

# Private Information Retrieval

Dima Kogan

# Private information retrieval [CGKS95]

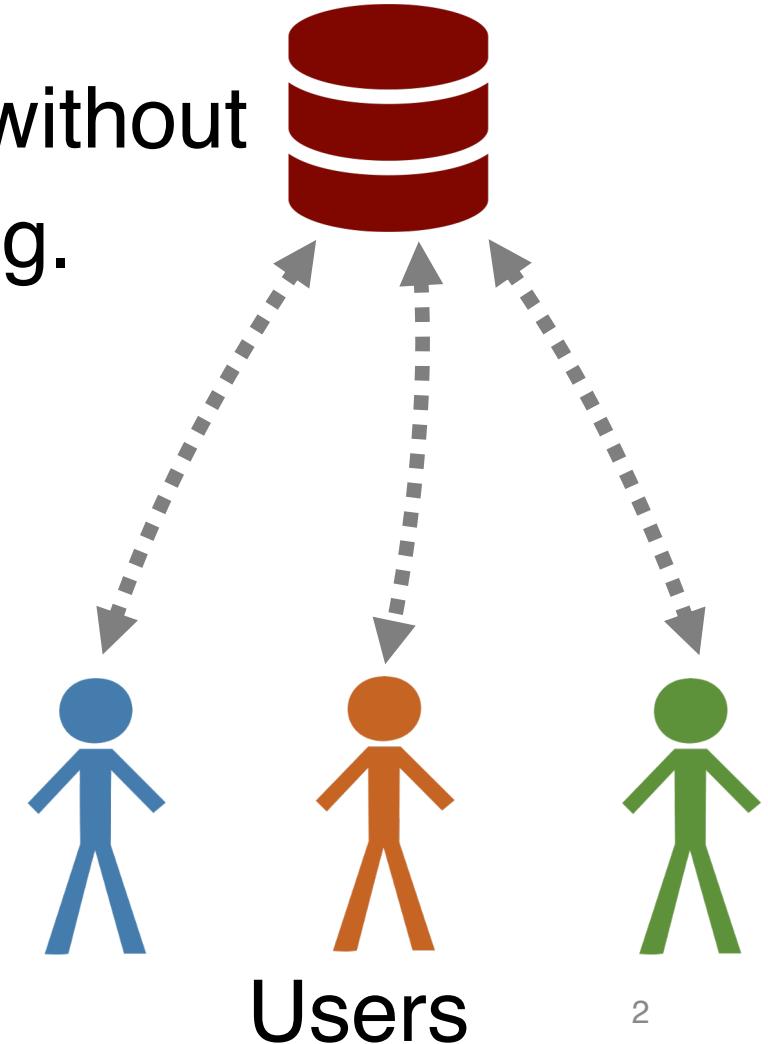
## Goal

Users can read records from a public DB without DB learning which records they are reading.

## Applications

medical encyclopedia, stocks  
privacy-preserving systems

Public Database



# Part I: introduction to PIR

The model

Basic schemes

Extensions

# Part II: reducing computation in PIR

# PIR requirements

## Correctness

Client learns its bit of interest (with overwhelming prob.)

## Security

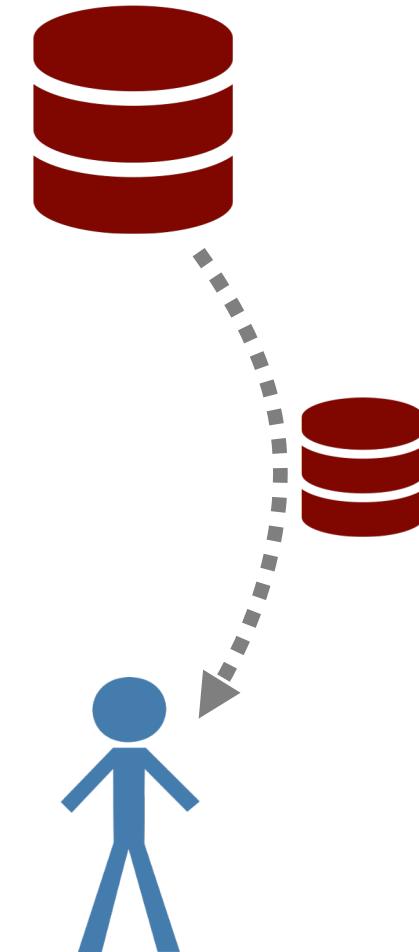
(Malicious) server “learns nothing” about client’s desired bit

# Naïve solution

Client downloads entire DB and queries it locally

Perfect correctness and privacy

**Downside:** large communication cost



# PIR requirements

## Correctness

Client learns its bit of interest (with overwhelming prob.)

## Security

(Malicious) server “learns nothing” about client’s desired bit

## Minimize communication

# Main approaches to PIR

## **Multi-server PIR** [CGKS95]

Replicate DB on non-colluding servers

## **Single-server PIR** [KO97]

Requires cryptographic assumptions

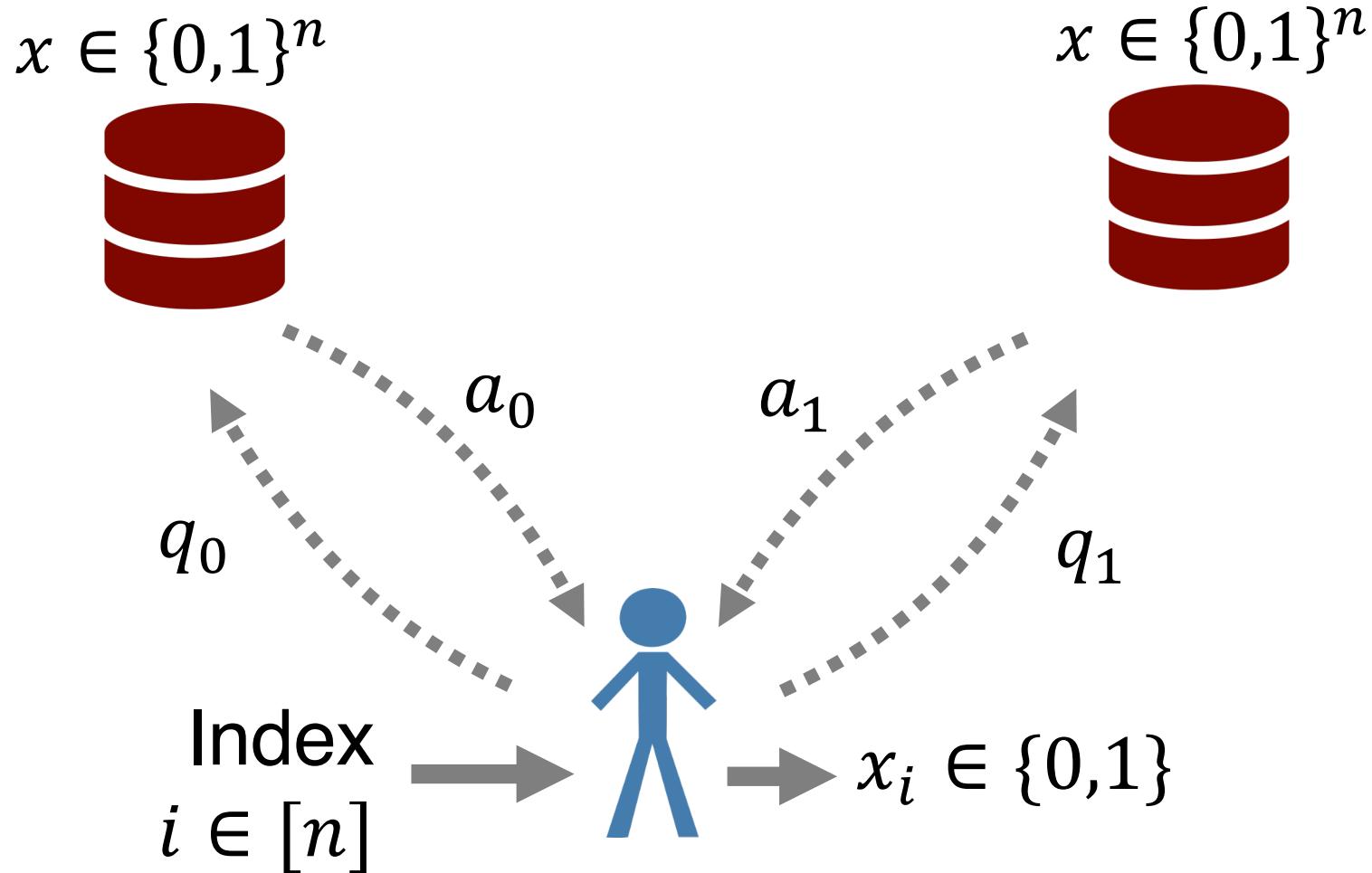
# 2-server PIR

Three efficient algorithms:

$$\text{Query}(1^n, i) \rightarrow (q_0, q_1)$$

$$\text{Answer}(x, q) \rightarrow a$$

$$\text{Reconstruct}(a_0, a_1) \rightarrow x_i$$



# PIR requirements

## Correctness

Client learns its bit of interest:  $\forall n \in \mathbb{N}, x \in \{0,1\}^n, i \in [n]$ :

$$\Pr \left[ \begin{array}{c} q_0, q_1 \leftarrow \text{Query}(1^n, i) \\ \text{Reconstruct}(a_0, a_1) = x_i : a_0 \leftarrow \text{Answer}(x, q_0) \\ a_1 \leftarrow \text{Answer}(x, q_1) \end{array} \right] = 1$$

# PIR requirements

## Security

(Malicious) server “learns nothing” about client’s desired bit

$\forall n \in \mathbb{N}, i, j \in [n], s \in \{0,1\},$

$$\{q_s : q_0, q_1 \leftarrow \text{Query}(1^n, i)\} \approx_c \{q_s : q_0, q_1 \leftarrow \text{Query}(1^n, j)\}$$

$$\left\{ \begin{array}{l} \text{View of server } s \text{ when} \\ \text{client reads bit } i \end{array} \right\} \approx_c \left\{ \begin{array}{l} \text{View of server } s \text{ when} \\ \text{client reads bit } j \end{array} \right\}$$

Non-colluding servers

Can be  $\equiv$  for perfect security

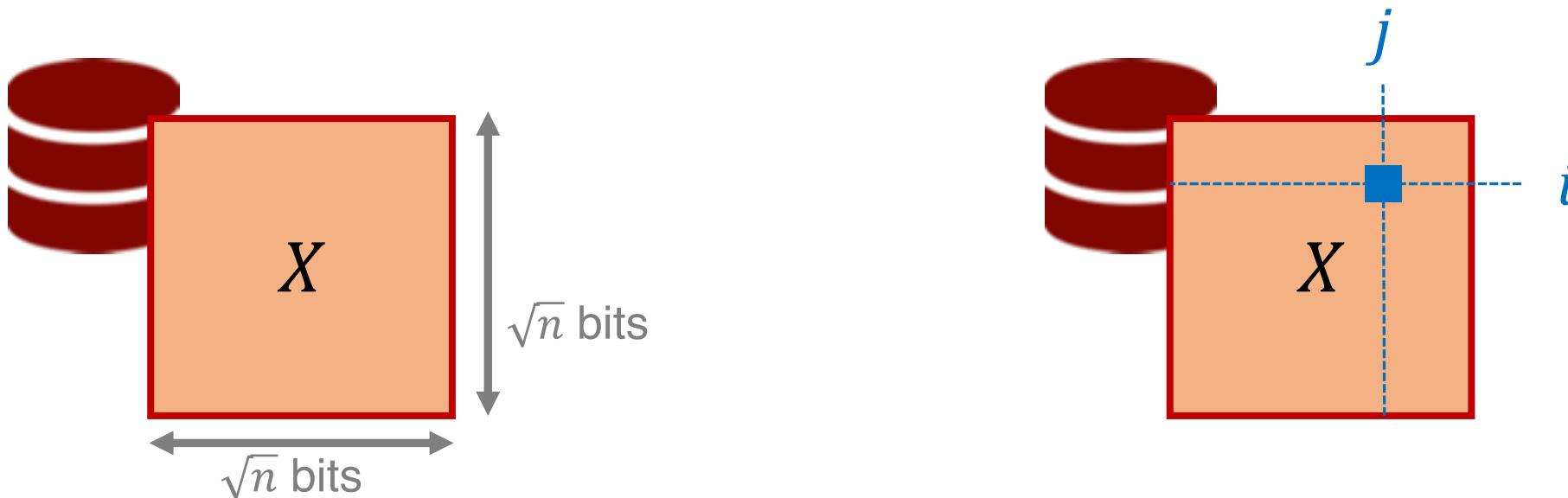
# Part I: introduction to PIR

The model

Basic schemes

Extensions

# A two-server PIR scheme [CGKS95]

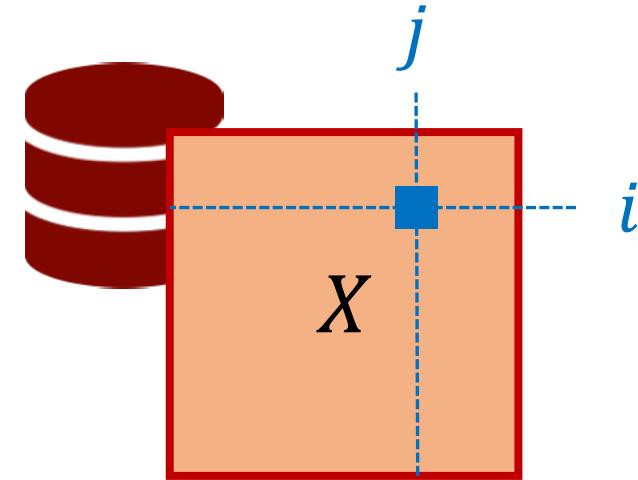
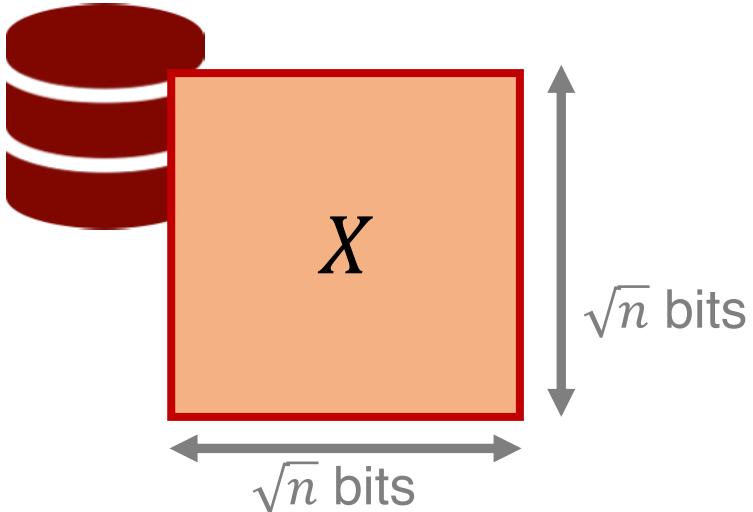


View DB  $x \in \{0,1\}^n$  as matrix  $X \in \mathbb{Z}_2^{\sqrt{n} \times \sqrt{n}}$



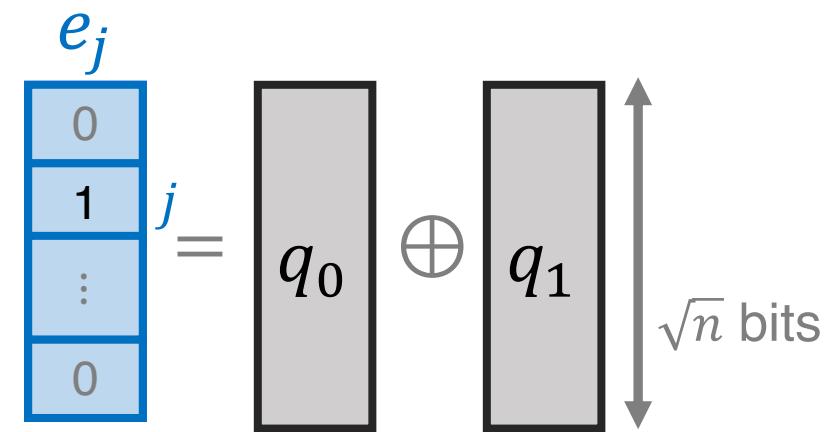
Client wants to read  $(i, j) \in [\sqrt{n}] \times [\sqrt{n}]$

# A two-server PIR scheme

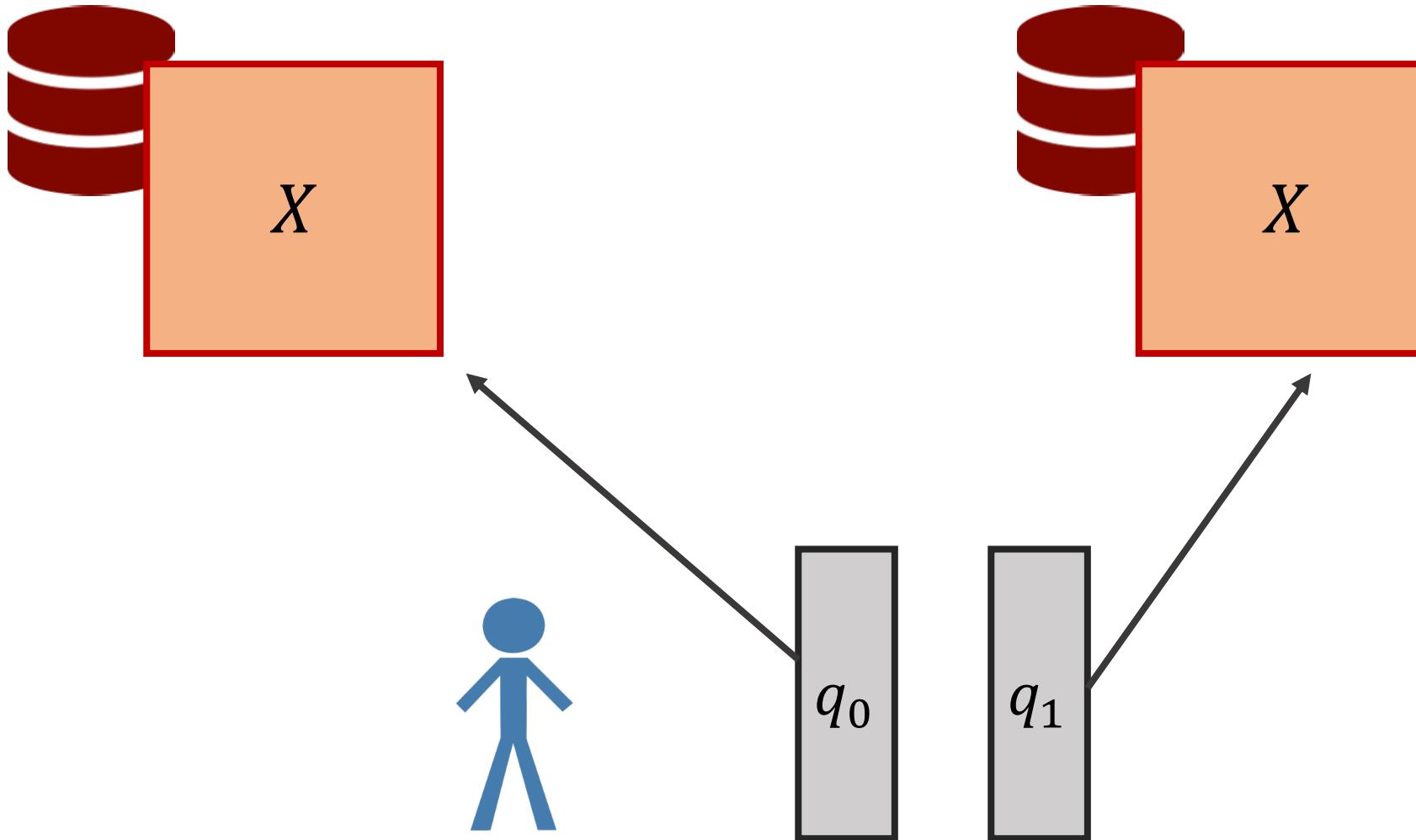


Query( $1^n, (i, j)$ ):

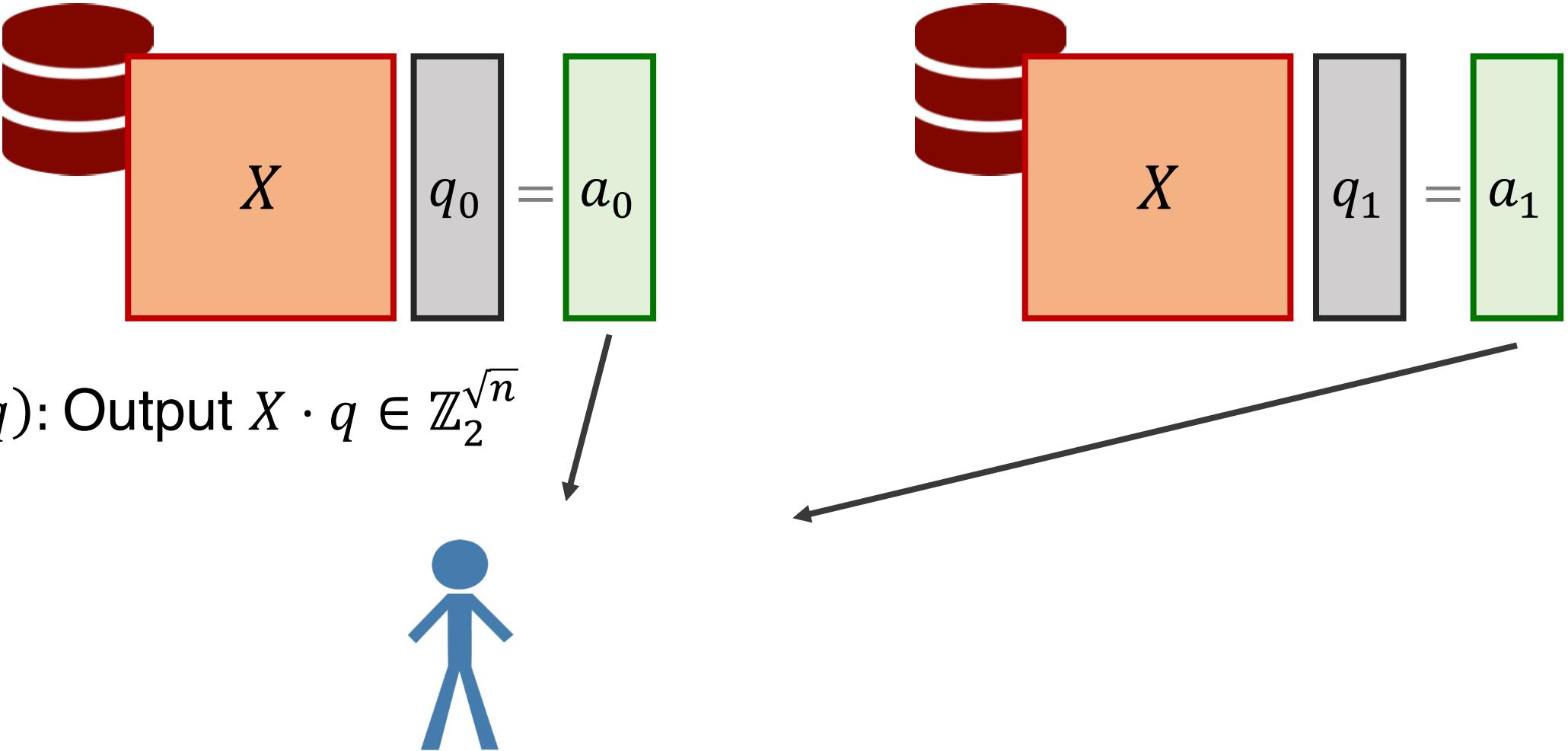
Sample  $q_0, q_1 \leftarrow_R \mathbb{Z}_2^{\sqrt{n}}$  s.t.  $q_0 + q_1 = e_j$



# A two-server PIR scheme



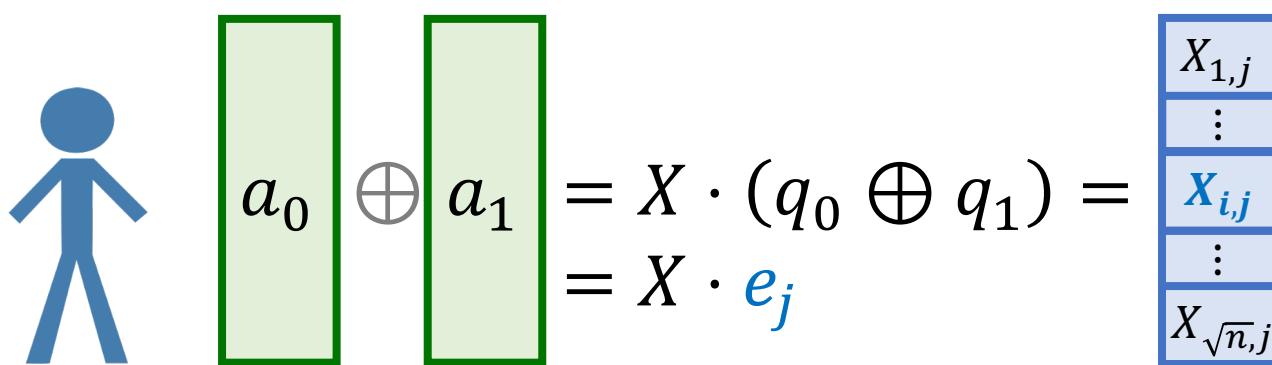
# A two-server PIR scheme



# A two-server PIR scheme



Reconstruct( $a_0, a_1$ ): Output  $(a_1 + a_1)_i$



# A two-server PIR scheme [CGKS95]

## Correctness:

$$(a_0 \oplus a_1)_j = (Xq_0 \oplus Xq_1)_i = (X(q_0 \oplus q_1))_i = (Xe_j)_i = x_{ij}$$

## Security:

Each of  $q_0, q_1$  on its own is a random vector,  
independent of client's input  $(i, j)$

## Efficiency: $|q_0| + |q_1| + |a_0| + |a_1| = 4\sqrt{n}$

Can reduce to  $O(\log n)$  using Distributed Point Functions.  
Security then computational.

# The state of PIR

## Multi-server PIR

[CGKS95, Amb97, CG97, BI01, BIKR02, Yek08, Efr12, ...]):

- Information-theoretic security:  $n^{o(1)}$  communication [DG16]
- Computational security:  $O(\log n)$  communication [GI14, BGI15]

## Single-server PIR [KO97]

# Single-server PIR [KO97]

Idea: Encrypt query instead of secret sharing it

Uses: Linearly Homomorphic Encryption

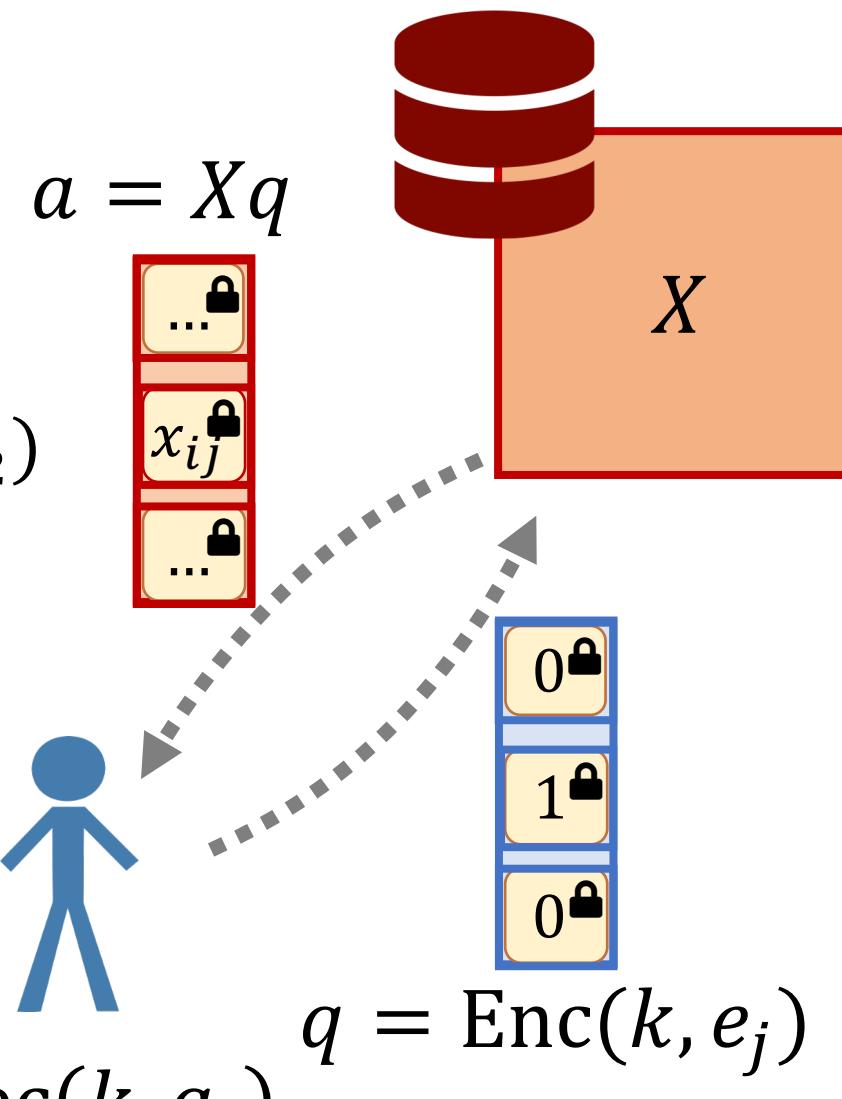
$$\text{Enc}(k, m_1) + \text{Enc}(k, m_2) = \text{Enc}(k, m_1 + m_2)$$

Can build from DDH, QR, LWE,...

**Security:** follows from security of LHE

**Communication:**  $2\sqrt{n}$  ciphertexts

Output  $\text{Dec}(k, a_i)$



# The state of PIR

## Multi-server PIR [CGKS95]

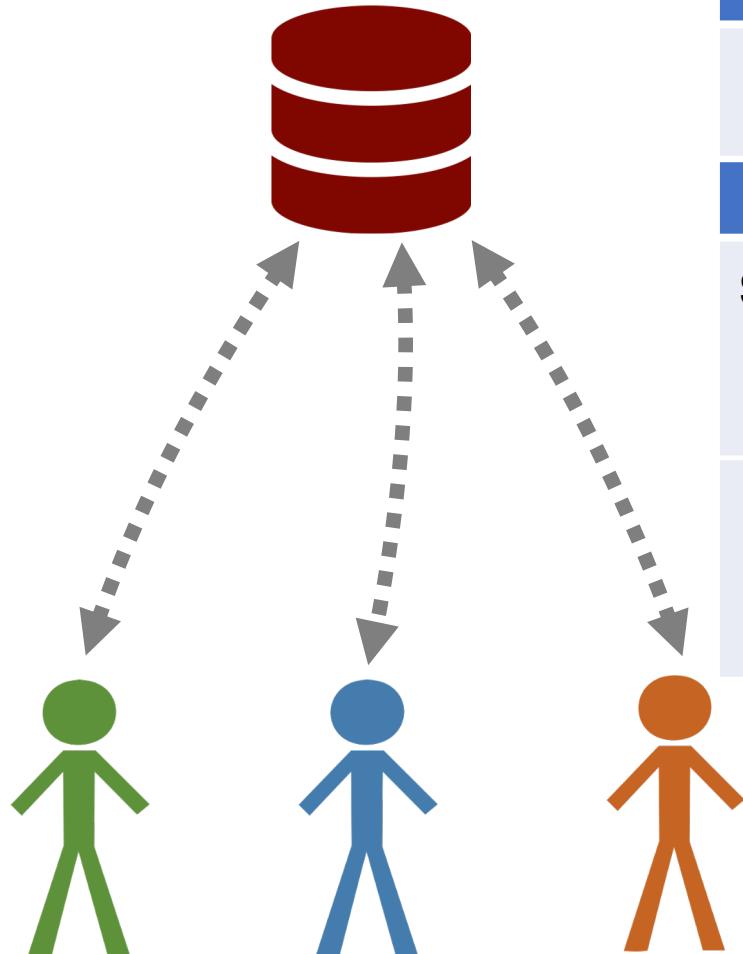
- Replicate DB on non-colluding servers
- State of the art (following [Amb97, CG97, BI01, BIKR02, Yek08, Efr12, ...]):
  - Information-theoretic security:  $n^{o(1)}$  communication [DG16]
  - Computational security:  $O(\log n)$  communication [GI14, BGI15]

## Single-server PIR [KO97]

- Requires cryptographic assumptions
- State of the art:
  - $\text{polylog}(n)$  communication [CMS99, Lip05, ...]

# PIR (today)

private access to public data

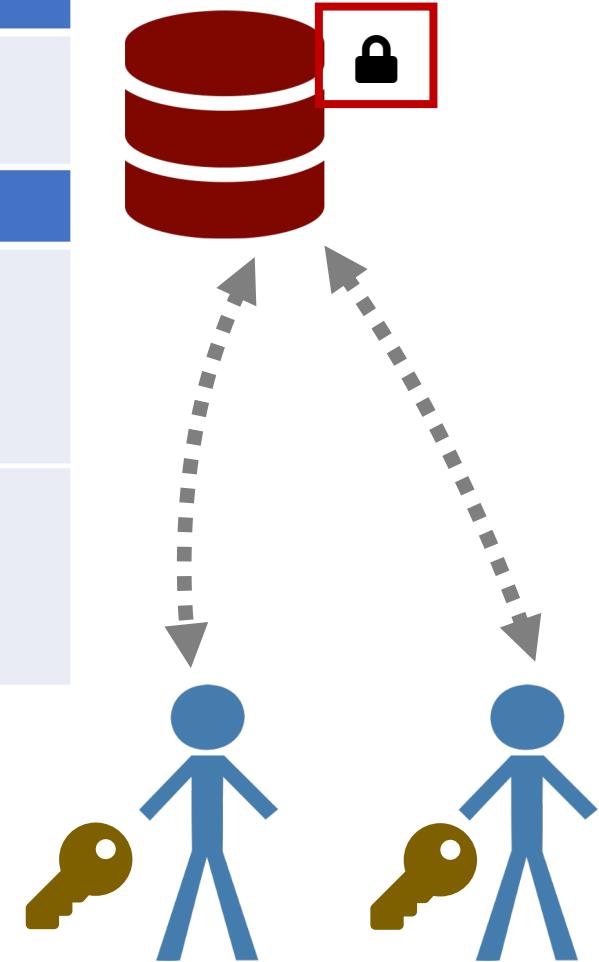


VS

# ORAM (Sunday)

private access to private data

Similarity	
both hide the data access pattern	
Differences	
<b>Supports many clients</b> No need in trust or coordination	<b>Single client</b> Operations require private state
<b>Read only</b>	<b>Read &amp; write</b>



# Part I: introduction to PIR

The model

Basic schemes

Extensions

# Making PIR more realistic...

Most interesting databases are not n-bit arrays

- Extension 1: longer DB rows
- Extension 2: key-value lookups

Servers also have privacy requirements

- Extension 3: privacy for the server

# Extension 1: longer DB records

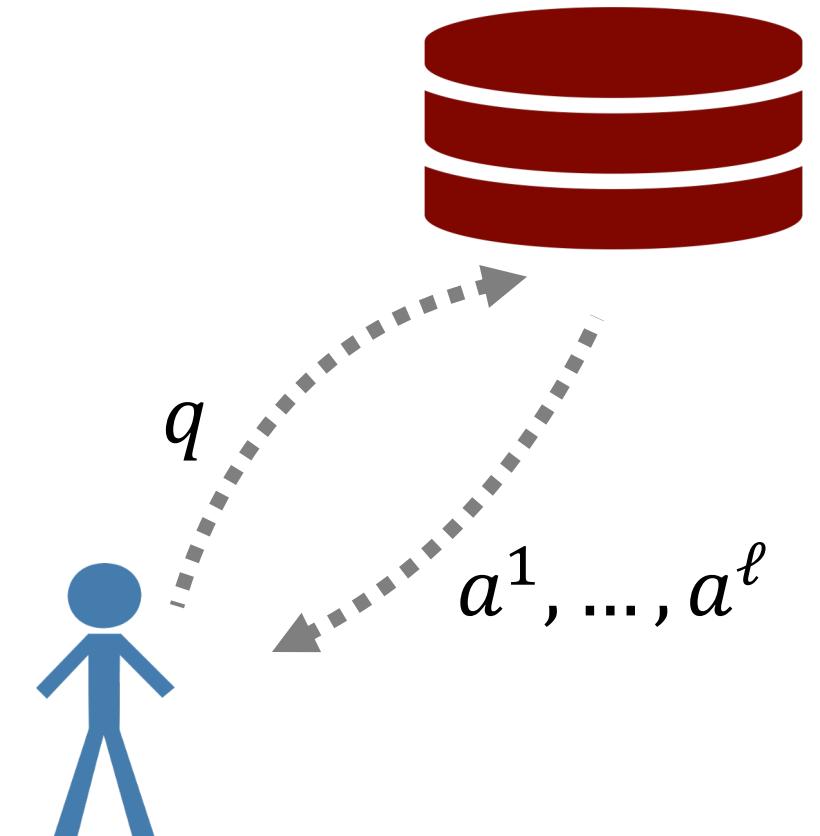
**Want:** PIR for  $\ell$ -bit records:  $X \in \{0,1\}^{n \times \ell}$

On the same  
server(s)

$$X^1, \dots, X^\ell \in \{0,1\}^n$$

**Solution:** query each bit of the record  
using PIR for 1-bit records

**Communication:**  $U(n) + \ell \cdot D(n)$



Reconstruct bit by bit

# Extension 2: PIR by keywords [CGN98]

Standard PIR: DB is array  $x \in \{0,1\}^{n \times \ell}$

Many applications want a **set**  $S = \{k_1, \dots, k_n\}$   
or a **dictionary**  $D = \{(k_1, v_1), \dots, (k_n, v_n)\}$

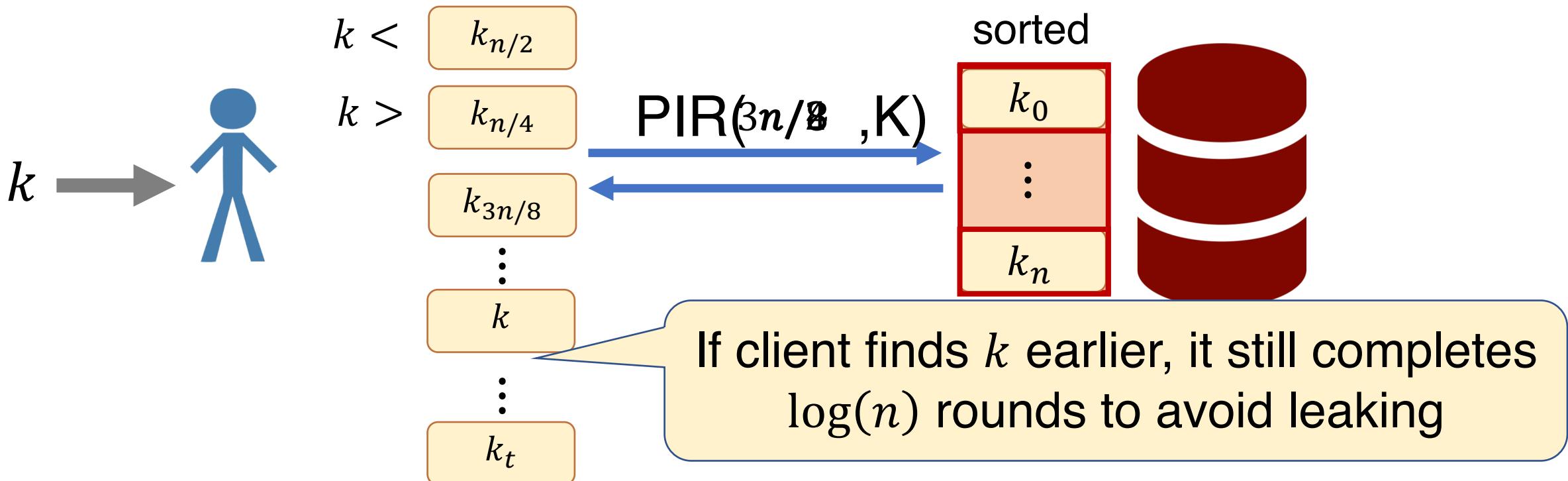
## Idea:

- Build a standard data structure (e.g., binary search tree, hash table,...)
- Wrap each probe into data-structure with PIR

# Extensions : PIR by keywords [CGN98]

Simple example: store set  $S$  as a sorted array

Client runs binary search for  $k$  using PIR



Efficiency:  $\log(n)$  PIR queries (can reduce to 2 queries by using perfect hashing)

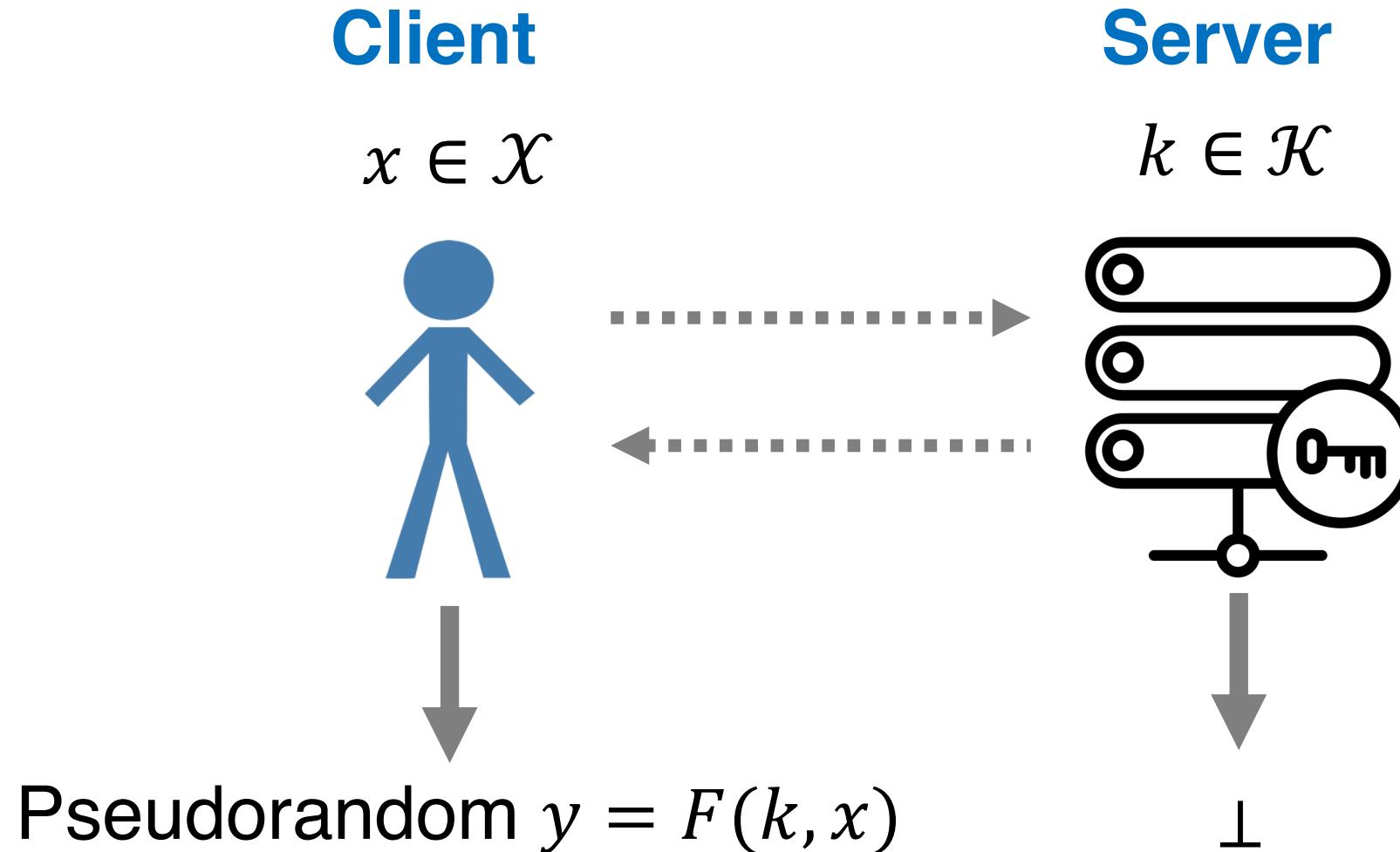
# Extension 3: symmetric PIR [GIKM97]

Symmetric PIR – privacy also for the server

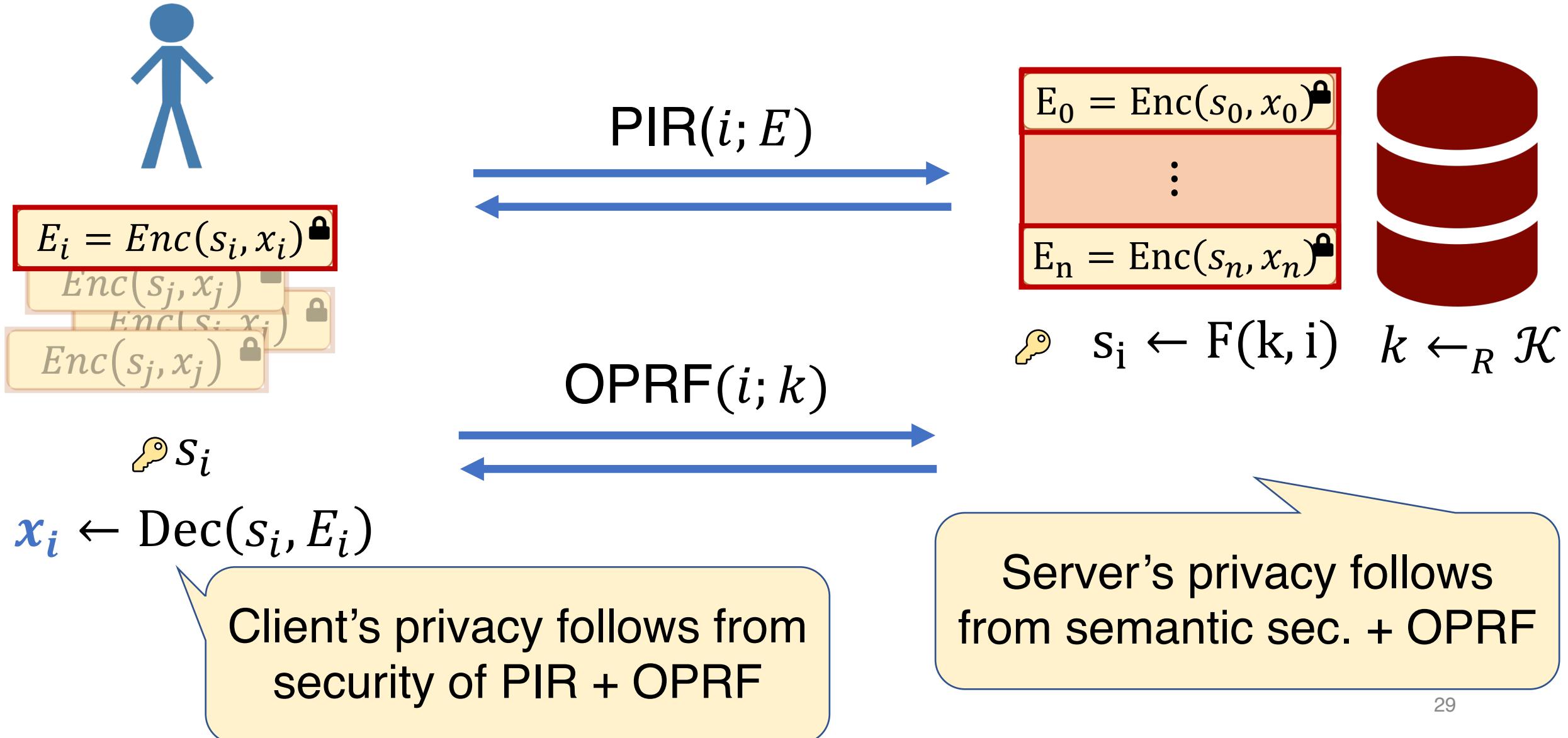
**Privacy requirement (informal):** after each execution of the protocol, the client should learn at most one record in the database

**Tool:** Oblivious Pseudorandom Function (OPRF)

# Tool: Oblivious PRF [NR97,FIPR05]



# Symmetric PIR from OPRF [FIPR05]



# Modern PIR requires lots of computation

Server linearly scans the entire DB to respond to a query  
⇒ a barrier to deployment

Server **must do  $\Omega(n)$  work** to respond to a query [BIM04]

- Intuition: If server doesn't touch bit  $i$ , client isn't reading bit  $i$
- Holds even if you have **many non-colluding servers**
- Holds irrespective of **cryptographic assumptions**

# Part II: reducing computation in PIR

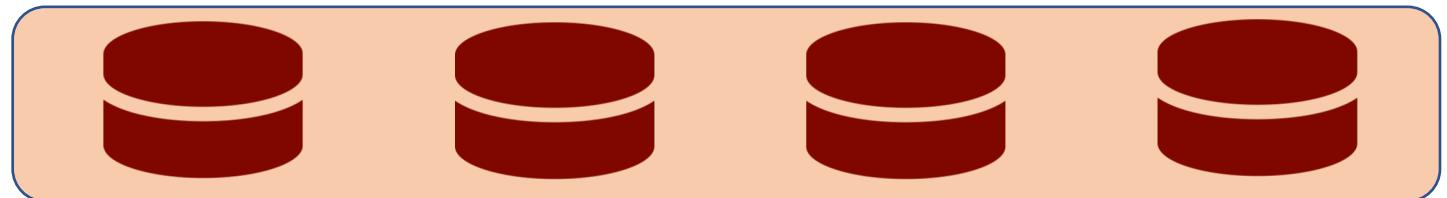
Two classic approaches

Recent approach: offline/online PIR

# Batch PIR [IKOS04, IKOS06, LG15, Hen16, ACLS18]



$i_1, \dots, i_Q$



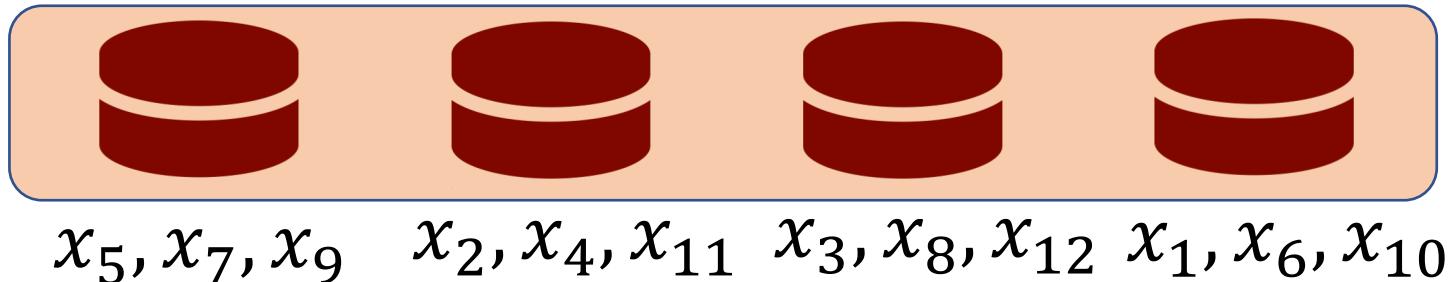
Server & client agree on a random partition of DB into  $Q$  buckets

For  $i_1, \dots, i_Q \in [n]$ , w.h.p. no bucket contains more than  $\lambda \log n$  indices

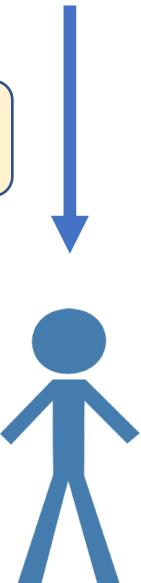
Client queries each bucket  $\lambda \log n$  times

# Batch PIR [IKOS04, IKOS06, LG15, Hen16, ACLS18]

$$x \in \{0,1\}^{12}$$
$$Q = 4$$

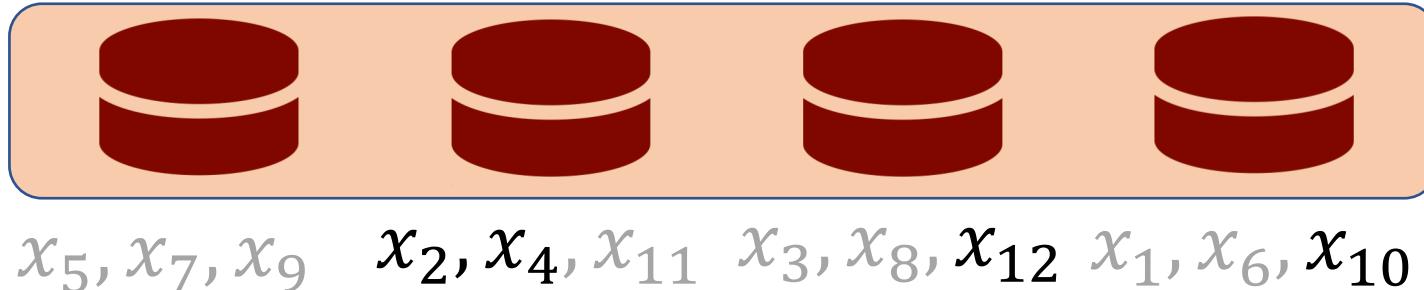


description of  $h$



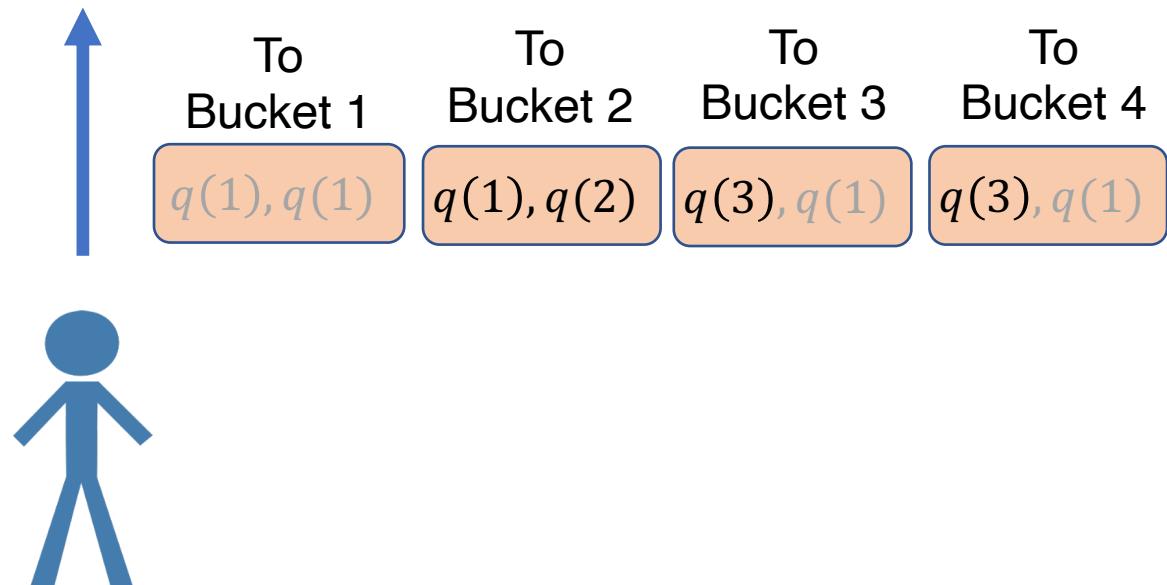
# Batch PIR [IKOS04, IKOS06, LG15, Hen16, ACLS18]

$$x \in \{0,1\}^{12}$$
$$Q = 4$$



Only works if client makes a non-adaptive batch of queries

10,2,4,12



# Batch PIR – analysis

## Correctness:

W.h.p each bucket contains  $\leq \lambda \log n$  input indices

Correctness follows from correctness of base PIR scheme

**Security:** follows from security of base PIR scheme (server sees  $\lambda \log n$  queries for each bucket, regardless of client's inputs)

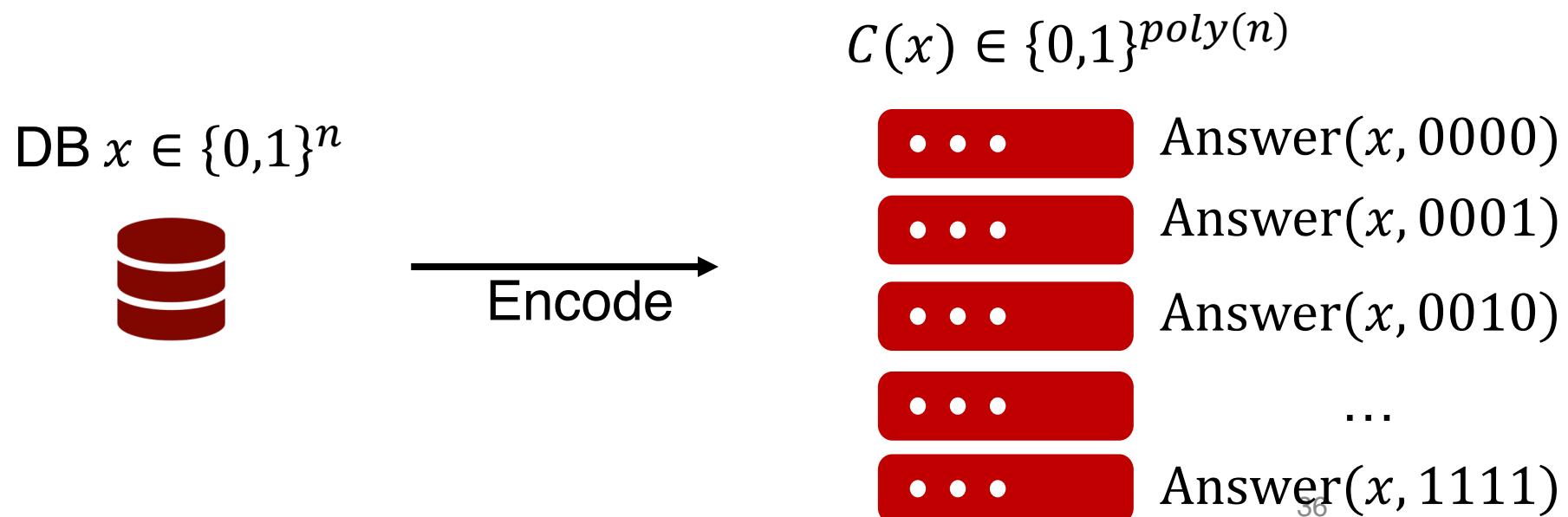
## Efficiency:

$$\text{Server time: } Q \cdot \lambda \log n \cdot T\left(\frac{n}{Q}\right) = Q \cdot \lambda \log n \cdot \frac{n}{Q} = n \cdot \lambda \log n$$

#buckets      time/query      indep. of Q  
                  queries/bucket

# PIR with preprocessing [BIM04]

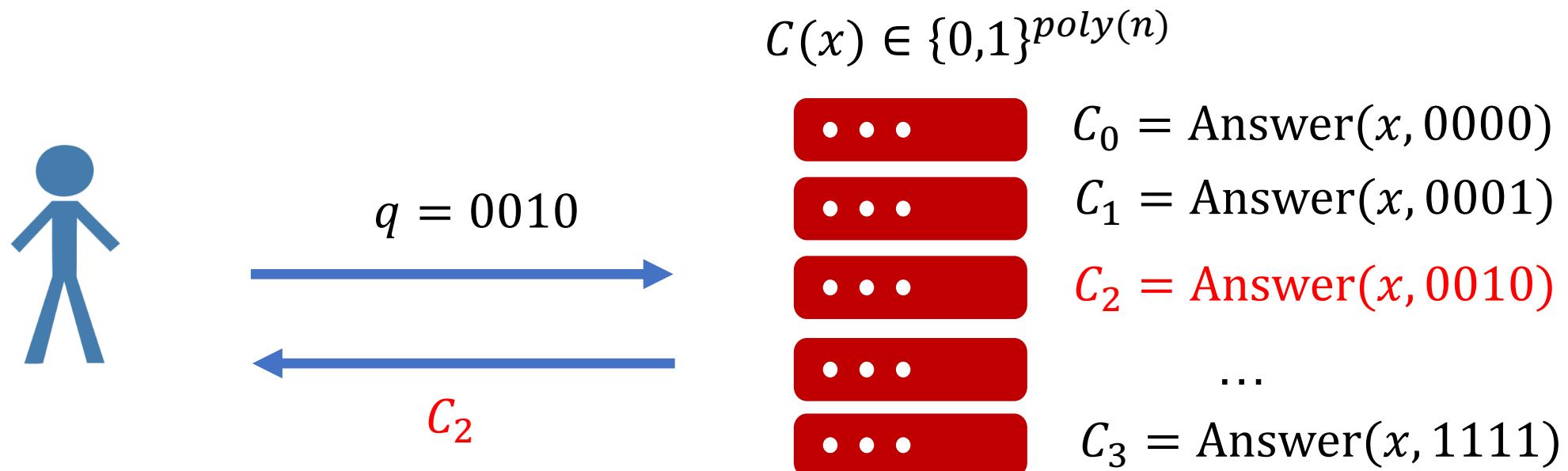
Main idea: servers precompute the responses to all possible queries of a standard PIR scheme.



# PIR with preprocessing [BIM04]

Main idea: servers precompute the responses to all possible queries of a standard PIR scheme.

In the online phase: the server looks up the response from a table.



# PIR with preprocessing [BIM04]

Base scheme has communication  $U(n)$  and  $D(n)$

⇒ new scheme with space  $S = D(n) \cdot 2^{U(n)}$  and time  $T = U(n) + D(n)$

## Examples (2-server PIR):

- $\sqrt{n}$  scheme ⇒ PIR with preproc.  $S = 2^{\sqrt{n}}$  and  $T = \sqrt{n}$
- [BI01]  $U(n) = O(\log n)$  and  $D(n) = n^{1/2+\epsilon}$   
⇒ PIR with preproc.  $S = \text{poly}(n)$  and  $T = n^{1/2+\epsilon}$

**Advantage:** sublinear server time (e.g.,  $n^{0.6}$ )

**Disadvantage:** superlinear server storage (e.g.,  $n^3$ )

Also: DEPIR [BIPW17, CHR17], PANDA [HOWW18] for the 1-server setting

# Part II: reducing computation in PIR

Two classic approaches

Recent approach: offline/online PIR

Based on joint works with Henry Corrigan-Gibbs and  
Alexandra Henzinger [CK20,KC21,CHK22]

# The model

## Step 1: Offline phase

 $x \in \{0,1\}^n$  $O(n)$  time $x \in \{0,1\}^n$ 

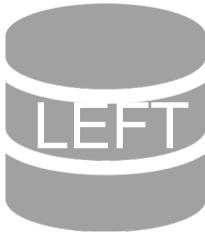
Hint

 $\approx \sqrt{n}$  bits $\circ(n)$  bits

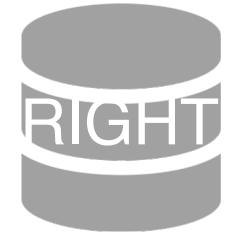
The offline phase runs in linear time.  
But only once per client.

# The model

## Step 1: Offline phase



$x \in \{0,1\}^n$



$x \in \{0,1\}^n$

Client stores hint

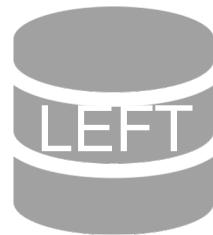


Hint

# The model

[DIO01, BIM04, BLW17, PPY18]

## Step 2: Online phase – reading $x_i$



$x \in \{0,1\}^n$



$x \in \{0,1\}^n$

Sublinear online time

$o(n)$  time

Client can repeat online phase to read multiple items **adaptively**

Index  
 $i \in [n]$



$x_i \in \{0,1\}$

$o(n)$  bits

Hint

# Offline/online PIR [CK20, KC21, CHK22]

up to poly( $\lambda, \log n$ ) factors for length- $n$  DB and sec. parameter  $\lambda$

## Two-server scheme

- $\sqrt{n}$  communication and online time (from any PRG)
- Client uses  $\sqrt{n}$  time and storage (using PRPs)

# Offline/online PIR [CK20, KC21, CHK22]

up to  $\text{poly}(\lambda, \log n)$  factors for length- $n$  DB and sec. parameter  $\lambda$

## Two-server scheme

- $\sqrt{n}$  communication and online time (from any PRG)
- Client uses  $\sqrt{n}$  time and storage (using PRPs)

## Single-server scheme

- $n^{3/4}$  online time (from DDH, DCR,...)
- $\sqrt{n}$  from FHE
- No public-key operations in the online phase

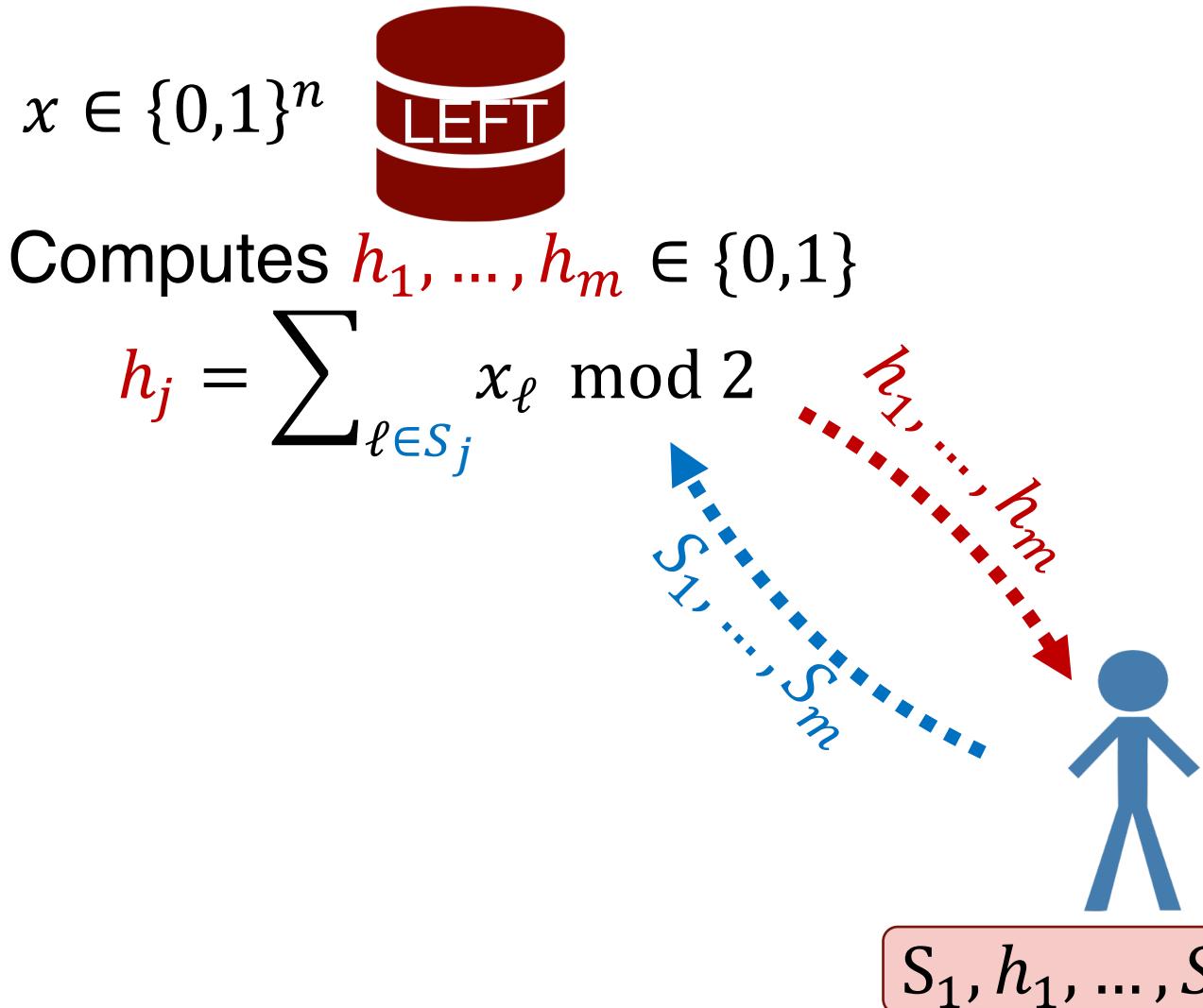
Our  $\sqrt{n}$  schemes achieve optimal comm–online time tradeoff

## Lower bound

- For offline/online schemes that store DB in its original form
- Communication  $C$  and online time  $T$  must be  $C \cdot T \geq n$

# Our scheme

## Step 1: Offline phase



$x \in \{0,1\}^n$

Random subsets  $S_1, \dots, S_m \subset [n]$   
each of size  $|S_j| = \sqrt{n}$

# Our scheme

## Step 2: Online phase – reading $x_i$

If  $i \notin S_1 \cup \dots \cup S_m$ , output “fail”

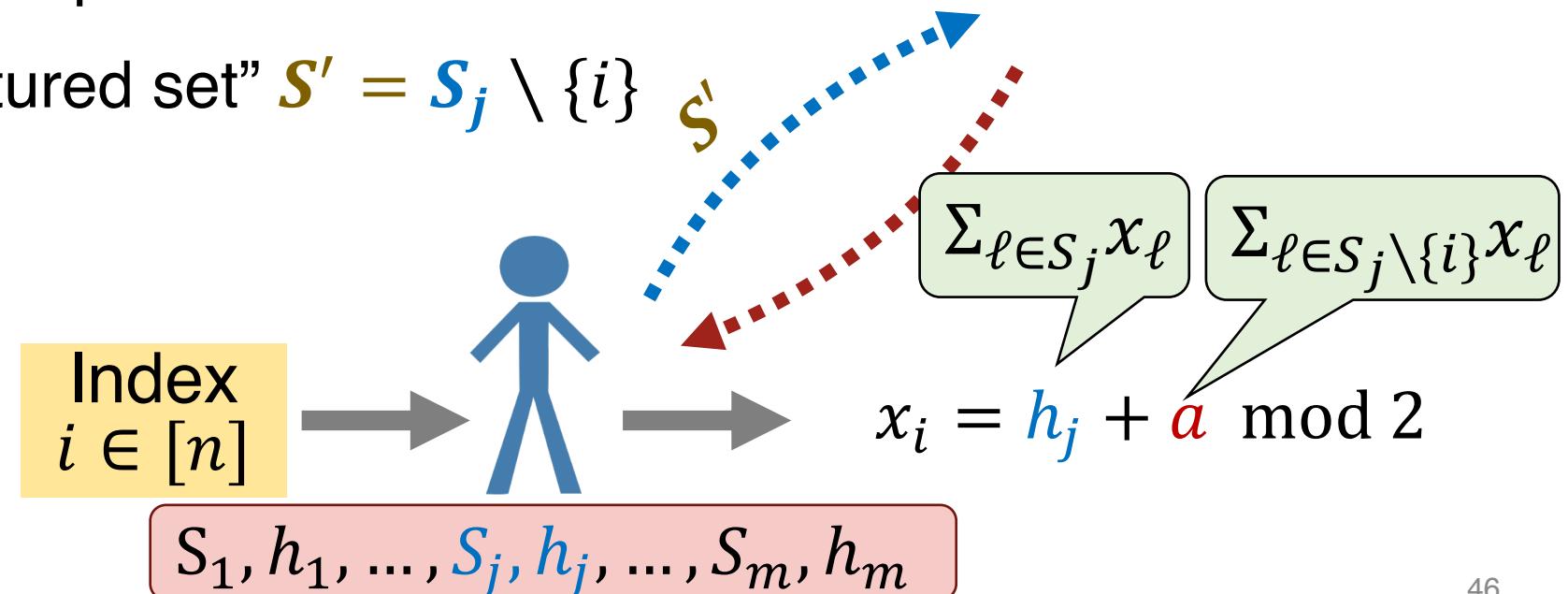
Else,  $i \in S_j$ ,

- With prob  $\frac{\sqrt{n}-1}{n}$ , send a random set  $S'$  containing  $i$  and output “fail”
- Else, send “punctured set”  $S' = S_j \setminus \{i\}$



$$x \in \{0,1\}^n$$

$$a = \sum_{\ell \in S'} x_\ell \bmod 2$$



# Our scheme Correctness

If  $i \notin S_1 \cup \dots \cup S_m$ , output “fail”

Else,  $i \in S_j$ ,

- With prob  $\frac{\sqrt{n}-1}{n}$ , send a random set  $S'$  containing  $i$  and output “fail”
- Else, send  $S' = S_j \setminus \{i\}$

Index  
 $i \in [n]$

$S_1, h_1, \dots, S_j, h_j, \dots, S_m, h_m$

Choose

$$m \approx \lambda \sqrt{n}$$

Then:

$$\Pr[\text{Fail}_1] \leq \text{negl}(\lambda)$$

( $\lambda n$  balls into  $n$  bins)

$$a = x_\ell \bmod z$$

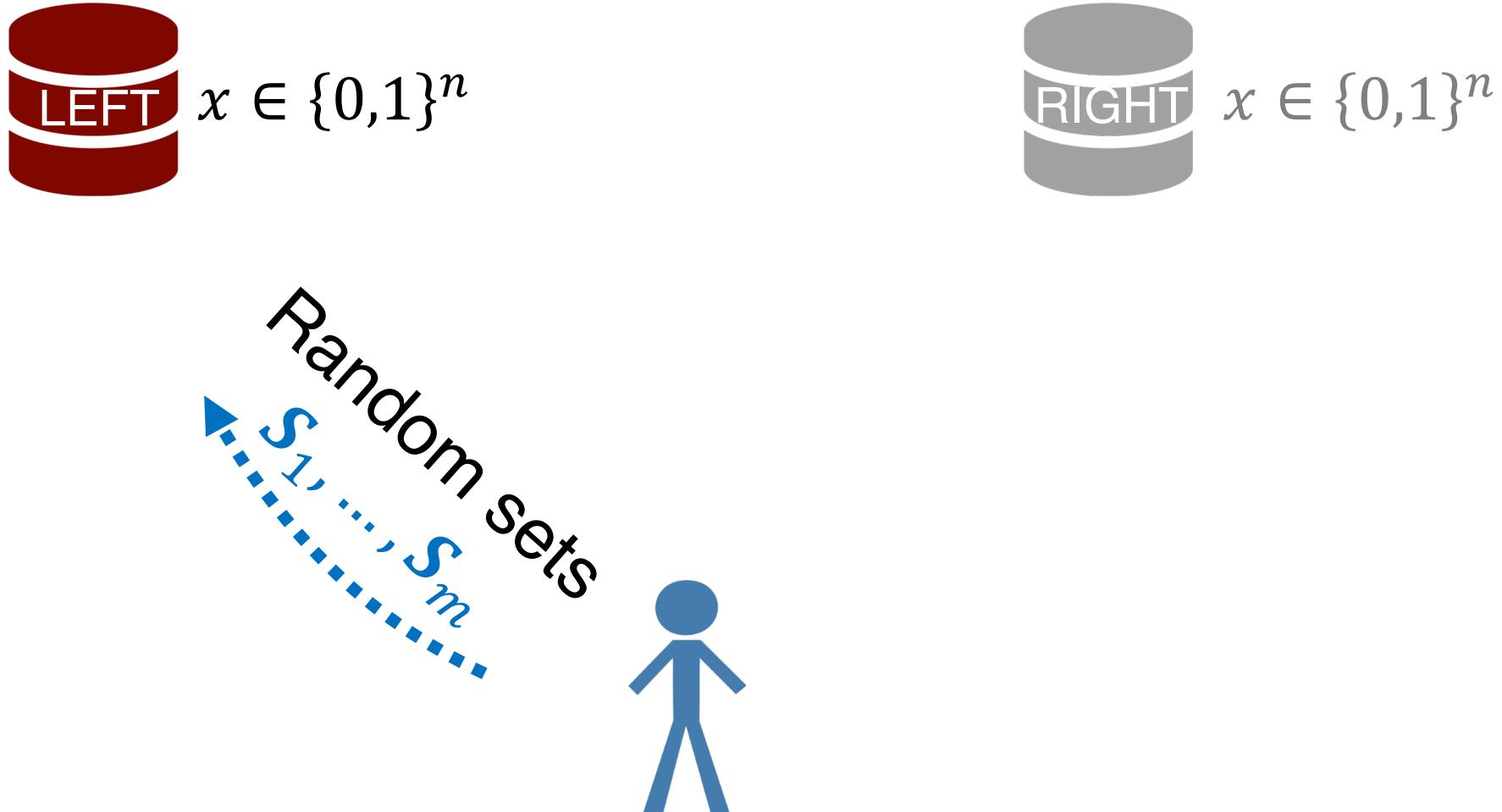
$$\Pr[\text{Fail}_2] \leq 1/\sqrt{n}$$

Repeat  $\lambda$  times to drive down failure prob.

can avoid this using more involved scheme

# Our scheme

## Security – left server



# Our scheme Security – right server

$$x \in \{0,1\}^n$$

- With prob  $\frac{\sqrt{n}-1}{n}$ , send a random set  $S'$  containing  $i$ , output “fail”
- Else, send set  $S' = S_j \setminus \{i\}$

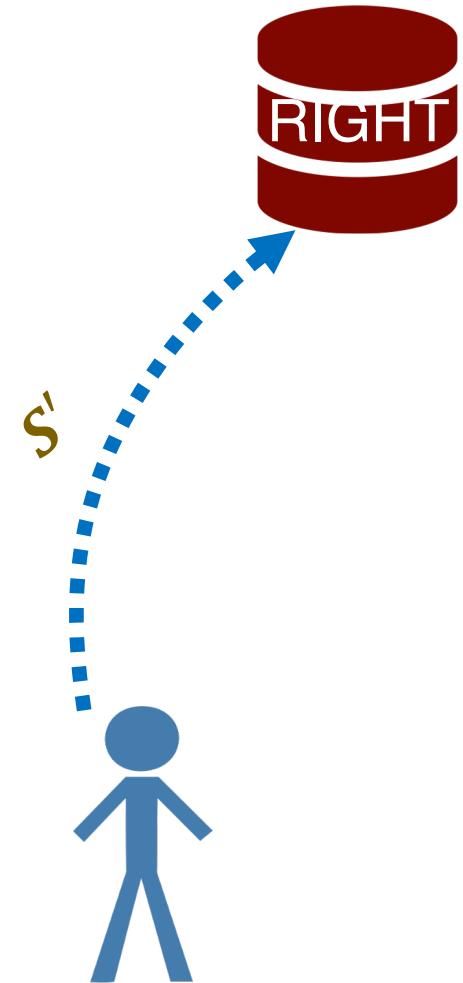
uniformly random size- $(\sqrt{n} - 1)$  subset of  $[n]$

w.p.  $p = \frac{\sqrt{n}-1}{n}$

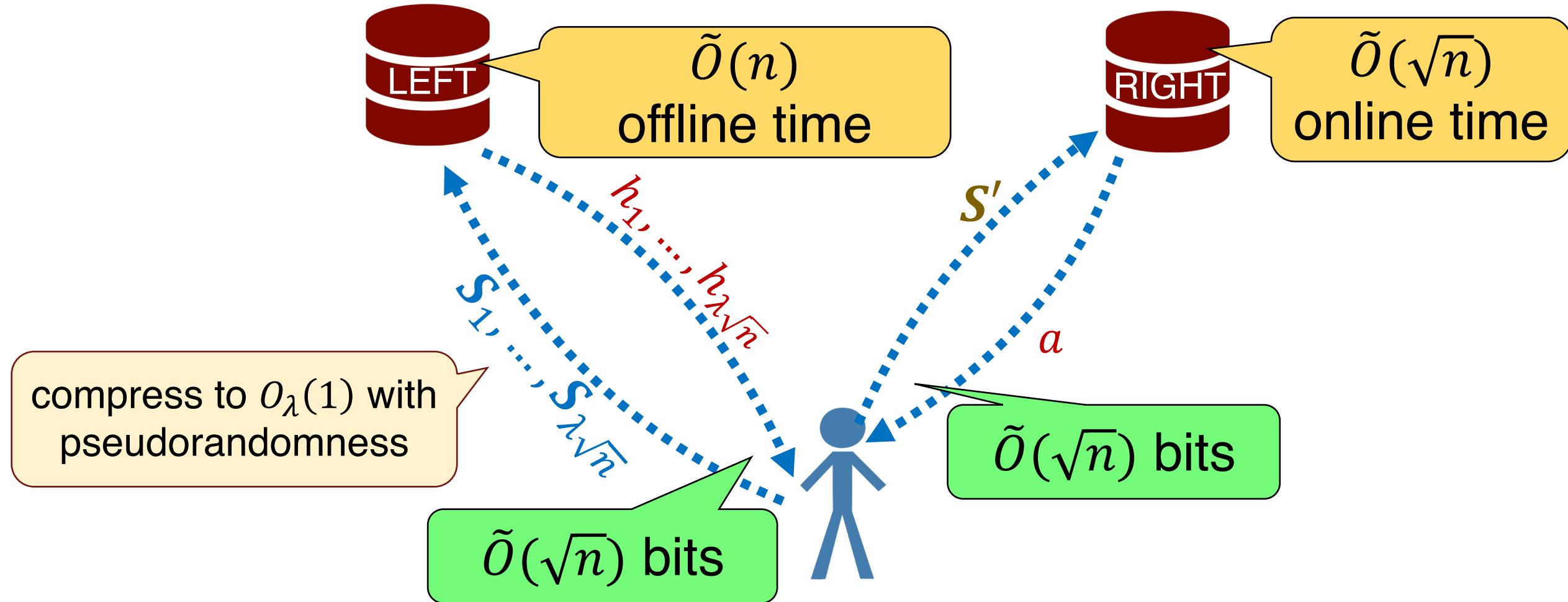
random set containing  $i$

w.p.  $1 - p$

random set without  $i$



# Our scheme Efficiency



## Two-server scheme summary

- $\sqrt{n}$  communication, online time, amortized total time per-query
- Client uses  $\sqrt{n}$  time and storage (using PRPs)

## Extensions (see paper)

- Trade-off communication for online time
- Statistical-security variant:  $n^{2/3}$  communication and client time
- Reducing online communication to **log n**
  - Using short description of ‘Puncturable sets’
  - Client storage and time increase to  $n^{5/6}$

[SACM20] show how to avoid this blowup using a new construction from privately puncturable PRFs (from LWE)

# Offline/online PIR

## Two-server scheme

- $\sqrt{n}$  communication and online time (from any PRG)
- Can reuse a single offline interaction for many online queries

## Single-server scheme

- $n^{3/4}$  communication and online time (from DDH, DCR,...)
  - $\sqrt{n}$  from FHE
- No public-key operations in the online phase

## Lower bound

- For offline/online schemes that store DB in its original form
- Communication  $C$  and online time  $T$  must be  $C \cdot T \geq n$

# From 2 servers to 1

Run both offline and online phases with the same server

Single server homomorphically evaluates offline query

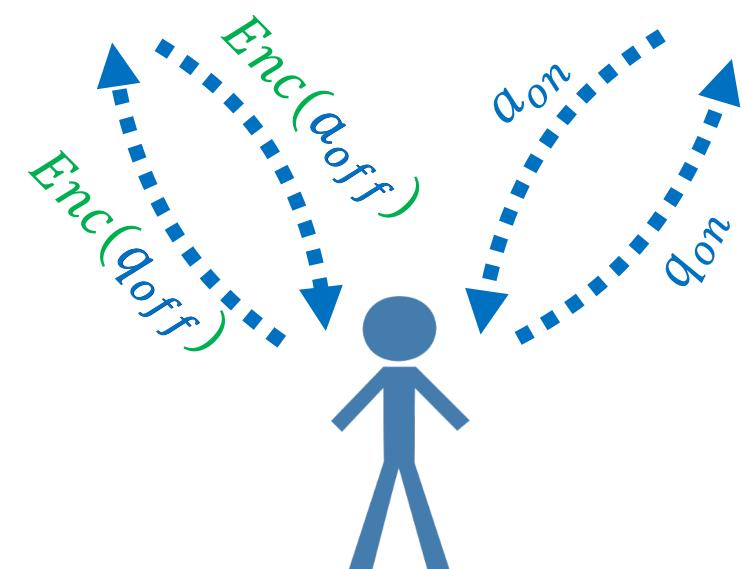
- Option 1: Fully HE
  - $\sqrt{n}$  communication and online time
- Option 2: Additively HE
  - $n^{3/4}$  communication and online time

Security only holds if server does not see both offline and online queries



Expensive offline phase

Fast online phase



# From 2 servers to 1 – linearize the offline phase

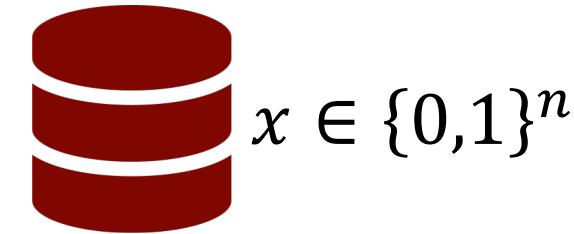
Make server's offline response **linear** in its query:

- Client sends characteristic vectors of sets
- Server computes matrix-vector product

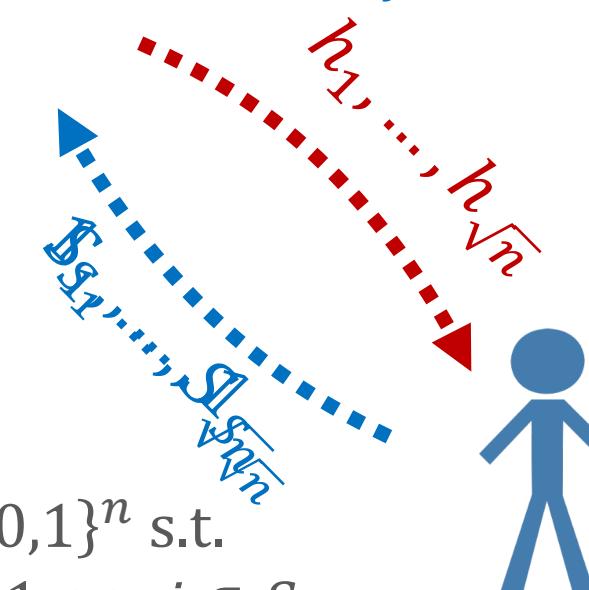
$$\begin{bmatrix} \mathbb{1}_{S_1} \\ \vdots \\ \mathbb{1}_{S_{\sqrt{n}}} \end{bmatrix} \cdot \begin{matrix} x \end{matrix} = \begin{bmatrix} h_1 \\ \vdots \\ h_{\sqrt{n}} \end{bmatrix}$$

Recall: in offline phase of 2-server scheme:

- Server sends  $\sqrt{n}$  sets of size  $\sqrt{n}$
- Server responds with parity of each set



$$h_j = \sum_{i \in S_j} x_i \bmod 2$$



$$\begin{aligned} \mathbb{1}_S &\in \{0,1\}^n \text{ s.t.} \\ (\mathbb{1}_S)_i &= 1 \Leftrightarrow i \in S \end{aligned}$$

# From 2 servers to 1 – encrypt the offline phase

Make server's offline response **linear** in its query

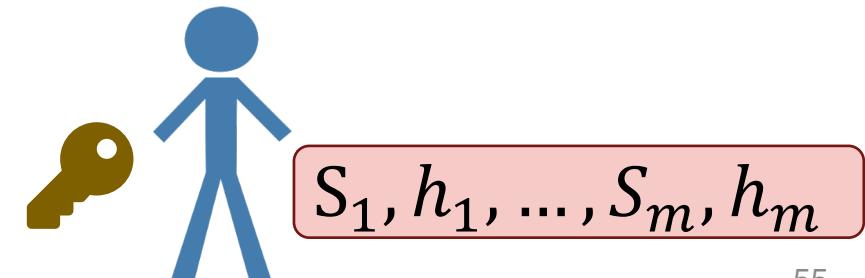
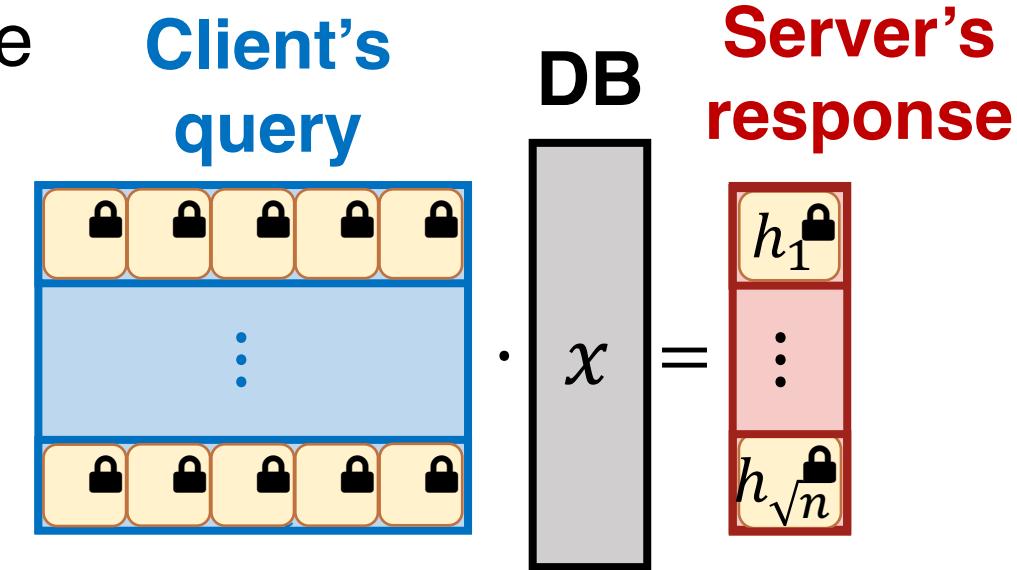
Client encrypts offline query component-wise  
with **linearly homomorphic encryption**

- E.g., Paillier

Server homomorphically computes the hint

Client decrypts hint

Online phase remains unchanged



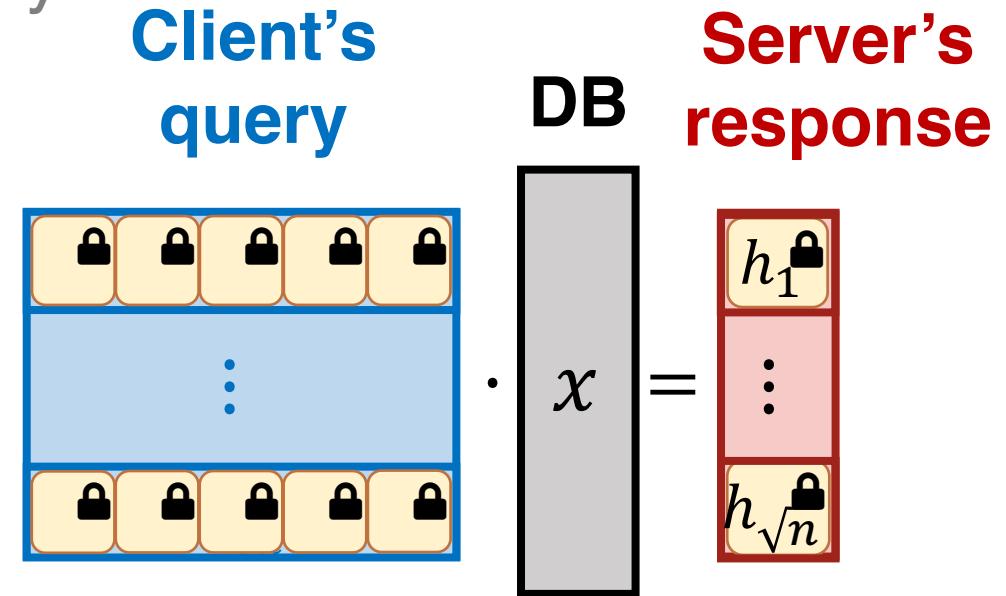
# From 2 servers to 1 – efficiency

Make server's offline response **linear** in its query

Client encrypts offline query component-wise  
with **linearly homomorphic encryption**

**Efficiency**: as described, the client

- uploads  $\sqrt{n} \times n$  matrix  $\rightarrow n^{3/2}$  ciphertexts
- downloads  $\sqrt{n}$  ciphertexts



Can use standard PIR rebalancing [CGKS95] to get  **$n^{3/4}$  total communication**

Idea: Divide DB into blocks of size  $\sqrt{n}$  of instead of single bits

More tricks to get  $n^{2/3}$  communication and quasilinear server time

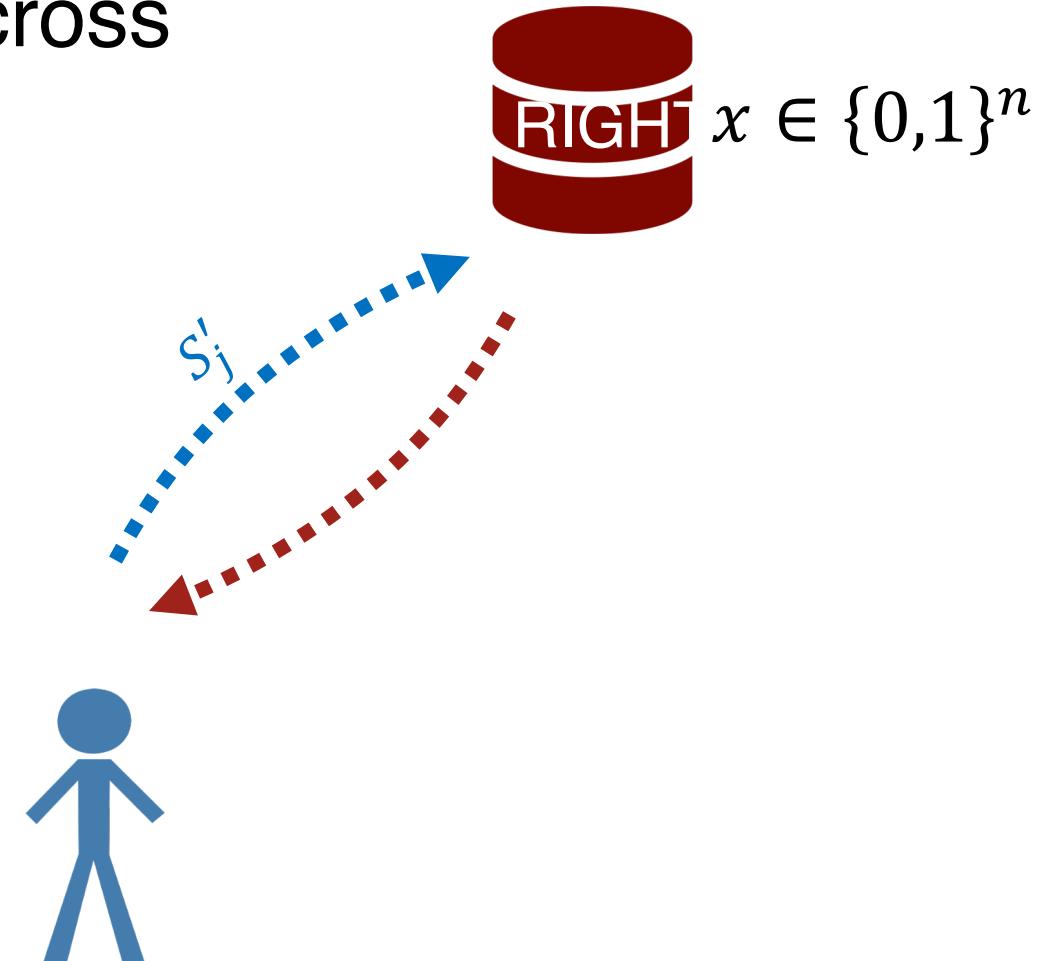
# Multiple queries

**Goal:** amortize cost of offline phase across multiple **adaptive** online queries

**Problem:** cannot reuse  $S_j$

Given  $S_j^1$  and  $S_j^2$ , server knows  $i_1 = S_j^2 \setminus S_j^1$

We will see a generic approach  
(using “Batch PIR” ideas)



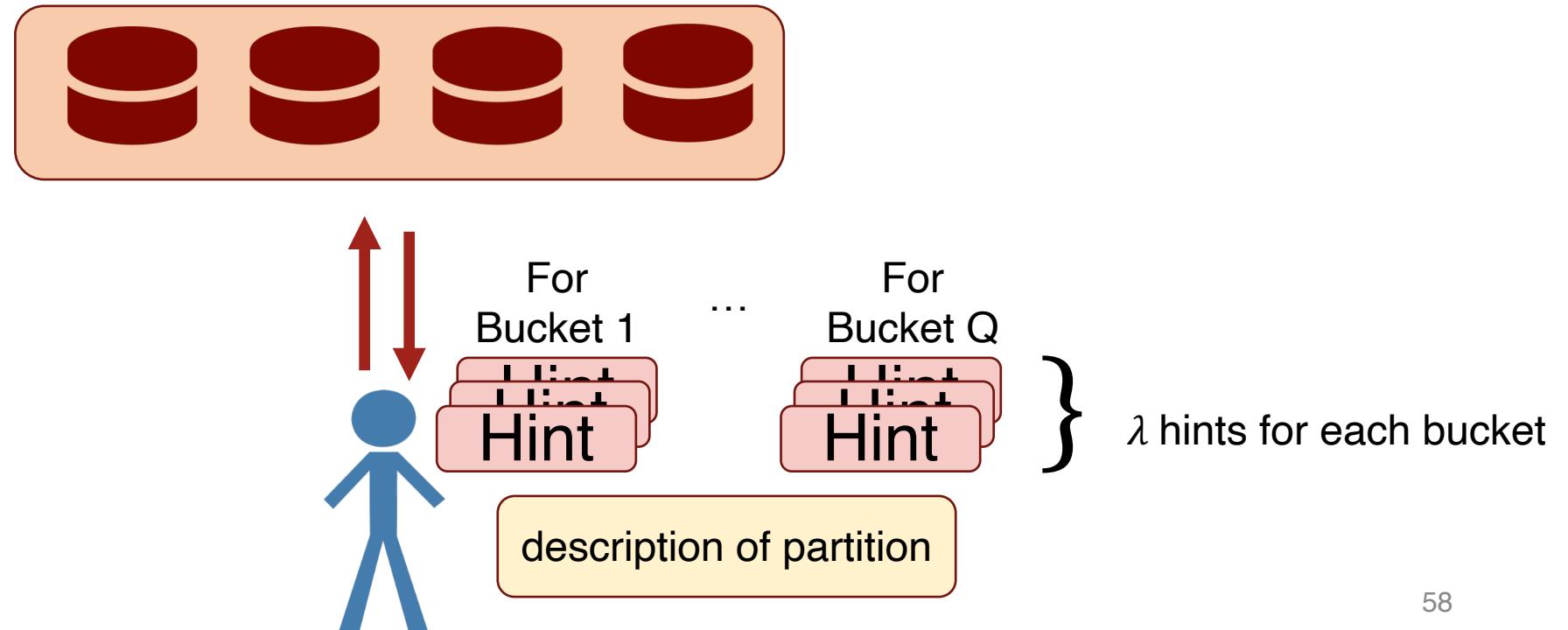
Hint  $S_1, h_1, \dots, S_j, h_j, \dots, S_m, h_m$

# Multiple adaptive queries – offline phase

Partition the DB into  $Q$  random buckets

Run the offline phase  $\lambda$  times for each bucket

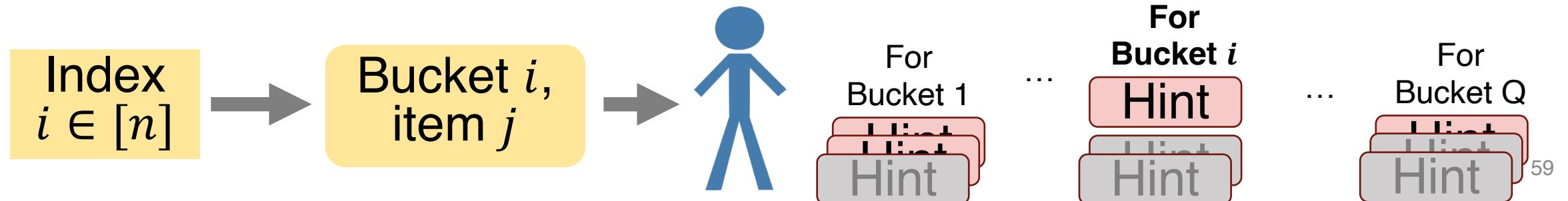
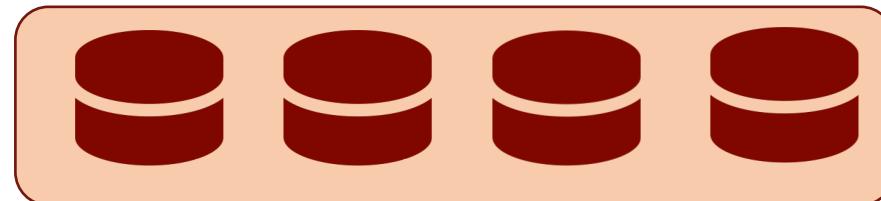
$x \in \{0,1\}^n$   
 $Q$  buckets



# Multiple adaptive queries – online phase

- Find matching bucket and unused hint

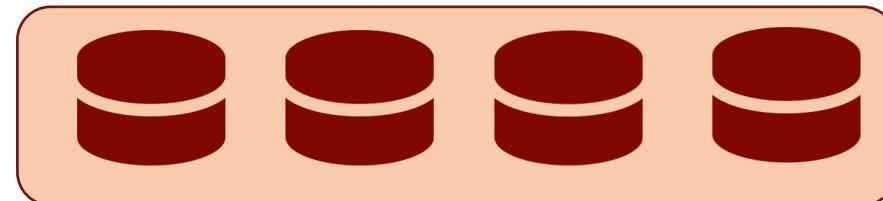
$x \in \{0,1\}^n$   
 $Q$  buckets



# Multiple adaptive queries – online phase

- Find matching bucket and unused hint
- Query using unused hint, server answers w.r.t. each bucket

$x \in \{0,1\}^n$   
 $Q$  buckets



$q \leftarrow \text{Query}(h_i, j)$



$a_1, \dots, a_Q$

Index  
 $i \in [n]$

Bucket  $i$ ,  
item  $j$



For  
Bucket 1  
Hint  
Hint  
Hint

For  
Bucket  $i$   
Hint

Hint  
Hint

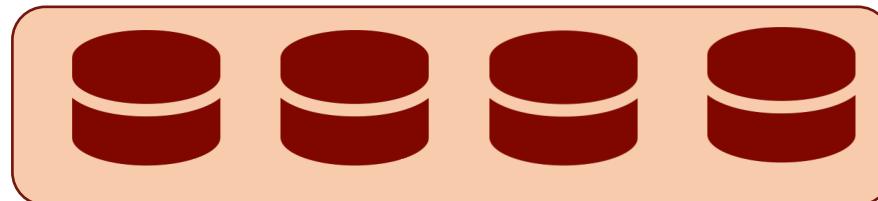
...  
...

For  
Bucket Q  
Hint  
Hint

# Multiple adaptive queries – online phase

- Find matching bucket
- Query using unused hint, server answers w.r.t. each bucket
- Reconstruct using matching response

$x \in \{0,1\}^n$   
 $Q$  buckets



$q \leftarrow \text{Query}(h_i, j)$



$a_1, \dots, a_Q$

Index  
 $i \in [n]$

Bucket  $i$ ,  
item  $j$



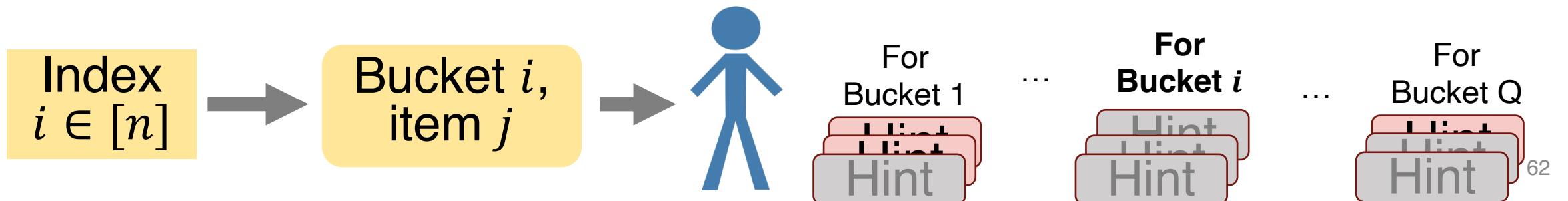
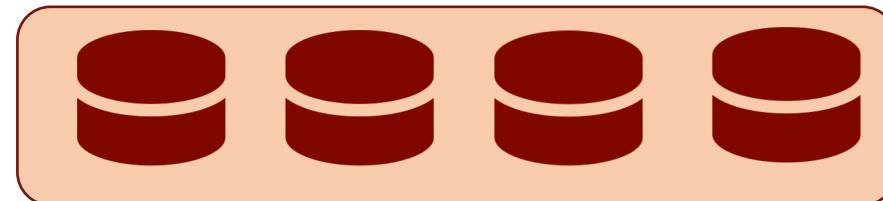
$x_{ij} \leftarrow \text{Reconstruct}(h_i, a_i)$

# Multiple adaptive queries – online phase

- Find matching bucket
- Query using unused hint, server answers with respect to each bucket
- Reconstruct using matching response
- Mark the hint as used

$x \in \{0,1\}^n$

$Q$  buckets



# Multiple adaptive queries – analysis

## Correctness:

The probability that a bucket overflows is negligible.

If no overflows, correctness follows from base PIR scheme.

## Security:

Server sees a single base query each time,  
security follows from base PIR scheme.

# Multiple adaptive queries – analysis

## Efficiency:

Online server time is  $T'(n) = Q \cdot T\left(\frac{n}{Q}\right)$

**2-server PIR:** our  $T = \sqrt{n}$  scheme and  $Q = n^{1/3} \Rightarrow T' = n^{2/3}$

- Different (non-generic) approach gives:  $T' = \sqrt{n}$  without additional assumption (only PRGs)

**1-server PIR:** our  $T = n^{2/3}$  scheme and  $Q = n^{1/4} \Rightarrow T' = n^{3/4}$

- Different (non-generic) approach gives:  $T' = \sqrt{n}$  from FHE

# Offline-online PIR with DB updates

When a database is repeatedly changing:

- standard PIR “just works”
- offline-online PIR requires the client to get an up-to-date hint

Naïve approach: rerun the offline phase after each change

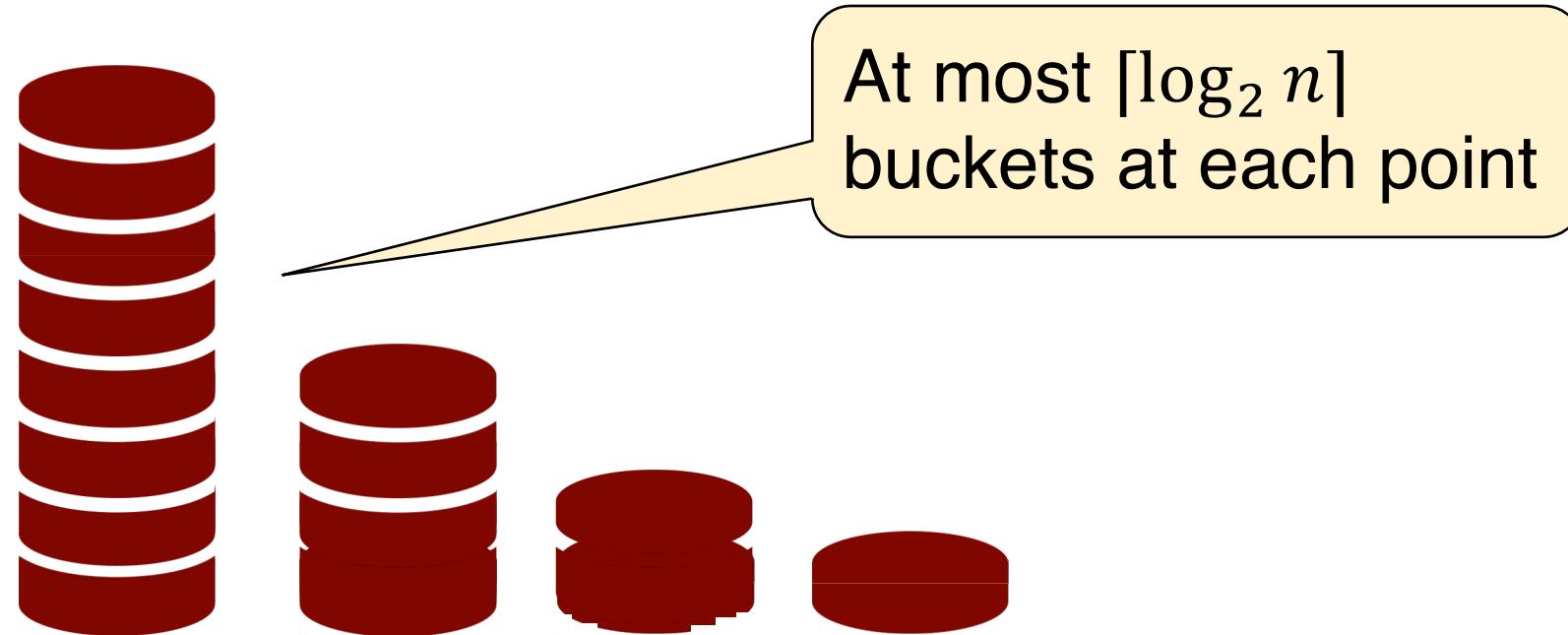
- **Linear** amount of server work on each change

Refined approach: incremental preprocessing

- **Logarithmic** (amortized) amount of server work on each change

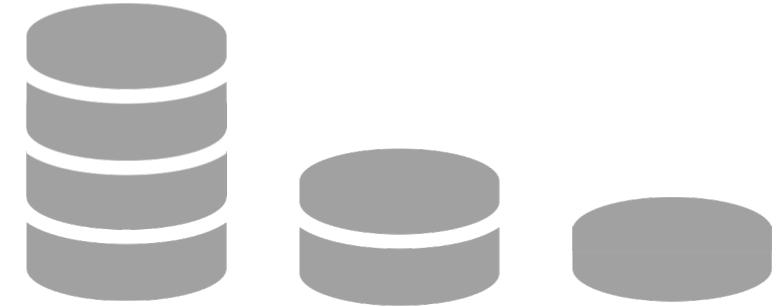
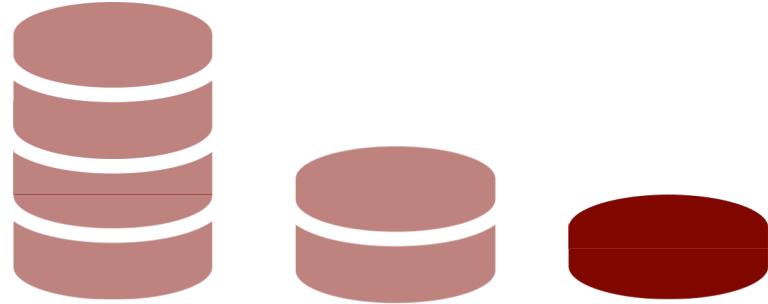
# Incremental preprocessing

- Database divided into **buckets**
- New record gets a new bucket
- Merge buckets with equal size



# Incremental preprocessing

## Offline phase



Hint

$\log n$  overhead – each record is preprocessed at most once in each bucket

Hint for each bucket



Hint

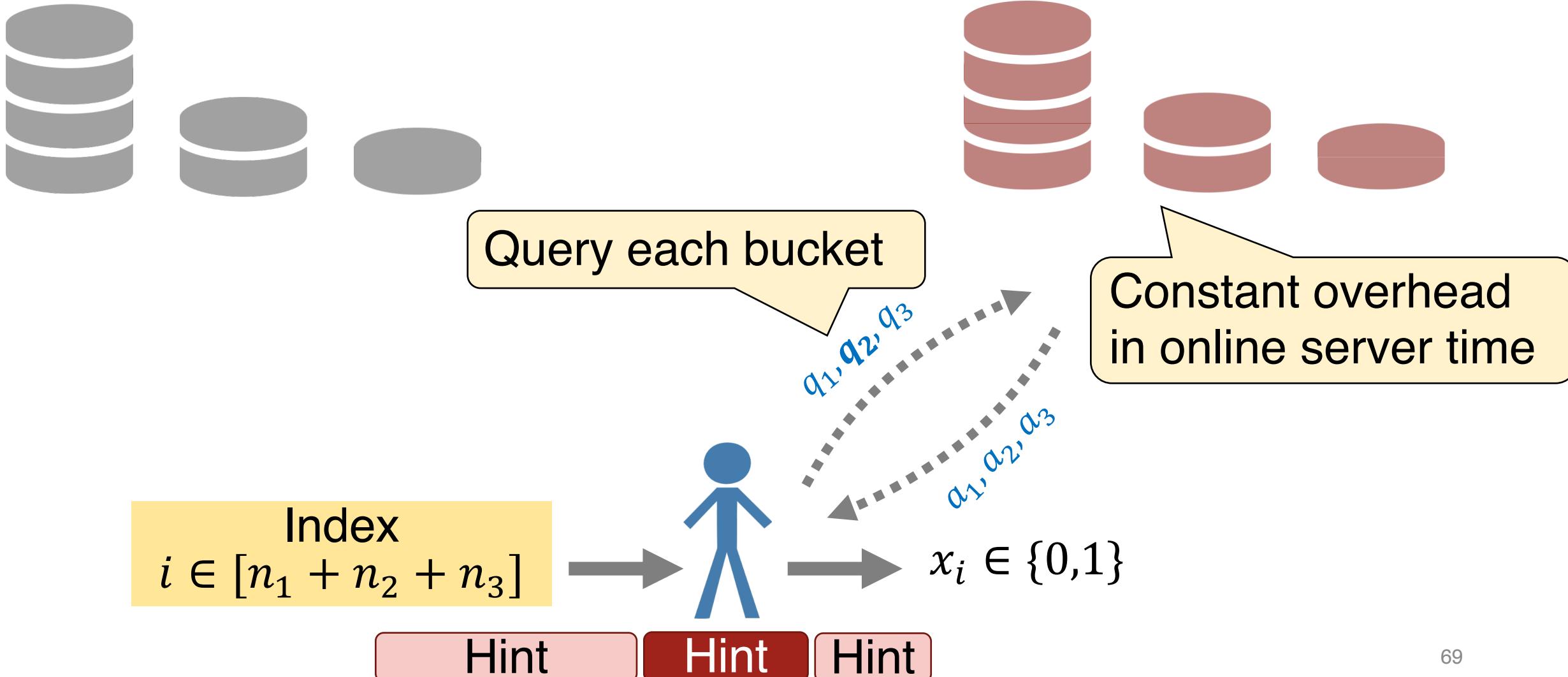
Hint

Hint

# Incremental preprocessing

## Online phase

[MZRA22] More efficient  
incremental offline/online PIR

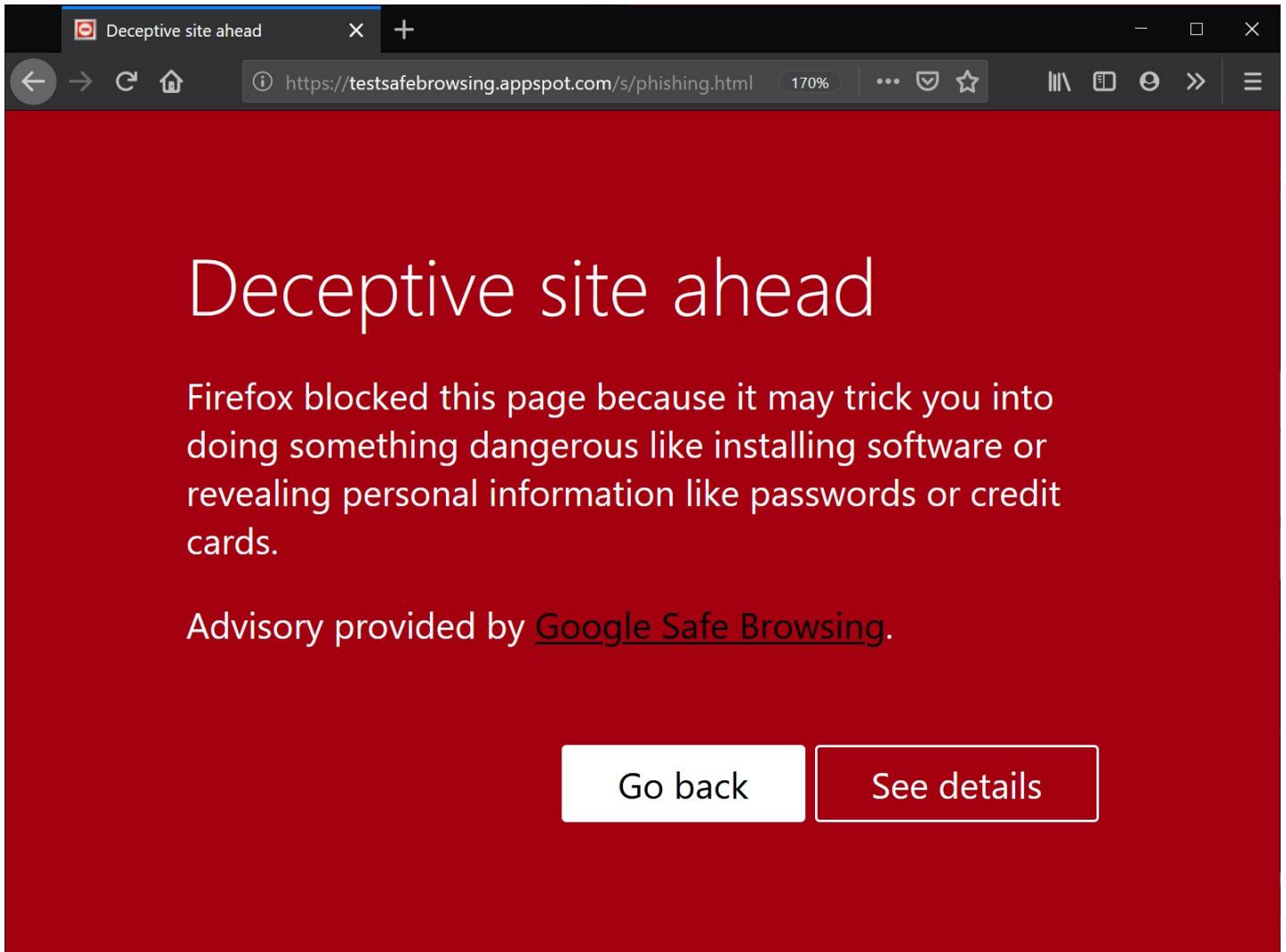


# Putting Offline-Online PIR to the test

Built a system for private Safe Browsing queries

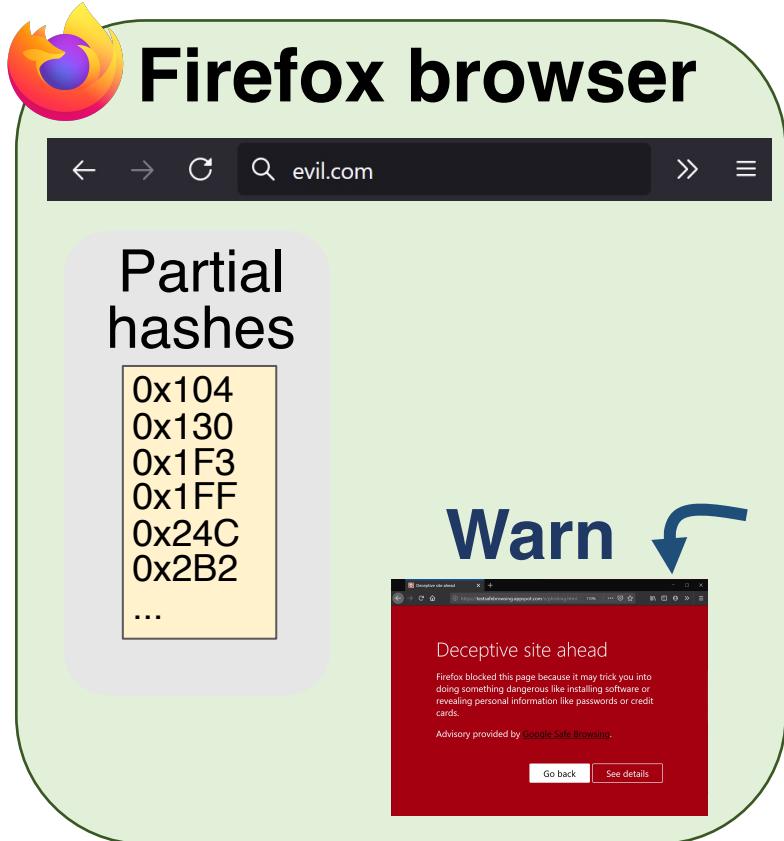
Integrated and evaluated with Firefox browser

Results: compare offline-online & DPF-based PIR



# Safe Browsing Today

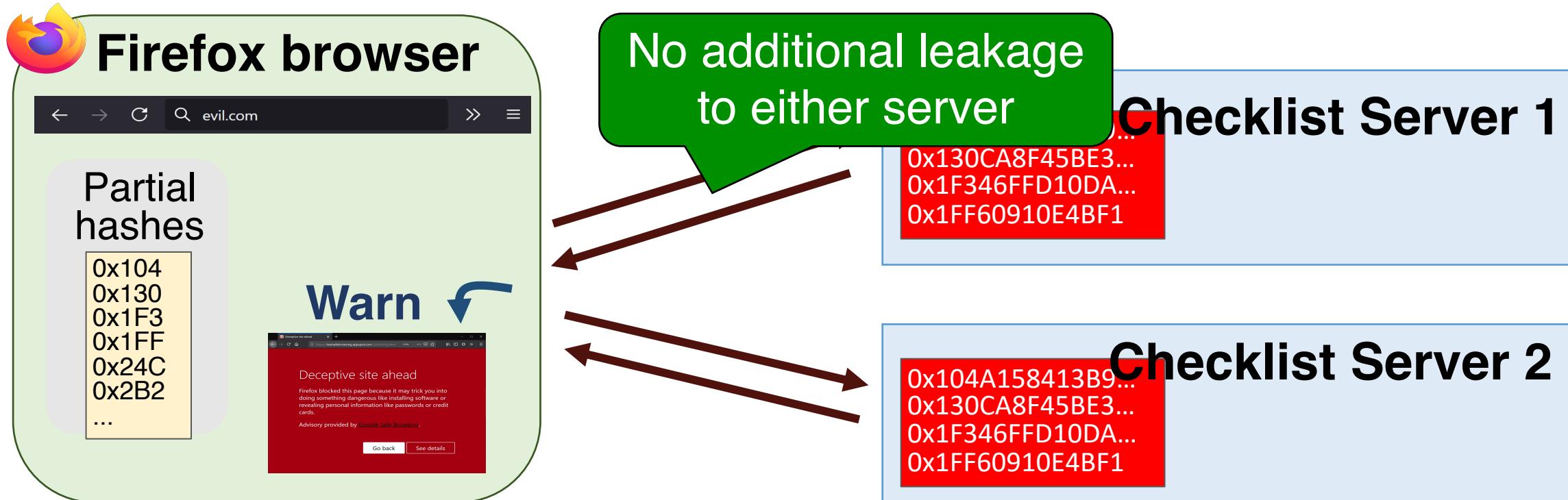
Leaks information about user's browsing history [GKL16,Per19,BK20]



# Checklist: A system for private blocklist lookups

(applied to Safe Browsing)

- Roughly 2500 lines of Go + 500 lines of C
- Supports new offline/online and standard (DPF-based) PIR



# Estimating the Safe Browsing parameters

**Monitor** a week of Safe Browsing requests and responses

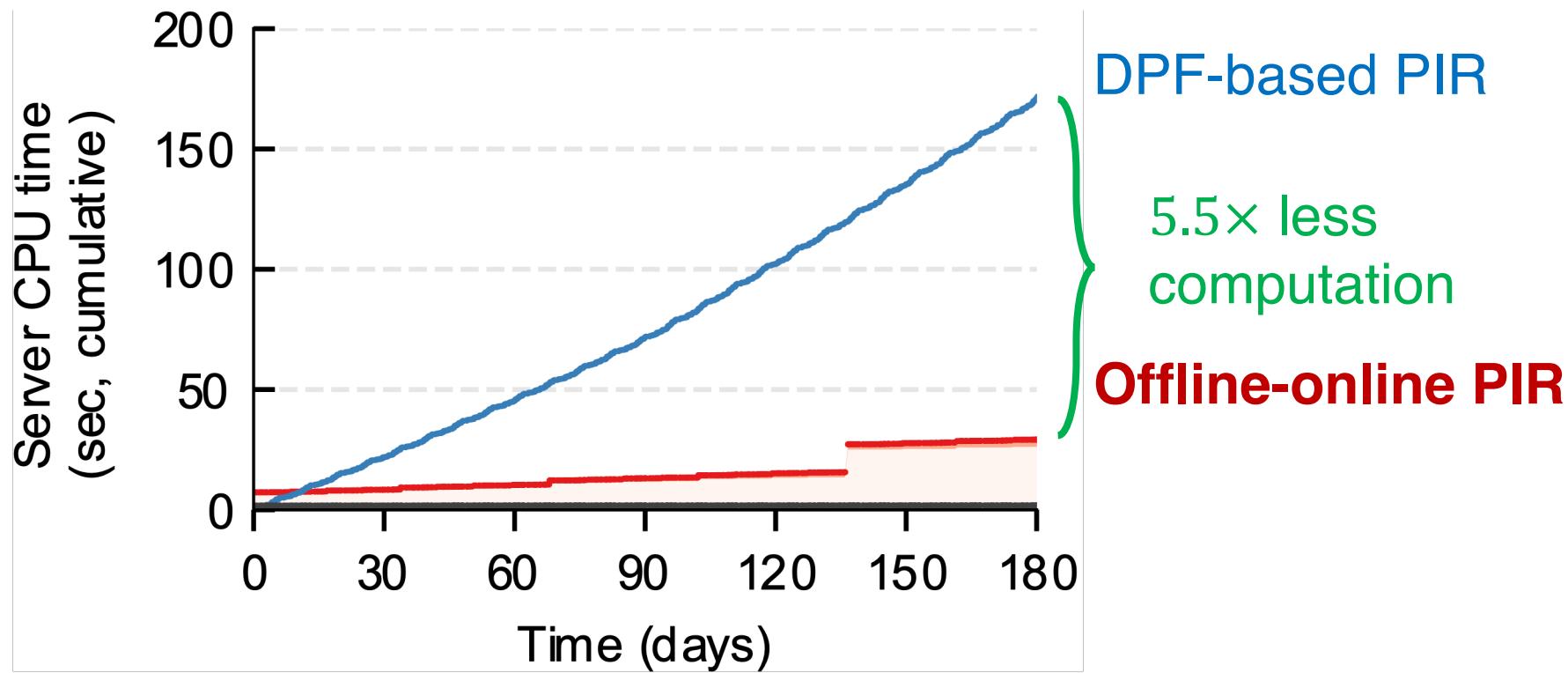
Local proxy forwards requests to the real Safe Browsing API

**Deduce:**

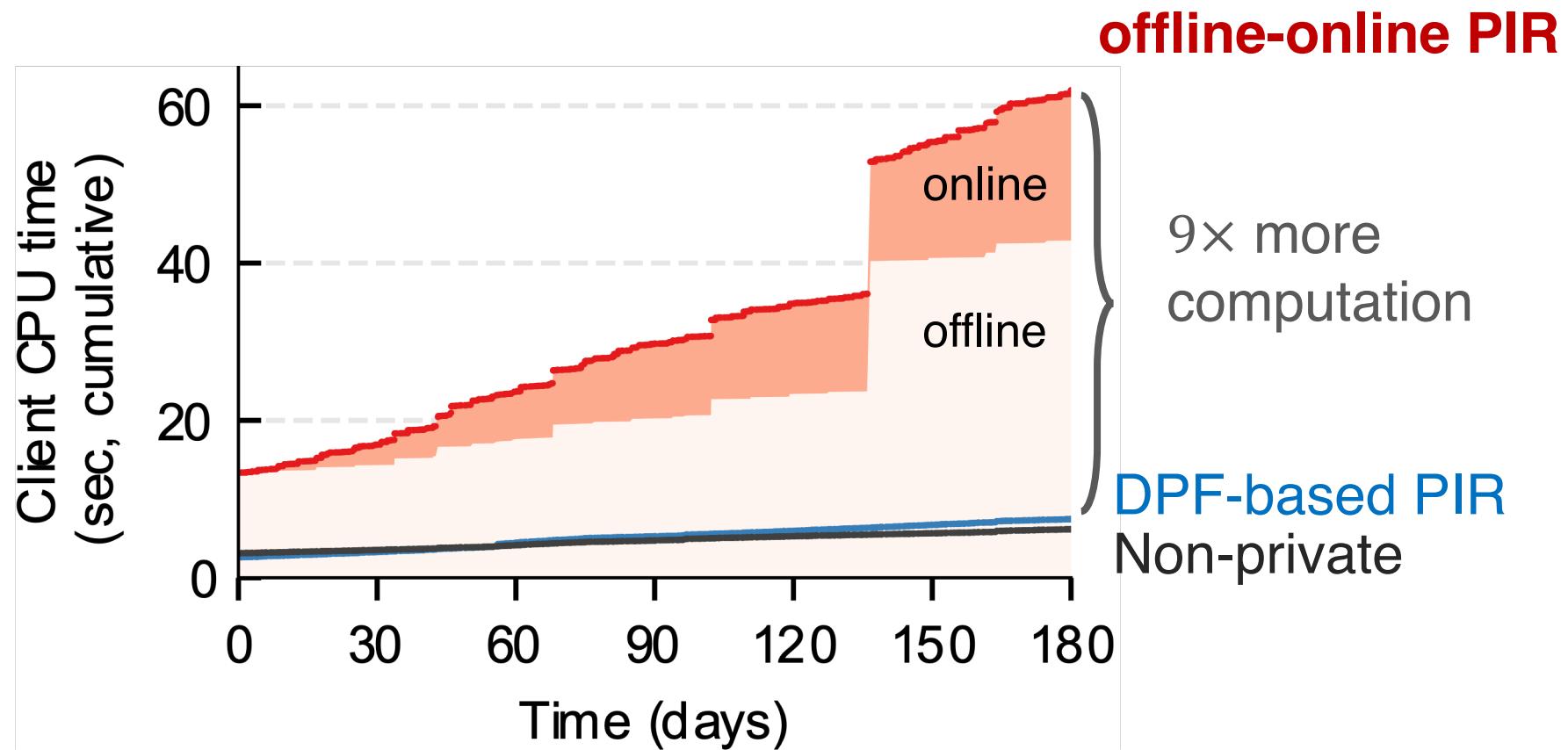
- Frequency of lookups
- Frequency of updates
- Database growth

# Offline/online PIR is more efficient for the server than traditional PIR

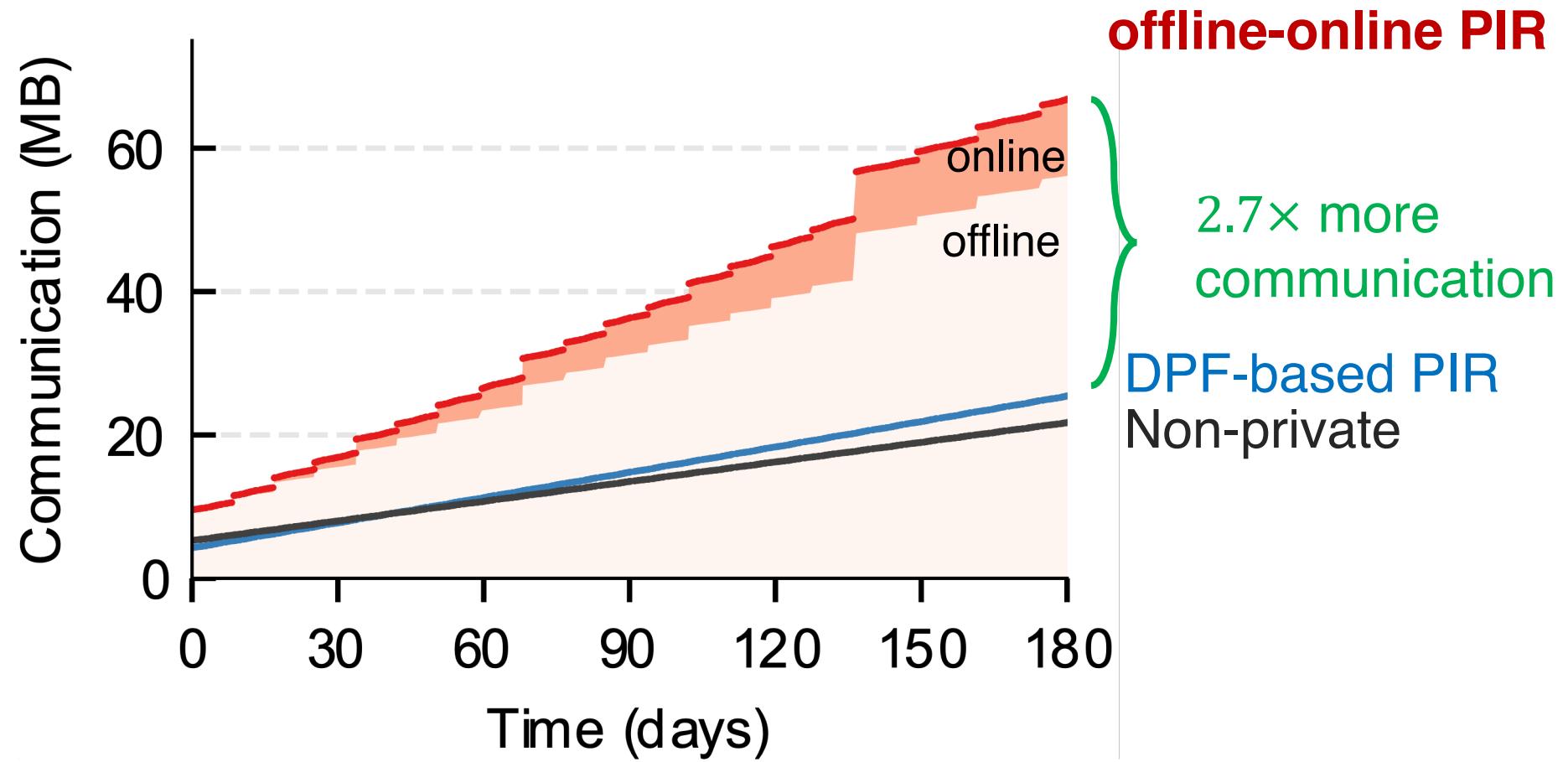
Replay recorded trace of a single user



# Traditional PIR is more efficient for the client than offline/online PIR



# Traditional PIR requires less communication than offline/online PIR



# Conclusion

- PIR has an array of applications to privacy-preserving systems
- Costs of PIR are now within the realm of practicality
  - Several alternatives that allow for different trade-offs

**Next steps:** Make it work in the real world

[dkogan@cs.stanford.edu](mailto:dkogan@cs.stanford.edu)

<https://cs.stanford.edu/~dkogan/>