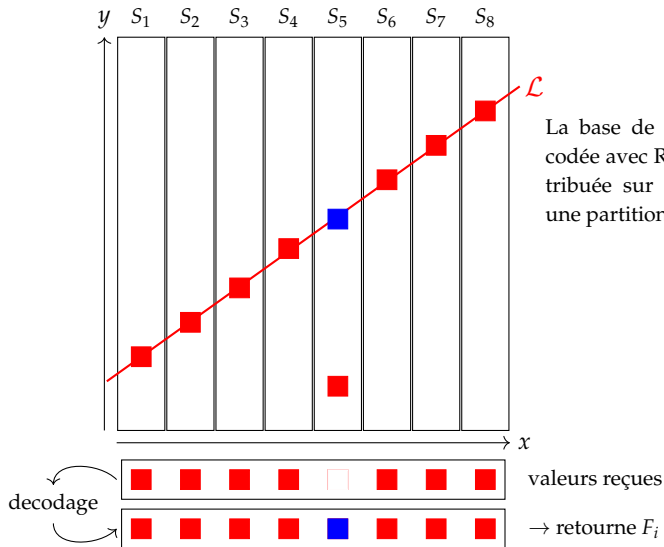
La base de données  $F$  est encodée avec  $RM_q(m, r)$ , puis distribuée sur les serveurs selon une partition en hyperplans.

# Un protocole de PIR construit sur les codes de Reed–Muller

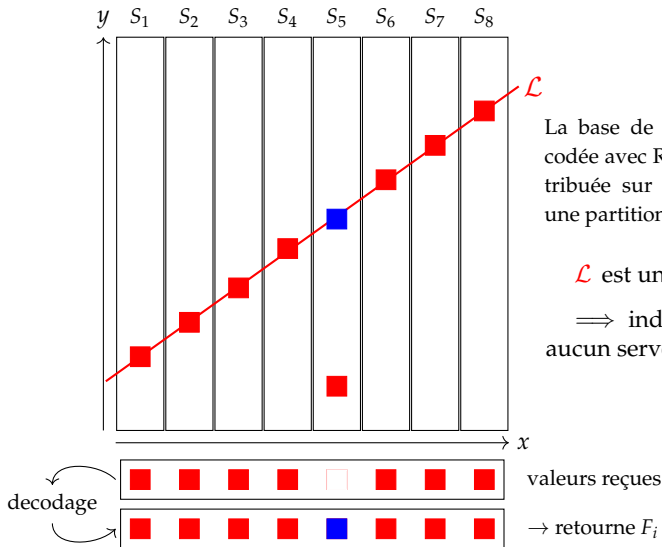


La base de données  $F$  est encodée avec  $RM_q(m, r)$ , puis distribuée sur les serveurs selon une partition en hyperplans.

# Un protocole de PIR construit sur les codes de Reed–Muller



# Un protocole de PIR construit sur les codes de Reed–Muller



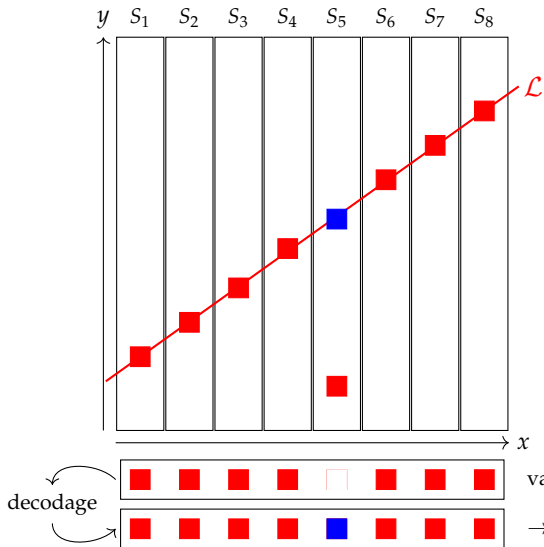
La base de données  $F$  est encodée avec  $RM_q(m, r)$ , puis distribuée sur les serveurs selon une partition en hyperplans.

$\mathcal{L}$  est une droite affine

$\implies$  individuellement, aucun serveur ne sait où est



# Un protocole de PIR construit sur les codes de Reed–Muller



La base de données  $F$  est encodée avec  $RM_q(m, r)$ , puis distribuée sur les serveurs selon une partition en hyperplans.

$\mathcal{L}$  est une droite affine

⇒ individuellement, aucun serveur ne sait où est



⇒ le protocole est confidentiel


Pour le code  $\text{RM}_q(m, r)$  de longueur  $n = q^m$ .

- ▶ communication :  $(m - 1)q \log q$  bits envoyés,  $q \log q$  bits reçus
- ▶ calcul :
  - ▶ **1 lecture pour chaque serveur (optimal)**
  - ▶ un décodage pour l'utilisateur
- ▶ stockage : le ratio  $k/n$  du code  $\text{RM}_q(m, r)$  avec  $r \leq q - 1$  est

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


⇒ On voudrait un code avec des propriétés similaires, mais une meilleure dimension.

Les “codes relevés” sont les plus grands codes ayant les mêmes propriétés locales que les codes de Reed–Muller.

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




Les “codes relevés” sont les plus grands codes ayant les mêmes propriétés locales que les codes de Reed–Muller.

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

**Définition.** Le code de Reed–Solomon **relevé  $m$  fois** de degré  $r$  sur  $\mathbb{F}_q$  est :

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

Les “codes relevés” sont les plus grands codes ayant les mêmes propriétés locales que les codes de Reed–Muller.


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

**Définition.** Le code de Reed–Solomon **relevé  $m$  fois** de degré  $r$  sur  $\mathbb{F}_q$  est :

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

Ces codes sont parfois bien plus grands que les codes de Reed–Muller (c’est surprenant !).

Les “codes relevés” sont les plus grands codes ayant les mêmes propriétés locales que les codes de Reed–Muller.

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

**Définition.** Le code de Reed–Solomon **relevé  $m$  fois** de degré  $r$  sur  $\mathbb{F}_q$  est :

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

Ces codes sont parfois bien plus grands que les codes de Reed–Muller (c’est surprenant !).

**Exemple.** Pour  $m = 2, r = 2$ , on considère  $f(X, Y) = X^2Y^2$  sur  $\mathbb{F}_4$ .

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

Alors,

$$\text{ev}_2(X^2Y^2) \in \text{Lift}_4(2, 2) \quad \text{mais} \quad \text{ev}_2(X^2Y^2) \notin \text{RM}_4(2, 2).$$

Les “codes relevés” sont les plus grands codes ayant les mêmes propriétés locales que les codes de Reed–Muller.

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

**Définition.** Le code de Reed–Solomon **relevé  $m$  fois** de degré  $r$  sur  $\mathbb{F}_q$  est :

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

Ces codes sont parfois bien plus grands que les codes de Reed-Muller (c'est surprenant !).

**Exemple.** Pour  $m = 2, r = 2$ , on considère  $f(X, Y) = X^2Y^2$  sur  $\mathbb{F}_4$ .

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

Alors,

$$\text{ev}_2(X^2Y^2) \in \text{Lift}_4(2, 2) \quad \text{mais} \quad \text{ev}_2(X^2Y^2) \notin \text{RM}_4(2, 2).$$

**Prop.** Pour tout  $m$  fixé, les codes relevés peuvent atteindre un ratio  $k/n$  arbitrairement proche de 1.

**Question :** comment gérer les coalitions de serveurs ?

**Question :** comment gérer les coalitions de serveurs ?

Simplifions et prenons  $m = 2$ .

**Définition.** Une  $t$ -courbe est :

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

**Question :** comment gérer les coalitions de serveurs ?

Simplifions et prenons  $m = 2$ .

**Définition.** Une  $t$ -courbe est :

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

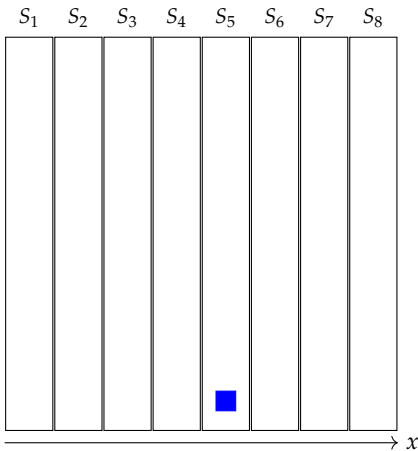
**Définition.** Le **relèvement pondéré du code de Reed–Solomon** de degré  $r$  et de poids  $t$  sur  $\mathbb{F}_q$  est :

$$\text{WLift}_q(t, r) := \{\text{ev}_2(f) \mid f \in \mathbb{F}_q[X, Y] \text{ et } \forall t\text{-courbe } \mathcal{L} \subset \mathbb{A}^2, \deg(f|_{\mathcal{L}}) \leq r\}$$

**Conséquence :** pour tout mot  $c \in \text{WLift}_q(t, r)$  et pour toute  $t$ -courbe  $\mathcal{L}$ , on a :

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

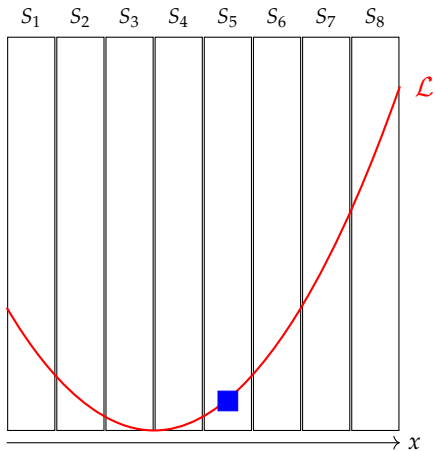
# Un protocole de PIR construit sur des relèvements pondérés



La base de données  $F$  est encodée avec le code  $WLift_q(t, r)$ , puis distribuée parmi les serveurs selon une partition en droites verticales

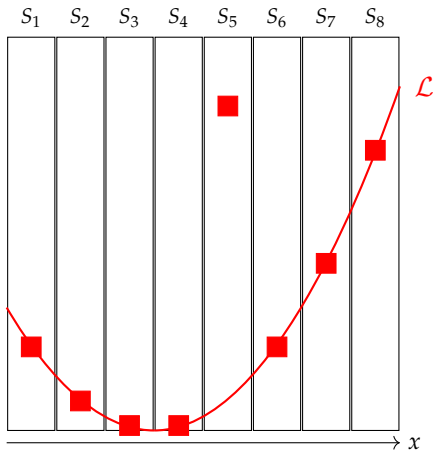


# Un protocole de PIR construit sur des relèvements pondérés



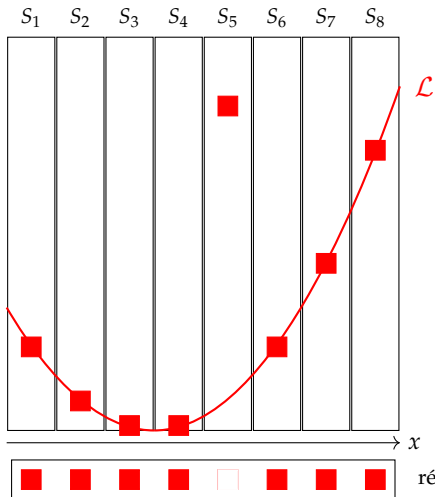
La base de données  $F$  est encodée avec le code  $\text{WLift}_q(t, r)$ , puis distribuée parmi les serveurs selon une partition en droites verticales

# Un protocole de PIR construit sur des relèvements pondérés



La base de données  $F$  est encodée avec le code  $\text{WLift}_q(t, r)$ , puis distribuée parmi les serveurs selon une partition en droites verticales

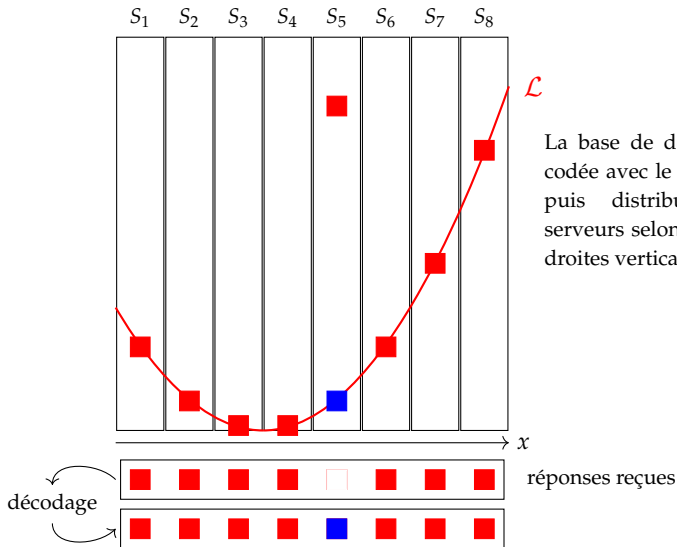
# Un protocole de PIR construit sur des relèvements pondérés



La base de données  $F$  est encodée avec le code  $WLift_q(t,r)$ , puis distribuée parmi les serveurs selon une partition en droites verticales

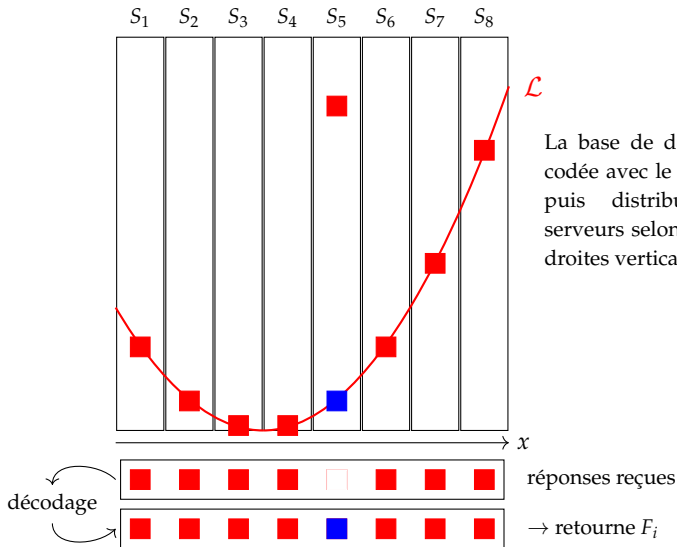
réponses reçues

# Un protocole de PIR construit sur des relèvements pondérés



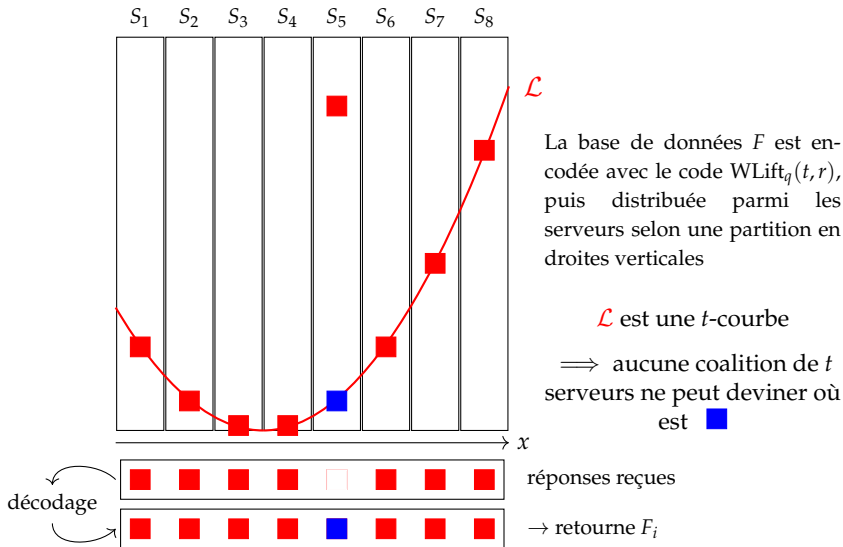
La base de données  $F$  est encodée avec le code  $WLift_q(t,r)$ , puis distribuée parmi les serveurs selon une partition en droites verticales


# Un protocole de PIR construit sur des relèvements pondérés



La base de données  $F$  est encodée avec le code  $\text{WLift}_q(t, r)$ , puis distribuée parmi les serveurs selon une partition en droites verticales


# Un protocole de PIR construit sur des relèvements pondérés



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

**Théorème.** Fixons  $p$  un nombre premier,  $t \geq 1$  et  $\alpha \geq 2$  des entiers. Alors, le ratio  $k/n$  du code  $\text{WLift}_{p^e}(t, p^e - \alpha)$  tend vers 1 pour  $e \rightarrow \infty$ .

**Corollaire :** pour un nombre de serveurs grandissant mais des tailles de coalitions constantes, on obtient des protocoles de PIR avec une redondance de stockage qui tend vers 0.

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

**Théorème.** Fixons  $p$  un nombre premier,  $t \geq 1$  et  $\alpha \geq 2$  des entiers. Alors, le ratio  $k/n$  du code  $\text{WLift}_{p^\alpha}(t, p^\alpha - \alpha)$  tend vers 1 pour  $e \rightarrow \infty$ .

**Corollaire :** pour un nombre de serveurs grandissant mais des tailles de coalitions constantes, on obtient des protocoles de PIR avec une redondance de stockage qui tend vers 0.

**Théorème.** Fixons  $p$  un nombre premier,  $t \geq 1$  et  $c \geq 1$  des entiers. Notons  $\gamma = 1 - p^{-c}$  et  $\mathcal{C}_e = \text{WLift}_{p^c}(t, \gamma p^e)$ . Alors, le ratio  $R_e = k/n$  du code  $\mathcal{C}_e$  vérifie

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

**Corollaire :** pour un nombre de serveurs grandissant avec une fraction constante de serveurs en coalition et/ou en échec, on obtient des protocoles de PIR avec une redondance de stockage bornée par une constante.



## Autres constructions de protocoles de PIR :

- ▶ sur des bases de données dont l'encodage n'est pas maîtrisé par l'utilisateur (RAID, Hadoop, Azure etc.),
- ▶ avec des patterns de coalitions particuliers,
- ▶ symétrique : l'utilisateur n'apprend pas plus que  $F_i$ ,
- ▶ etc.

## **Autres constructions de protocoles de PIR :**

- ▶ sur des bases de données dont l'encodage n'est pas maîtrisé par l'utilisateur (RAID, Hadoop, Azure etc.),
- ▶ avec des patterns de coalitions particuliers,
- ▶ symétrique : l'utilisateur n'apprend pas plus que  $F_i$ ,
- ▶ etc.

## **Autres applications de la notion de localité pour des protocoles cryptographiques :**

- ▶ preuves d'extractibilité d'une donnée,
- ▶ preuves (non-)interactives d'intégrité,
- ▶ etc.