# Lecture 8 - Private Information Retrieval (PIR)

CS 355 - Spring 2019
April 24, 2019
Henry Corrigan-Gibbs

# Logistics

* HW2 Due <u>Friday, 4/26 at 5pm</u>
  ↳ Come to OH if you need help!
* Please give feedback on PSETS
* Also, anonymous feedback form online.
  ↳ Anything that would improve course or make us better teachers.
* Grades for HW1 out now
  ↳ regrade policy

# Plan

* Recap: Multiparty Computation
* PIR: What it is, why it's amazing
  ↳ Formal defn's
* Constructions
  - Two-server PIR
  - One-server PIR

# A "perfect" cryptosystem

1) Has a beautiful theory

2) Works in practice

3) Solves a problem that people <ins>should?</ins> care about

\* When you're working on a problem, ask yourself how your work does against this rubric

# Today

- One of my favorite "almost perfect" ideas in crypto

- Lots of activity, even in last few years
  → Even today at Stanford....

- A classic crypto result
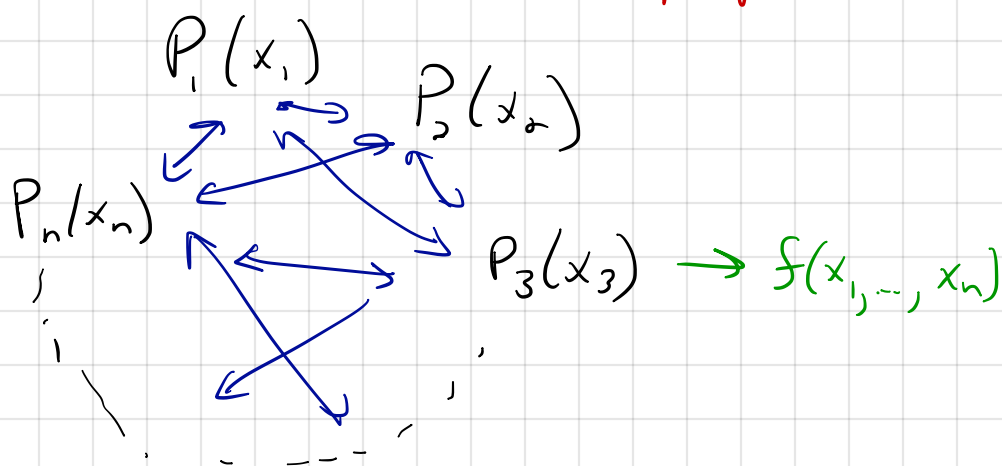  seems impossible, then turns out to be simple/elegant

The catch: For reasons we'll see, it's not quite practical yet.....

# Recap: MPC

- Each party $P_i$ holds secret input $x_i$

- Parties want to compute a joint function $f(x_1, ..., x_n)$ of their private inputs

... without leaking anything else!

→ "Best possible" result — can compute <u>any</u> fn in secure multiparty manner

$P_1(x_1)$
$P_2(x_2)$
$P_n(x_n)$
$P_3(x_3) \rightarrow f(x_1, ..., x_n)$

Why want this?

* Train a spam classifier over millions of peoples email w/o having to share mail

* Compute election results w/o having to publish votes

* Check if your password is in use in a database w/o leaking your password (or site leaking * DB)

⇒ Implies ZK ... e.g.

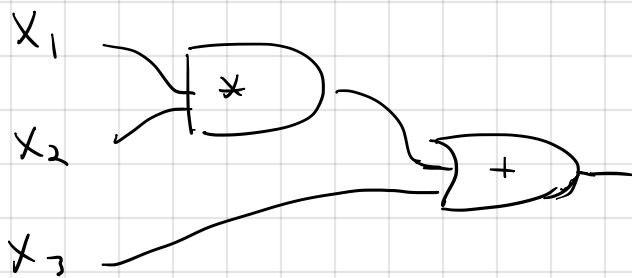$P(\text{Graph } G, \quad \text{3-coloring of } G)$          $V(\text{Graph } G)$

3 coloring of G is valid?

Many real-world complications

# Recap : MPC

Idea : View $f(x_1, ..., x_n)$ as an arithmetic ckt

gates are $+$ and $*$ mod $p$, wires are values in $\mathbb{F}_p$

$\uparrow$ Think: ints mod $p$

$x_1$
$x_2$
$x_3$

Note : Reexpressing computation $f(\cdot)$ as an arith ckt is without loss of generality.

$\hookrightarrow$ Any poly time computation has a poly-sized arith ckt

$\hookrightarrow$ If $f$ has Boolean ckt of size $S$, it has an arithmetic ckt of size $O(S)$.

## MPC Protocol (Ben-Or, Goldwasser, Wigderson '88)

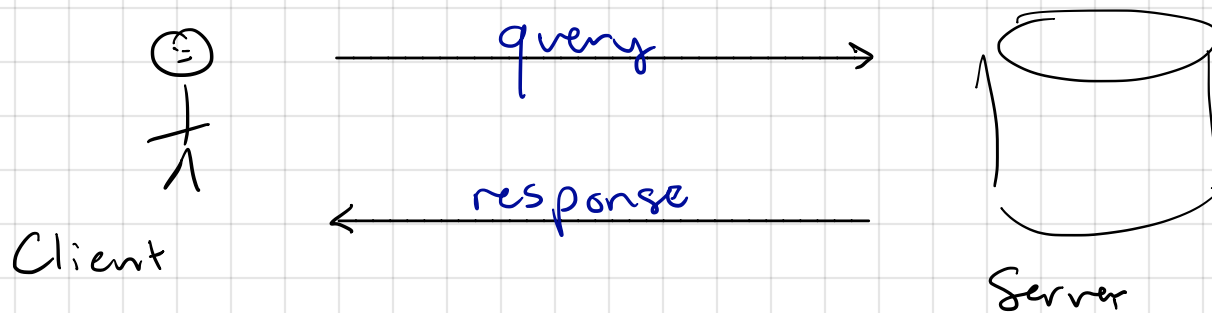* Parties start holding shares of input wires

* Parties jointly compute shares of internal wires

* Finally parties hold shares of output wire

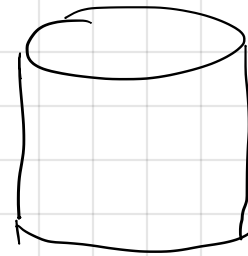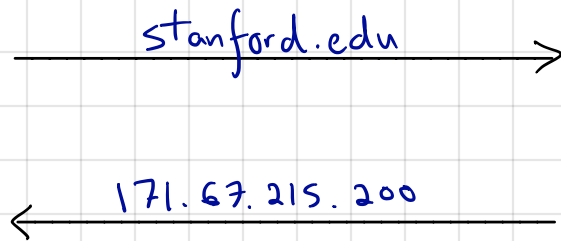$\hookrightarrow$ Publish shares to recover output $f(x_1, ..., x_n)$.

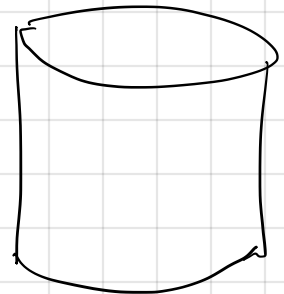# Private Information Retrieval

Every day on the Internet...



Client     query → ← response     Server

Examples:

① DNS

stanford.edu $\longrightarrow$

$\longleftarrow$ 171.67.215.200

TLD name server
for .edu

② WebMD

"fever" $\longrightarrow$

$\longleftarrow$ ⟨article about fever⟩

WebMD
web server

③ Many more: Querying stock price,
looking at courses in course catalog,
----

Notice: The client's query can be <u>sensitive</u>!
It can leak:
- What website you're visiting
- Your medical condition
- What stocks you're thinking about buying
⋮

Today, the client just sends this sensitive
data directly to the server!

# In systems today...



**You** — Your sensitive query → **Evil website** ← Useful response

Your sensitive query → **Data broker** ← $

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

**Question:** "Can you query a database without the database learning your query?"

**Trivial answer:** "Just download the entire DB."

The DB server doesn't learn your query...

Still, this scheme is somewhat unsatisfying

Let's ask a better question...

**Question:** "Can you query a database without the database learning your query *with communication <u>sublinear</u> in DB size?*"

**Answer:** Unconditionally.... No. ☹ [CGKS'95]

We won't prove this, but it's not too hard to show... See paper for 1-para proof.

What do we do when we get stuck in life?

**Option I:**

Change the model! (e.g. ROM)

What if we have two non-colluding copies of DB?

"two-server PIR"

[Can think of $k > 2$ non-colluding servers "k-server PIR"]

**Option II**

Make assumptions!

Under pretty basic assumptions (DDH, ....), we can build non-trivial single-server PIR

[Kushilevitz & Ostrovsky Focs'97]

Lets first consider two-server PIR...

# Private Information Retrieval  (Chor, Goldreich, Kushilevitz, Sudan) FOCS '95



Client

notation
$[n] = \{1, \ldots, n\}$

$i \in [n]$

$q_0$

$a_0$

$x_i \leftarrow \text{Recon}(a_0, a_1)$

$q_1$

$a_1$

Server 0

$x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \{0,1\}^n$

Server 1

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

To keep things simple, will focus on "one-round" PIR scheme
  ⌐→ One message from client to server ("query")
  ⌐→ One message from server to client ("response")

Also, we'll think about the DB as holding bits
  ⌐→ Can handle longer msgs

# Syntax : Three eff algs

$$(q_0, q_1) \leftarrow \text{Query}(i)$$

... for $i \in [n]$, where $n =$ length of database

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

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

# Properties

① Correctness: Client gets the bit it wants

$$\forall_{i \in [n]} \; \forall_{x \in \{0,1\}^n}$$

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

② Security: No single server learns the bit the client wants

$$\exists \; \text{eff Sim s.t} \quad \forall_\beta \in \{0,1\}$$

$$\left\{ \text{Sim}(\beta) \right\} \stackrel{\approx}{c} \left\{ q_\beta \; \middle| \; (q_0, q_1) \leftarrow \text{Query}(i) \right\}$$

<span style="color:green">Simulation! So useful!</span>

<span style="color:green">Can also be $\equiv$ for info. theoretic or "perfect" security.</span>
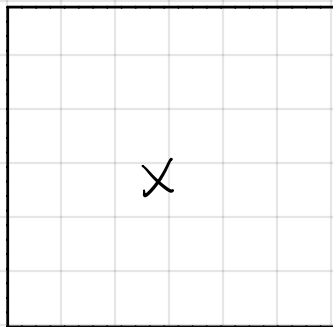
<span style="color:red">N.B. IF both servers collude and share their $q$'s all bets are off !</span>
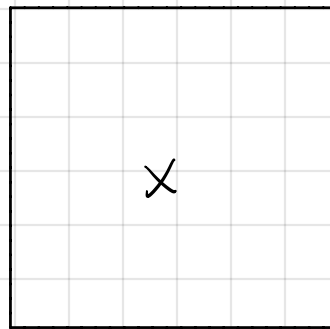
# An $O(\sqrt{n})$ - communication PIR scheme.

↳ Already very non-trivial!

View database as a matrix $X \in \mathbb{Z}_2^{\sqrt{n} \times \sqrt{n}}$

$q_0 = $ | 0 | 1 | 0 | 0 | 1 | 1 |

$q_1 = $ | 0 | 1 | 1 | 0 | 1 | 1 |

$X$

$X$

Client wants to read $X_{ij}$   $i, j \in [\sqrt{n}]$

## Query $(i, j) \rightarrow (q_0, q_1)$

Sample two random vectors $q_0, q_1 \in \mathbb{Z}_2^{\sqrt{n}}$
s.t.

$$q_0 + q_1 = (0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 0) \in \mathbb{Z}_2^{\sqrt{n}}$$

↳ $j$th position

return $(q_0, q_1)$

## Answer $(X, q) \rightarrow a$

Output $X \cdot q \in \mathbb{Z}_2^{\sqrt{n}}$.

## Reconstruct $(q_0, q_1) \rightarrow x_{ij}$

Compute $col_j \leftarrow q_0 + q_1 \in \mathbb{Z}_2^{\sqrt{n}}$
output $i$th element as $X_{ij}$.

# Why it works.

## Correctness:

$$a_0 + a_1 = Xq_0 + Xq_1$$

$$= X(q_0 + q_1)$$

$$= \left( X \begin{matrix} \vdots \\ x_j \\ \vdots \end{matrix} \right) \left( \begin{matrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \\ 0 \end{matrix} \right) \leftarrow j\text{th position}$$

$$= \left( x_j \right) \in \mathbb{Z}_2^{\sqrt{n}}$$

Then $i$th component of $a_0 + a_1$ gives $X_{ij}$.

## Security

$$\text{Sim}(\beta): \qquad q_\beta \xleftarrow{R} \mathbb{Z}_2^{\sqrt{n}}$$

$$\text{output } q_\beta.$$

Each query is distributed uniformly at random.

# What do we know about 2-server PIR?

* Without privacy, total communication $\geq \log n$ bits

* Best known lower bound (impossibility result) says for PIR
  Communication $\geq 5 \cdot \log n$ bits
  (Wehner & de Wolf '05)

* When I took CS355, best protocol had
  communication $\leq \tilde{O}(n^{1/3})$ (GKGS '98)
  (not too complicated)

* Relatively recent big result
  (Dvir & Gopi '15)    Communication $\leq n^{O(\sqrt{\log \log n / \log n})}$
  (not simple)

Today, we'll see a scheme that achieves
  Communication $\leq O(n^{1/2})$

$\Longrightarrow$ A good open Q: Are there better PIR schemes? $\Longleftarrow$

PIR scheme w/ comm complexity $\leq O(\log^2 n)$?

* <u>With</u> computational assumptions (PRGs), have very good
  <u>recent</u> schemes
  Communication $\leq O(\lambda \log n)$
                       $\uparrow$
                    Security
                    Parameter

  $\searrow$ These results are essentially
  best possible... but still
  good open Qs here.

# Single Server PIR: $O(\sqrt{n})$ communication

Say you have an additively hom semantically secure enc scheme
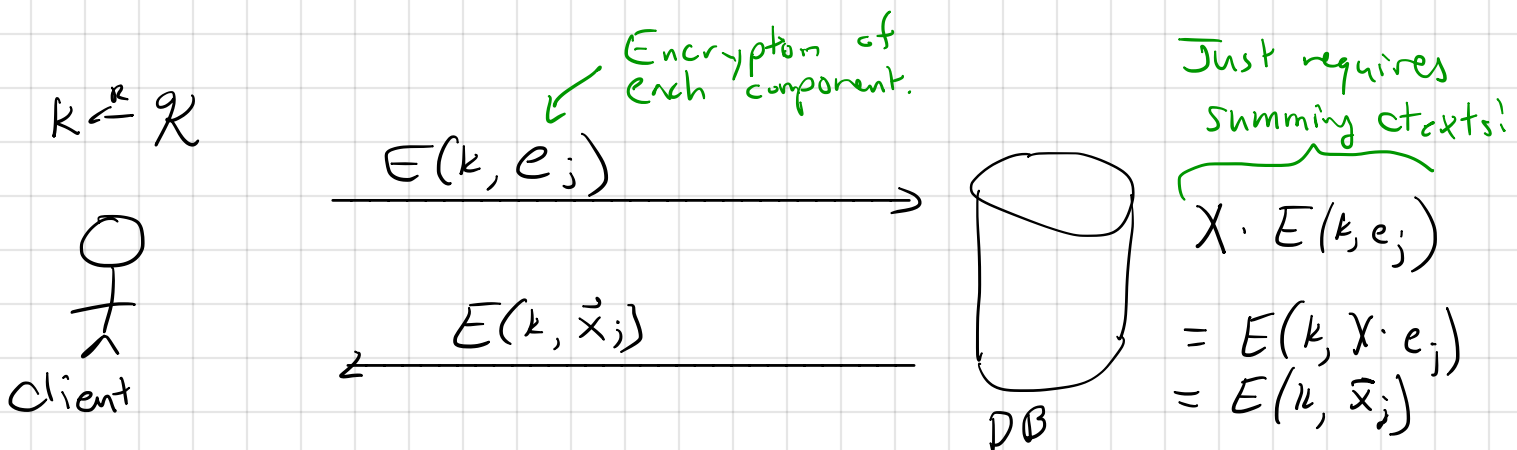
$$E(k, m_0) + E(k, m_1) = E(k, m_0 + m_1)$$

↳ Can build from DDH, Quadratic Residuosity, lattices, ...

IDEA: Client sends encrypted query vector, Server computes mut-ve product "under encryption"

Again write DB as matrix $X \in \mathbb{Z}_2^{\sqrt{n} \times \sqrt{n}}$

Client wants bit $(i,j) \in [\sqrt{n}]^2$.

Notation  $e_j = (0000 \cdots 010 \cdots 000)$
                            ↳ jth position

$k \xleftarrow{R} \mathcal{K}$

Encryption of each component.



$\xrightarrow{\quad E(k, e_j) \quad}$

Client

$\xleftarrow{\quad E(k, \vec{x}_j) \quad}$

DB

Just requires summing ctexts!

$X \cdot E(k, e_j)$

$= E(k, X \cdot e_j)$

$= E(k, \vec{x}_j)$
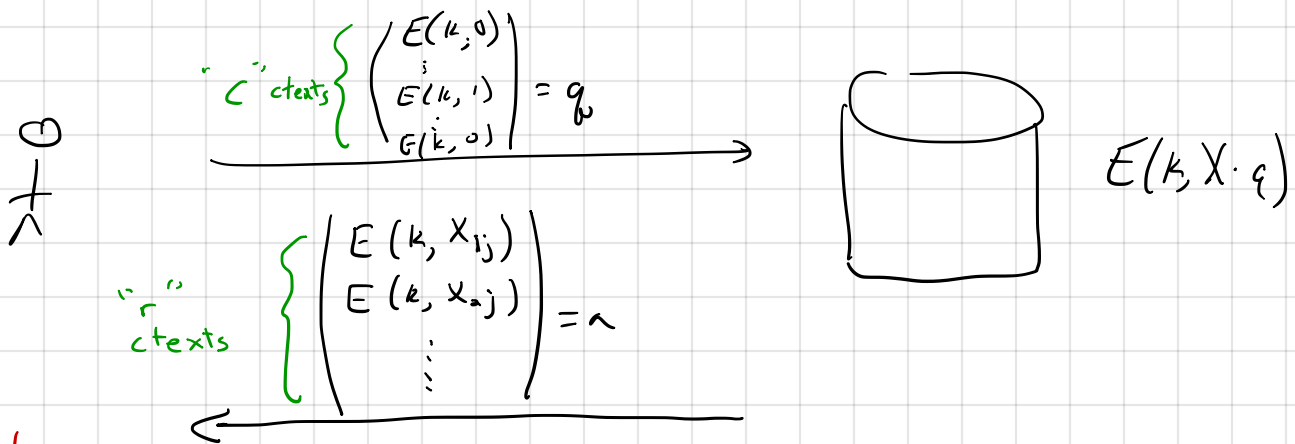
↳ Decrypt to recover $\vec{x}_j$ (jth column of X)
   ↳ Get bit $X_{ij}$

# Single-Server PIR: Reducing communication

Lets look more closely at our PIR scheme

$$X = \begin{array}{c} c \\ r \begin{pmatrix} & & \\ & & \\ & & \end{pmatrix} \end{array}$$

"$c$" ctexts $\left\{ \begin{pmatrix} E(k, 0) \\ \vdots \\ E(k, 1) \\ \vdots \\ E(k, 0) \end{pmatrix} = q \right.$

$E(k, X \cdot q)$

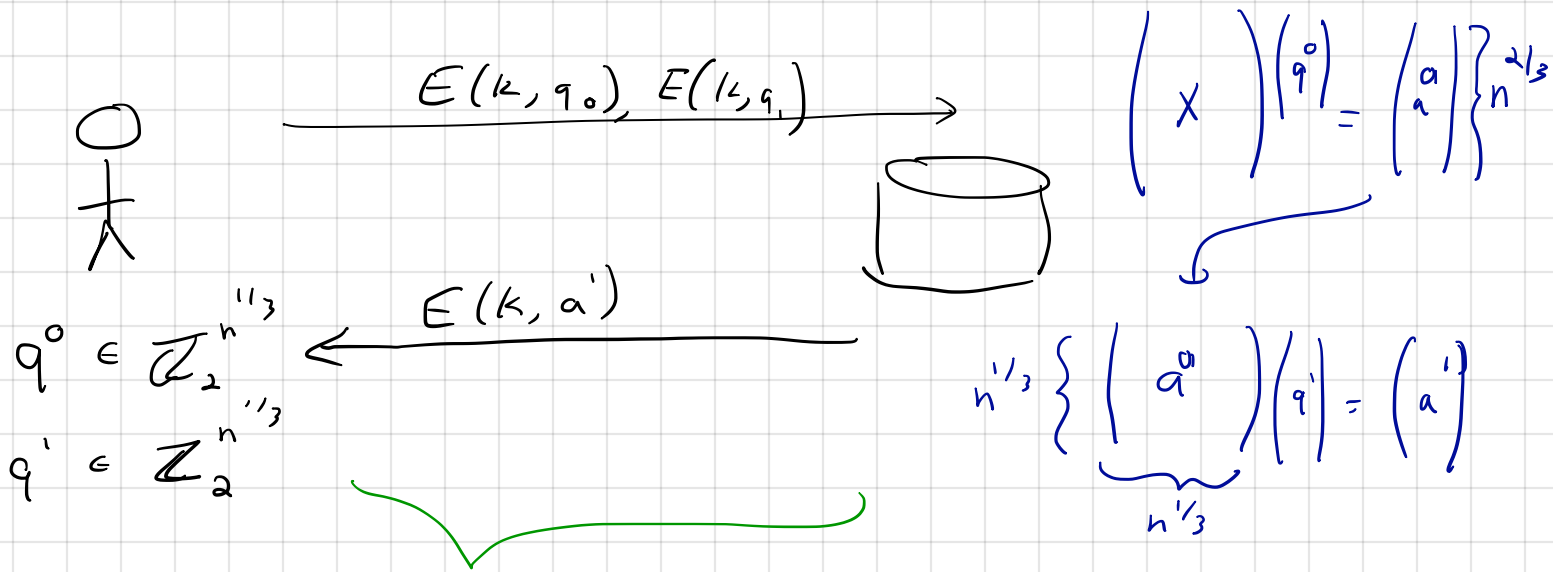"$r$" ctexts $\left\{ \begin{pmatrix} E(k, X_{1j}) \\ E(k, X_{2j}) \\ \vdots \\ \vdots \end{pmatrix} = a \right.$

Client discards all but one of responses!

Idea: View answer $a$ as a database. Apply single-server PIR recursively to fetch $i$-th element of answer!

$$X = \begin{pmatrix} & n^{1/3} \\ & \end{pmatrix} n^{2/3}$$

$$E(k, q_0), \ E(k, q_1) \longrightarrow$$

$$\begin{pmatrix} X \end{pmatrix}\begin{pmatrix} q^0 \end{pmatrix} = \begin{pmatrix} a \\ a \end{pmatrix} \Big\} n^{2/3}$$

$$\xleftarrow{\quad E(k, a') \quad}$$

$$q^0 \in \mathbb{Z}_2^{n^{1/3}}$$

$$q^1 \in \mathbb{Z}_2^{n^{1/3}}$$

$$n^{1/3} \Bigg\{ \underbrace{\begin{pmatrix} a^0 \end{pmatrix}}_{n^{1/3}}\begin{pmatrix} q^1 \end{pmatrix} = \begin{pmatrix} a' \end{pmatrix}$$

Total communication is $O(n^{1/3})$!

OR takes $2^{\tilde{\Omega}(n^{1/3})}$ time

Under reasonable assumptions can continue recursion to get complexity $2^{O(\sqrt{\log n} \ \log \log n)}$

Under slightly crazier assumptions ($\Phi$-hiding), can get polylog($n$) communication... (See Ostrovsky & Skeith survey).

# Extensions

- PIR by keyword
- PIR writing