

Lecture 3: Private Information Retrieval

MIT- 6.893

Fall 2020

Henry Corrigan Gibbs

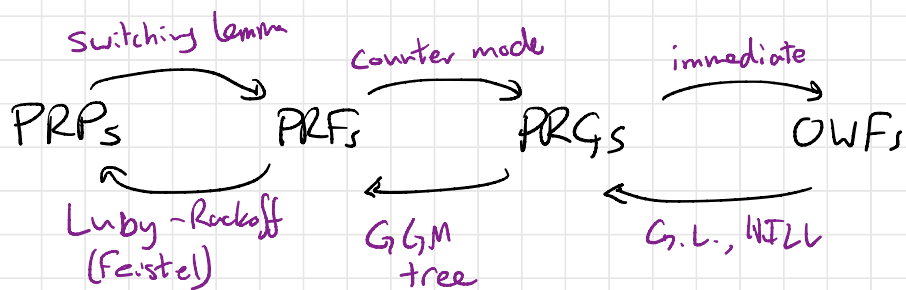
Plan

- * Recap: Preliminaries
- * PIR: What it is, why it's amazing
- * Stretch break
- * Constructions
 - Two-server PIR
 - One-server PIR

Logistics

- * HW1 due this Friday 9/18 @ 5pm Boston
via Gradescope
↳ You must use Latex template
- * OH: W 3-4:30pm on Zoom (link on Piazza)
- * Please give feedback on psets
- * Anonymous feedback form

Recap: Fundamental Primitives



→ All imply each other.
→ None imply key exchange.

- If $P = NP$, none exist.
- If $P \neq NP$, ?

Asymptotic view:

security parameter λ

efficient = $\text{poly}(\lambda)$

small = $\text{negl}(\lambda)$

Concrete view

efficient = runs on your computer in reasonable time budget

small = 2^{-128} .

A "perfect" research result

← Adopted from talk by Dan Spielman, who adopted it from someone else (I think)

- 1) Has a beautiful theory
- 2) Works in practice
- 3) Solves a problem that people care about.

⇒ It's a rare piece of work that meets this rubric. But aim high.

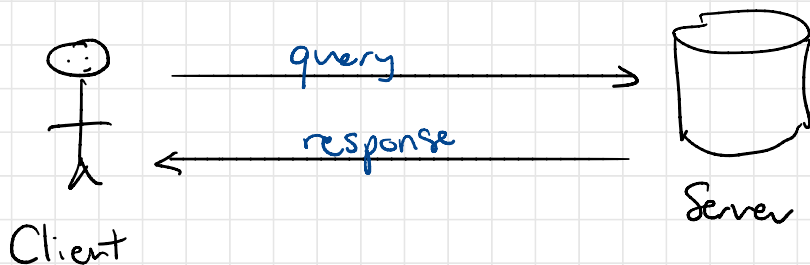
Today

- One of my favorite "almost perfect" ideas in crypto.
- Lots of activity, more ongoing even today.
 - ↳ Will cover recent results next week.
- A classic crypto result: seems impossible, then is simple.

Bad news: For reasons we'll see, it's not quite practical yet.

Private Information Retrieval

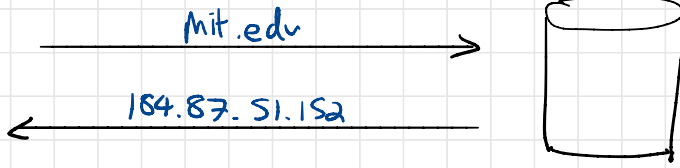
Every day on the Internet



Examples

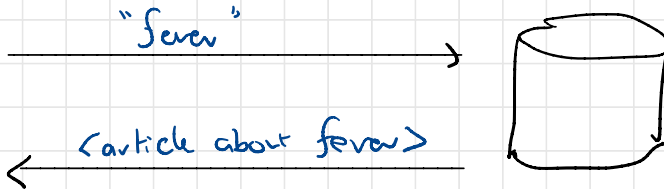
①

DNS



TLD name server
for .edu

②



WebMD

③

Many more....

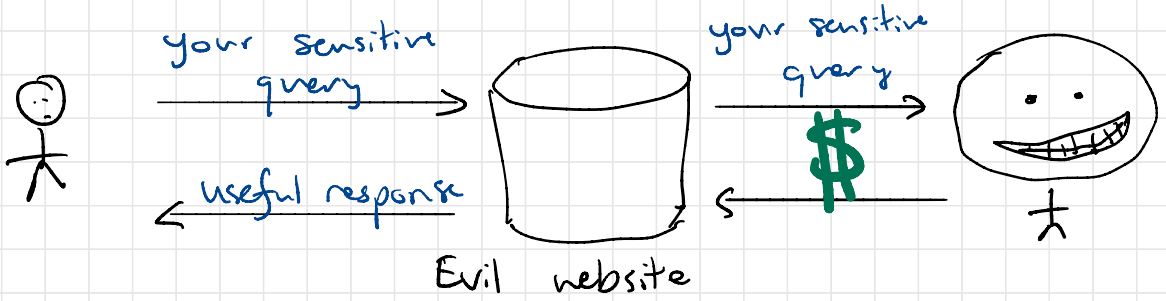
- Searching Google
- Looking news articles
- Fetching data from social media networks
- Looking for Airbnb properties

Nota: The client's query can be sensitive!

It can leak:

- what website you're visiting
- your health conditions
- your travel plans
- political interests
- ...

↳ Today, server learns all of these things!



Question:

"Can you query a database without the database learning you query?"

Trivial answer:

"Yes, just download the entire database."

→ DB server doesn't learn your query.

Still, this is unsatisfying.

Let's ask a better question....

Better Question:

"Can you query a database without the database learning you query..."

... With communication sublinear in the database size?

Answer: Unconditionally, no. [CGKS '95]

We won't prove this, but there is a clean info. theoretic argument in the original PIR paper.

What do we do when we are stuck in life?

Option I

Change the model.
(ROM, etc.)

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



"two-server PIR"

[CGKS '95]

Option II

Make assumptions!

Under basic "public-key" assumptions (DDH, Paillier, ...) can build non-trivial PIR w/ 1 server

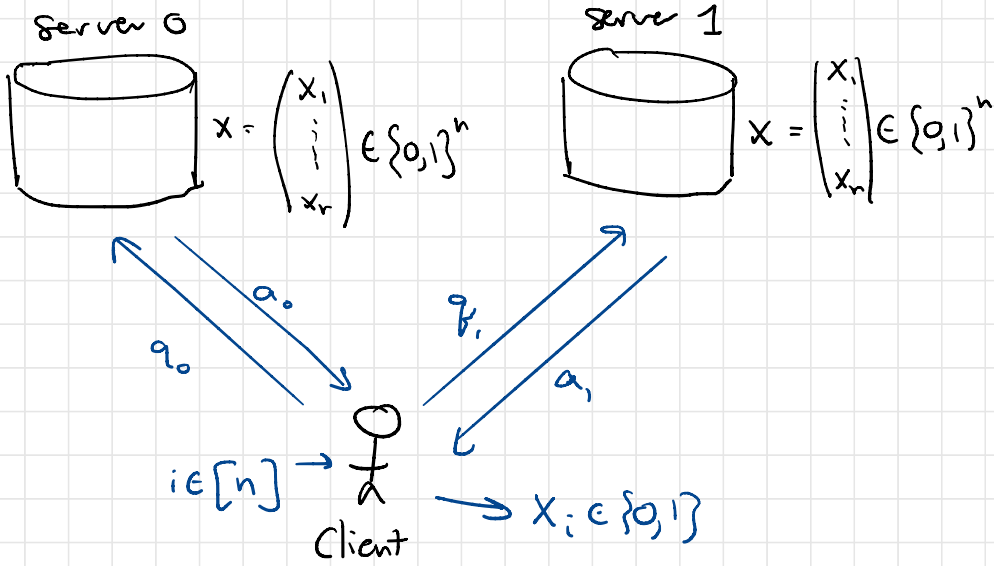
"single-server PIR"

[Kushilevitz & Ostrovsky '97]

↖ (k-server PIR uses k DB replicas)

Two-server PIR

Model



Important! Security only holds if servers do not collude (i.e., one of the two servers is honest).

Non-essential simplifications

* DB is an array of bits (can extend to handle longer rows)

* DB lookup is by index (can implement a key-value map)

later this week

More formally:

Two-server PIR consists of three eff algs:

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

$$\text{Answer}(x, q_\beta) \rightarrow a_\beta$$

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

Properties

① Correctness: Client gets the bit it wants.

$$\forall n \in \mathbb{N}, \forall i \in [n], \forall x \in \{0,1\}^n$$

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

② No single server learns anything about client's bit.

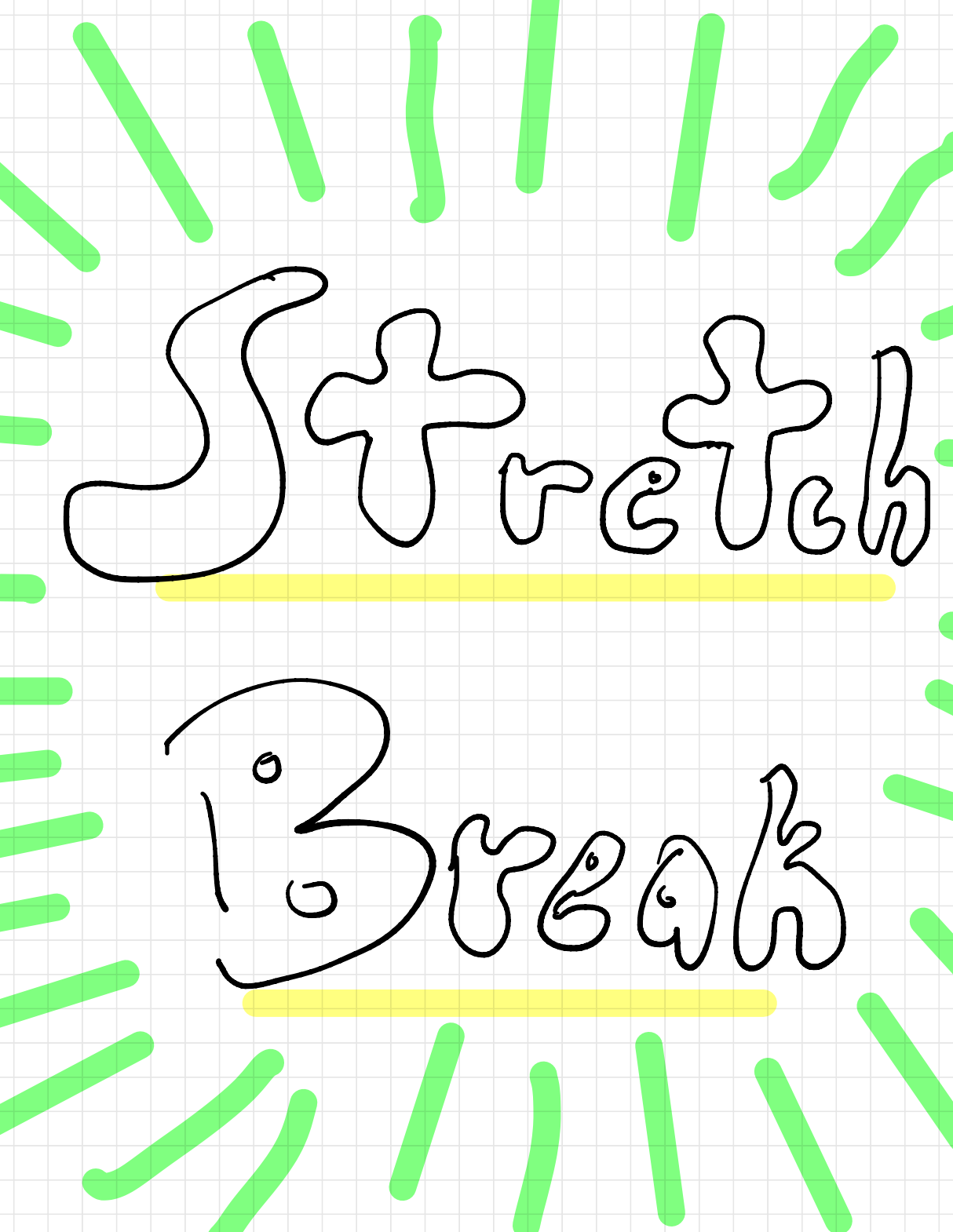
$$\forall n \in \mathbb{N} \quad \forall i, i' \in [n], \quad \forall \beta \in \{0,1\}$$

Can be \equiv for perfect privacy.

$$\{q_\beta : (q_0, q_1) \leftarrow \text{Query}(1^n, i)\} \stackrel{c}{\approx} \{q_\beta : (q_0, q_1) \leftarrow \text{Query}(1^n, i')\}$$

Non-collusion is captured by our requirement that the marginal distributions are indist.

→ In info-theoretic setting (explain), (q_0, q_1) will leak secret index i .



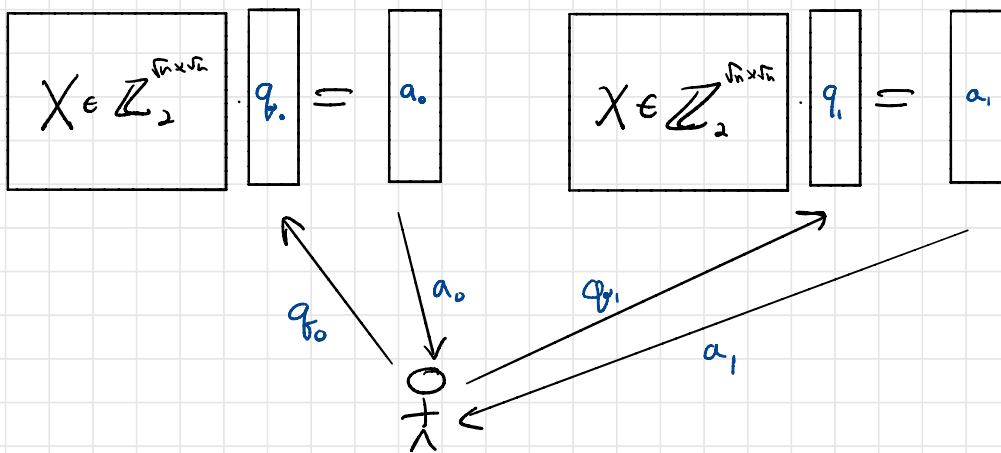
Stretch

Break

Two-server PIR scheme with $O(\sqrt{n})$ communication.

Idea: View Database
a matrix in $\mathbb{Z}_2^{\sqrt{n} \times \sqrt{n}}$.

↑ Already very non-trivial.



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

Query $(1^n, i, j) \rightarrow (q_0, q_1)$

Sample random $q_0, q_1 \in \mathbb{Z}_2^{\sqrt{n}}$ st.

$$q_0 + q_1 = e_i \in \mathbb{Z}_2^{\sqrt{n}}$$

Answer $(x, q) \rightarrow X \cdot q \in \mathbb{Z}_2^{\sqrt{n}}$

Reconstruct $(a_0, a_1) \rightarrow (a_0 + a_1)_i = x_{ij}$

All-zeros with 1 in j -th position



↑ Select value from i -th position

① Correctness.

$$\begin{aligned} a_0 + a_1 &= (Xq_0 + Xq_1)_i \\ &= (X(q_0 + q_1))_i \\ &= (Xe_j)_i \\ &= x_{ij} \end{aligned}$$

② Security

q_0 is a uniform random vector
(independent of $(:;j)$).

↳ Same for q_1 .

Notice: No computation assumptions here!

Efficiency: Upload: $2\sqrt{n}$ bits

Download: $4\sqrt{n}$ bits.

⇒ Total: $4\sqrt{n}$ bits.

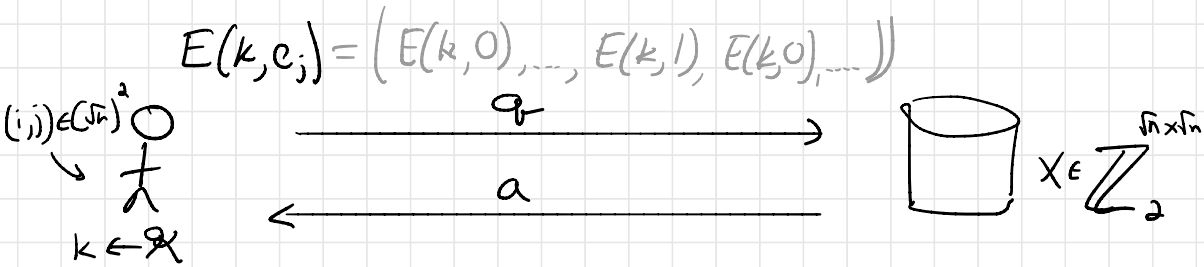
Single-server PIR

Linearly homomorphic encryption scheme:

$$E(k, m_0) + E(k, m_1) = E(k, \underbrace{m_0 + m_1}_{\text{mod } 2})$$

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

Idea: Client sends encryption of its query vector rather than using secret sharing.



Output $x_{i,j} \leftarrow (\text{Dec}(k, a))_i$

$$= \text{Dec}(k, \text{Enc}(k, \vec{x}_j))$$
$$= x_{i,j}$$

$$\begin{aligned} a &\leftarrow X \cdot q \\ &= X \cdot \text{Enc}(k, e_j) \\ &= \text{Enc}(k, X \cdot e_j) \\ &= \text{Enc}(k, \vec{x}_j) \end{aligned}$$

By adding ciphertexts

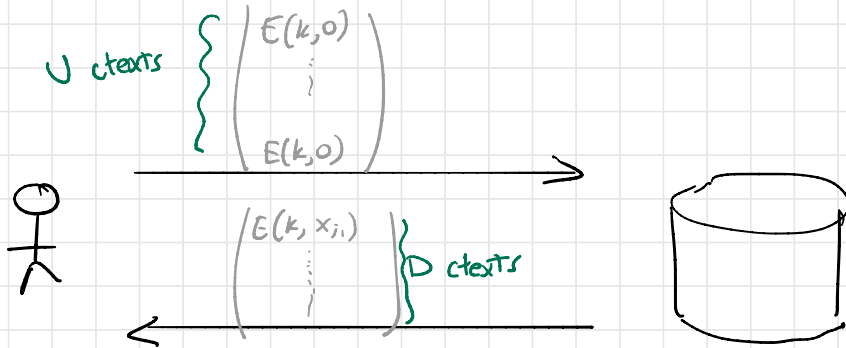
Communication: $2\sqrt{n}$ ciphertexts.

Reducing the Communication

Kushilevitz &
Ostrovsky '97

Let's look more closely at our PIR scheme...

$$\bar{X} = D \begin{pmatrix} u \\ \vdots \\ u \end{pmatrix}$$



Client throws away all
but one of the responses!

Idea: View answer to query as
another database and run
a second PIR on this DB.

The only catch is that each step of the recursion

n bits $\longrightarrow \sqrt{n}$ ciphertexts.

Under "reasonable" assumptions (OR takes $2^{\sqrt{n^{1/2}}}$ time),
get $2^{O(\sqrt{\log n \log \log n})}$ communication.

↳ With more esoteric cryptosystems,
(Based on Phi-hiding, Damgård-Jurik), can
drop comm cost to $\text{polylog}(n)$.

State of the art in PIR

* Two-server PIR

Information theoretic - $n^{O(\sqrt{\log n} / \log n)} = n^{o(1)}$
(Dvir & Gopi 2015)

↳ Do better schemes exist? With $O(\log^2 n)$ comm.?

Computational - $O(1 \log n)$ ← Boyle speaking here next week!
(Boyle, Gilman, Ishai 2015)

↳ Requires only PRGs. Concretely quite efficient.

* Single-server PIR

polylog(n) communication - from QR, DDH, LWE, ...

(Camenisch, Micali, Staller '99, Lipmaa '05, ...)

Computational Efficiency in PIR

In all above schemes, the servers run in time **linear** in the DB size.

⇒ Linear DB Scan per query.

For certain "natural" PIR schemes, this limitation is inherent. (Beimel, Ishai, Malkin '04)

→ We will see some ways ←
to reduce the computation
cost at the servers.