

PRIVATE INFORMATION RETRIEVAL

- TODAY :
- Motivation
 - Two-server PIR scheme
 - Computational single server PIR scheme

Last two lectures: "generic" MPC

works for any computation (described as a circuit)

can be costly!

Solutions we saw have communication complexity linear in the size of the circuit...

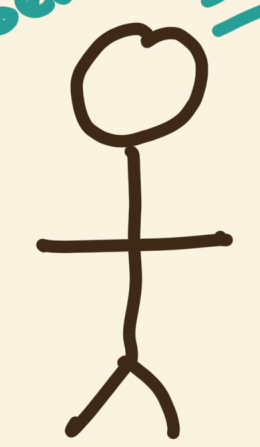
Can we do better? Say, for some specific class of computations?

Yes! And we'll see one such special-purpose MPC today for the computation known as **Private Information Retrieval**.

Motivation

E.g. Private browsing (simplified)

would like to search about sensitive data



query: peanut allergies

response: relevant info about peanuts

Server



Server learns your query... maybe you wish to keep your sensitive data private so info is not sold to ad platforms or insurance companies

We can perform the DB lookup privately using MPC techniques we've already seen.
But circuit for DB lookup is linear in size of DB \leftarrow can be really expensive :(

$$\text{MPC for } f_{\text{lookup}}(\text{DB}, i) = (\text{DB}[i], \perp)$$

\nearrow \times 's output \nwarrow \square 's output

Private Information Retrieval (PIR) is a specialized MPC which will allow us to query the DB without the server learning any info about the query!

wait a min... is this not just OT?
No! OT is in fact stronger: it required privacy from the user too. In PIR, we don't need to hide $\text{DB}[j]$ for $j \neq i$ from the user. It is one-sided privacy!

First (Trivial) Construction

Since, we don't care about DB's privacy for the user...

Just send the entire DB!

Perfectly
secure 

BUT

Utterly
inefficient 

We want to reduce communication between server and client.
Asymptotically, we want sublinear in the size of DB!

Is this ever possible with perfect security?

No. if $|comm.| < |DB|$, we must lose some info. about some entry j . $j \neq i$
since $DB[i]$ is what the user wishes to learn. Unbounded server can identify j
and now learns $i \neq j$.

We can get around this issue in two ways

Multiple (Non-colluding)
Servers [CGKS95]

- Replace single server with 2 or more non-colluding servers
 ↖ not allowed to communicate
- can still achieve perfect security

we'll cover both today!

Computational
Security [K097]

- Security against computationally bounded server
- Relies on cryptographic primitives

Two Server PIR

Two servers: S_0 & S_1

DB: bit-vector of length N : $DB \in \{0,1\}^N$

Goal: For input $i \in [N]$, user learns $DB[i]$ without revealing i to S_0 or S_1 .

Def. A **two server PIR scheme** is a tuple of eff. algs.:

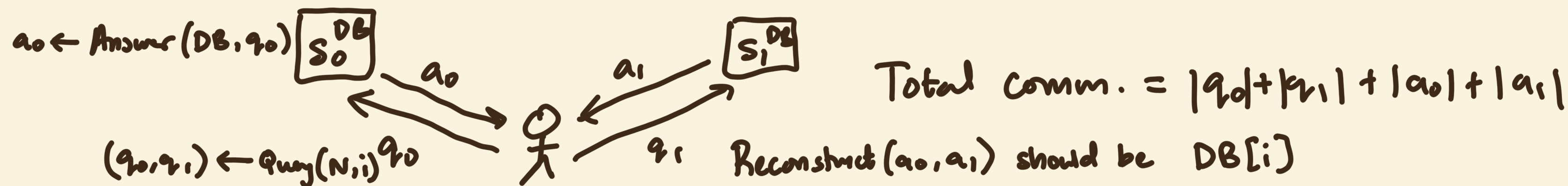
• $Query(N, i) \rightarrow (q_0, q_1)$: $N = |DB|$, i is the query index

query strings $\begin{cases} q_0 \xrightarrow{\text{sent to}} S_0 \\ q_1 \xrightarrow{\text{sent to}} S_1 \end{cases}$

• $Answer(DB, q) \rightarrow a$: $DB \in \{0,1\}^N$, q is a query string from the client/user, a is the server response

• $Reconstruct(a_0, a_1) \rightarrow b$: a_i are server responses, b is reconstructed entry

$S_0 : a_0 \xrightarrow{\quad} \text{Client}$
 $S_1 : a_1 \xrightarrow{\quad} \text{Client}$



Correctness. $\forall N \in \mathbb{N}, i \in [N], DB \in \{0,1\}^N,$

$$\Pr \left[b = DB[i] : \begin{array}{l} (q_0, q_1) \leftarrow \text{Query}(N, i) \\ a_0 \leftarrow \text{Answer}(DB, q_0) \\ a_1 \leftarrow \text{Answer}(DB, q_1) \\ b \leftarrow \text{Reconstruct}(a_0, a_1) \end{array} \right] = 1$$

Security. \forall poly-bounded $N \in \mathbb{N}, DB \in \{0,1\}^N,$ pairs $(i,j) \in [N]^2, b \in \{0,1\},$

$$\{q_b : (q_0, q_1) \leftarrow \text{Query}(N, i)\} \approx \{q_b : (q_0, q_1) \leftarrow \text{Query}(N, j)\}$$

↑ perfect/statistical/computational

Note. we're looking at query strings sent to each server individually.

A (nother) Trivial Construction

Query(N, i):

- Sample $q_0 \leftarrow \{0,1\}^N$ and set $q_1 := q_0 \oplus e_i$

- Output (q_0, q_1)

$$(0, \dots, 0, \overset{\text{i-th position}}{1}, 0, \dots, 0)^T$$

Answer(DB, q):

- Output $a \leftarrow \langle DB, q \rangle$

← recall inner product: $\sum_{i=1}^N DB[i] \cdot q[i] \pmod{2}$

Reconstruct(a_0, a_1):

- Output $a_0 \oplus a_1$

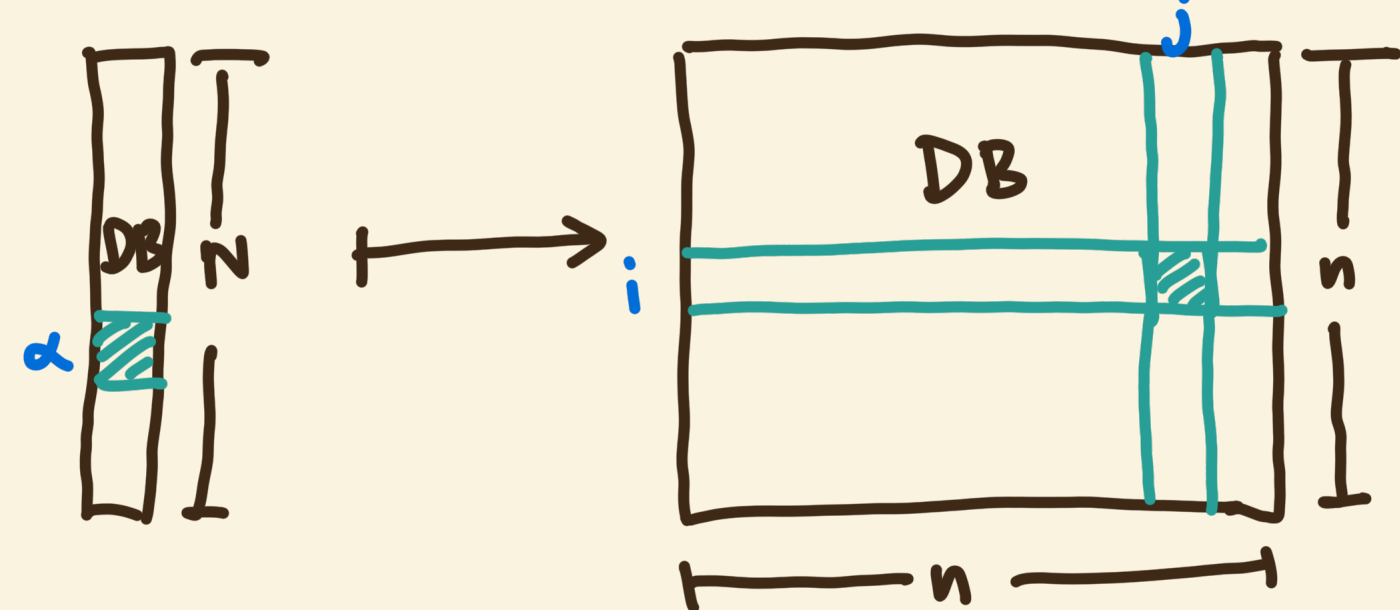
Correctness. $a_0 \oplus a_1 = \langle DB, q_0 \rangle \oplus \langle DB, q_1 \rangle$
 $= \langle DB, q_0 \oplus q_1 \rangle$
 $= \langle DB, e_i \rangle$
 $= DB[i]$

Security. Marginally, q_0 and q_1 are uniformly random N -bit vectors
 server S_i only sees q_i , and so, the distributions are identical!

But alas, communication is still linear: (← traded download (server → client)
 for upload (client → server)

[CGKS95] Construction

Assume $N = n^2$. We can view our DB as a matrix.



$i = \lfloor \alpha / n \rfloor$
 $j = \alpha \bmod n$ ($\alpha = n \cdot i + j$)

We will construct

$q_0 \oplus q_1 = e_j =$

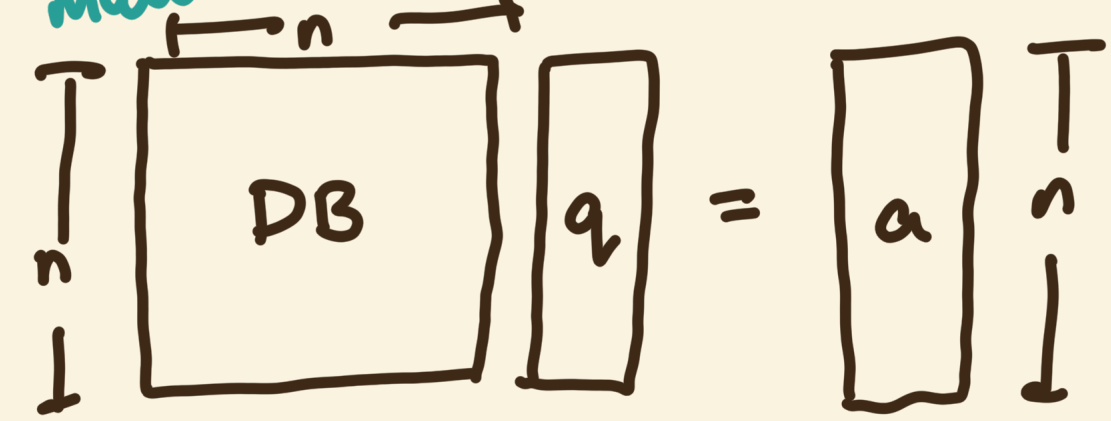
← j^{th} position

The diagram shows a vertical column representing the unit vector e_j . It consists of a sequence of bits: 0, ..., 0, 1, 0, ..., 0. The '1' is at the j^{th} position from the top, indicated by an arrow and the text '← j^{th} position'.

Query $(N, (i, j))$: interpreted as a pair

- Sample $q_0 \leftarrow \{0,1\}^n$, and set $q_1 := q_0 \oplus e_j$ vs. N from before
- Output (q_0, q_1)

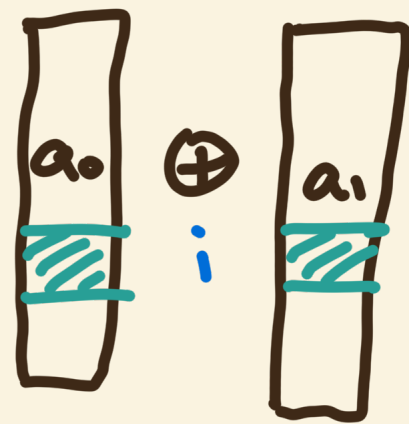
Answer (DB, q) :

- Compute $DB \cdot q \rightarrow a$ matrix-vector mult.
- 

- Output a .

Reconstruct (a_0, a_1) :

- Output $b \leftarrow (a_0 \oplus a_1)[i]$



Correctness. $a_0 \oplus a_1$

$$\begin{aligned} &= DB \cdot q_0 \oplus DB \cdot q_1 \\ &= DB \cdot (q_0 \oplus q_1) \\ &= DB \cdot e_j \\ &= DB_j \leftarrow j^{\text{th}} \text{ col of } DB \end{aligned}$$

$$\begin{aligned} \text{So, } (a_0 \oplus a_1)[i] &= DB_j[i] \\ &= DB[i, j] \checkmark \end{aligned}$$

Security. similar to before:

q_0 and q_1 are unif. random (marginally)

Communication

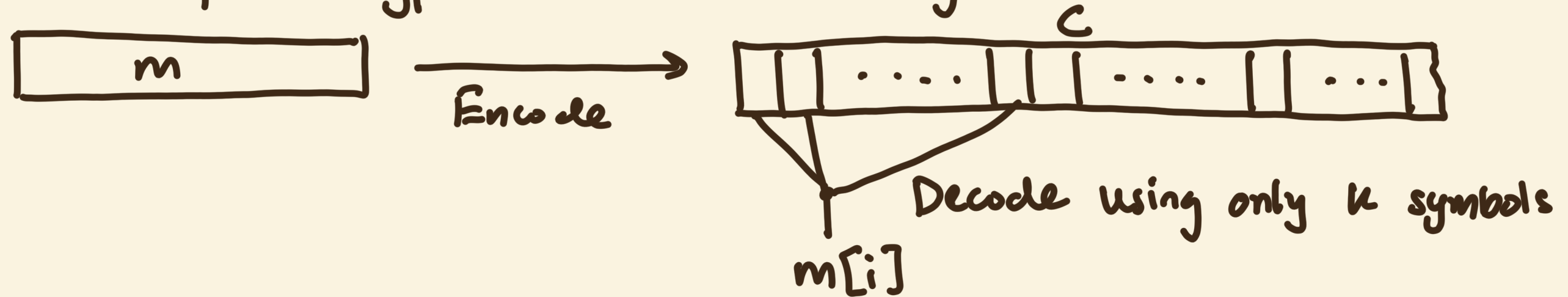
Complexity. $O(n) = O(\sqrt{N})$

bits! (upload = download)

(Note that client learns all values in the column but this is fine for PIR!)

Aside. Connection to Locally Decodable Codes (LDCs).

Perfectly secure (or information-theoretic) PIR is closely related to LDCs in coding theory. LDCs are a special type of error-correcting code:



Decoding any bit of the message does not require reading the whole codeword

Using LDCs, we can do even better: $O(N^{\sqrt{\log \log N / \log N}}) = O(N^{o(1)})$ comm. due to [DG15]

Computational Single Server PIR

Using two servers, we could hide the query from both servers?

Other ways to "hide" data? Encryption! What if we encrypt the query vector?

We wish to compute $DB \cdot q$, but standard encryption does not support this...



Enter
Linearly Homomorphic
Encryption

Briefly, Linear Homomorphic Encryption (LHE) is an enc. scheme that also satisfies the following property:

Linear Homomorphism.

$$\forall k \in \mathcal{K}, \forall m_0, m_1 \in \mathcal{M}, \forall c_0, c_1 \in \mathbb{Z}$$

← abelian group

$$c_0 \cdot \text{Enc}(k, m_0) + c_1 \cdot \text{Enc}(k, m_1) = \text{Enc}(k, c_0 m_0 + c_1 m_1)$$

group op. of ctxt space

group op. of msg space

[K097] use LHE to get Comp. Single Server PIR:

Idea. Do CGKS under LHE

E.g. ElGamal is a public-key LHE

Query(N, (i, j)):

• Sample $k \xleftarrow{\$} \mathcal{K}$

• Output $q := (\text{Enc}(k, e_j[1]), \text{Enc}(k, e_j[2]), \dots, \text{Enc}(k, e_j[n]))$

← vector of ciphertexts. Encryption of e_j vector component-wise

Answer(DB, q):

• Using the LHE, we can compute $a := \text{DB} \cdot q$

← scalars

← ctxts

$$a = \begin{matrix} \boxed{\text{DB}} & \begin{matrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{matrix} & = & \begin{bmatrix} \text{DB}_1[1] \cdot c_1 + \text{DB}_2[1] \cdot c_2 + \dots + \text{DB}_n[1] \cdot c_n \\ \vdots \\ \text{DB}_1[n] \cdot c_1 + \text{DB}_2[n] \cdot c_2 + \dots + \text{DB}_n[n] \cdot c_n \end{bmatrix} \end{matrix}$$

by the
linear homomorp-
hism property

$$\begin{aligned} &= \begin{bmatrix} \text{Enc}(k, \sum_{i=1}^n \text{DB}_i[i] \cdot e_j[i]) \\ \vdots \\ \text{Enc}(k, \sum_{i=1}^n \text{DB}_i[n] \cdot e_j[i]) \end{bmatrix} \\ &= \begin{bmatrix} \text{Enc}(k, \text{DB}_j[i]) \\ \vdots \\ \text{Enc}(k, \text{DB}_j[n]) \end{bmatrix} \end{aligned}$$

• Output a

Reconstruct(k, a):

• Decrypt all n ciphertexts in a

$$u \leftarrow (\text{Dec}(k, a[i]), \dots, \text{Dec}(k, a[n]))$$

• Output $u[i]$. \leftarrow which will be equal to $\text{DB}_j[i] = \text{DB}[i, j]$

Correctness. Follows from what we showed above.

Security (informal). Follows from semantic security of LHE.

Enc. of e_j is indist. from Enc of e_j ,

Communication. $|ctx| = O(\lambda)$ so overall comm. is $O(\lambda n) = O(\lambda \sqrt{N})$.

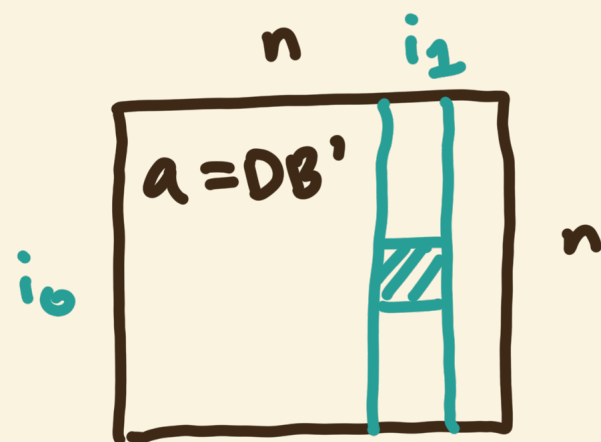
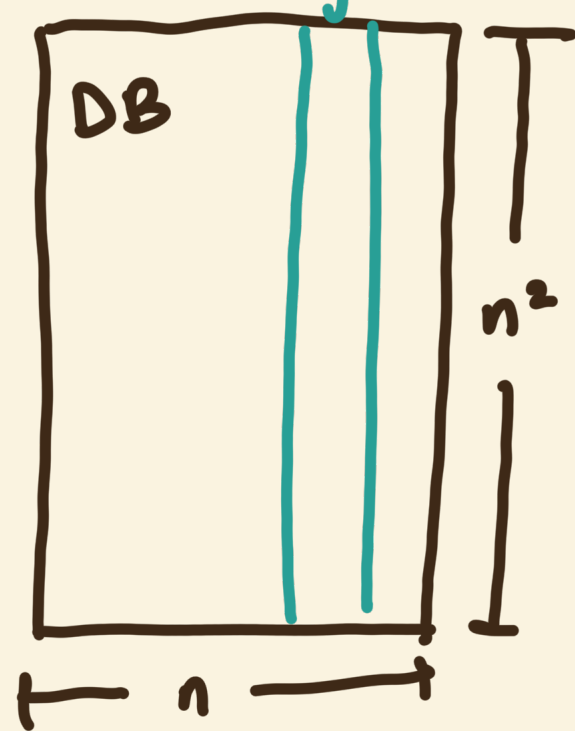
The question we keep asking: **Can we do better??**

Yes! We will use recursion :0.

The client is only interested in $DB[i, j]$ but the server is currently sending the entire column: DB_j .

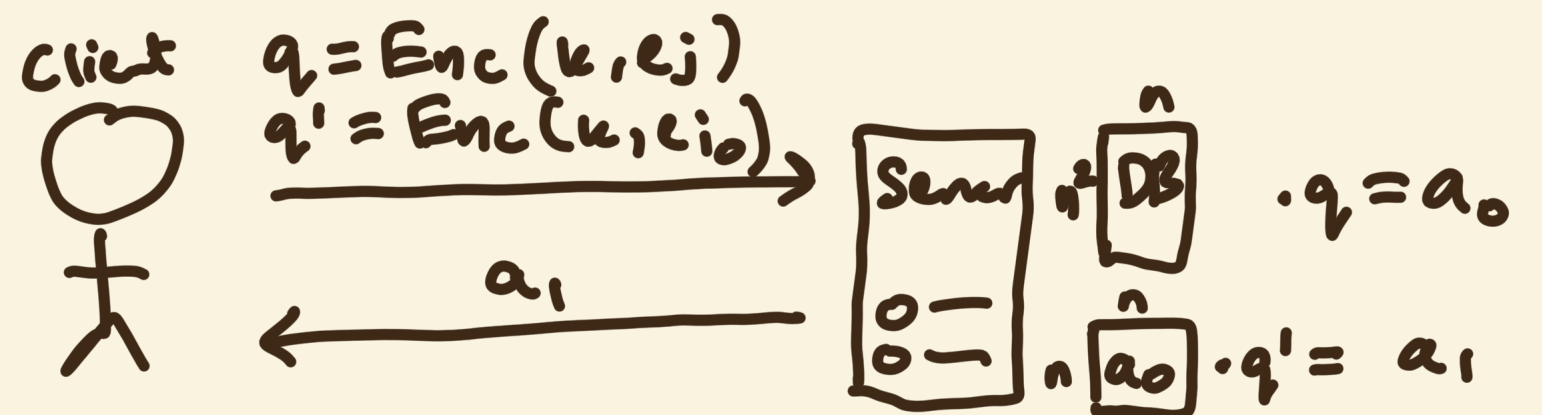
Idea. Think of DB_j as the database we do PIR on!

Suppose $N = n^3$.



$$d = (i_0 \cdot n + i_1) + j \quad i_0, i_1, j \in [n]$$

Client sends two selection vectors:



$$\text{Total comm.} = |q| + |q'| + |a_1|$$

$$= O(\lambda n)$$

$$= O(\lambda N^{1/3})$$

Nice!

Can we keep recursing?

Yes, but there is a tipping point.

Issue. $|entry| = 1 \text{ bit}$ but $|ctxt| = \lambda \text{ bit}$

Every time we recurse, we save on comm. w.r.t. to DB but blow up comm. by a factor of λ .

Why? Because we recurse on a sub-DB of ciphertexts NOT a DB of 1-bits...

So, ctxt of λ bits $\Rightarrow |a_0| = \lambda n$,

which means $|a_1| = \lambda (\lambda n) = \lambda^2 n$

\swarrow ctxt size \nwarrow sub-DB size

If you recurse l times, $|a_l| = \lambda^l n = \lambda^l N^{1/l}$ \leftarrow Can blow up very quickly...

Solution. Use LHE with (1) good "rate": $|m|/|c|$ is closer to 1

(2) large message space

E.g. Damgård - Jurik

State of the art. polylogarithmic in $|DB|$. See website for survey.