


# CS 350S: Privacy-Preserving Systems

## Private Information Retrieval III

# Recap: two-server PIR scheme

Write database as  $\sqrt{n} \times \sqrt{n}$  matrix




|          |          |          |          |
|----------|----------|----------|----------|
| $x_1$    | $x_2$    | $x_3$    | $x_4$    |
| $x_5$    | $x_6$    | $x_7$    | $x_8$    |
| $x_9$    | $x_{10}$ | $x_{11}$ | $x_{12}$ |
| $x_{13}$ | $x_{14}$ | $x_{15}$ | $x_{16}$ |



|          |          |          |          |
|----------|----------|----------|----------|
| $x_1$    | $x_2$    | $x_3$    | $x_4$    |
| $x_5$    | $x_6$    | $x_7$    | $x_8$    |
| $x_9$    | $x_{10}$ | $x_{11}$ | $x_{12}$ |
| $x_{13}$ | $x_{14}$ | $x_{15}$ | $x_{16}$ |

# Recap: two-server PIR scheme

Write database as  $\sqrt{n} \times \sqrt{n}$  matrix




$$\begin{matrix} x_1 & x_2 & x_3 & x_4 \\ x_5 & x_6 & x_7 & x_8 \\ x_9 & x_{10} & x_{11} & x_{12} \\ x_{13} & x_{14} & x_{15} & x_{16} \end{matrix}$$




$$\begin{matrix} x_1 & x_2 & x_3 & x_4 \\ x_5 & x_6 & x_7 & x_8 \\ x_9 & x_{10} & x_{11} & x_{12} \\ x_{13} & x_{14} & x_{15} & x_{16} \end{matrix}$$


# Recap: two-server PIR scheme




# Recap: two-server PIR scheme



# Recap: two-server PIR scheme



# Recap: two-server PIR scheme



Write database as  $\sqrt{n} \times \sqrt{n}$  matrix

$$\begin{matrix} x_1 & x_2 & x_3 & x_4 \\ x_5 & x_6 & x_7 & x_8 \\ x_9 & x_{10} & x_{11} & x_{12} \\ x_{13} & x_{14} & x_{15} & x_{16} \end{matrix} \times \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} x_2 \\ x_6 \\ x_{10} \\ x_{14} \end{bmatrix}$$

$$\begin{matrix} x_1 & x_2 & x_3 & x_4 \\ x_5 & x_6 & x_7 & x_8 \\ x_9 & x_{10} & x_{11} & x_{12} \\ x_{13} & x_{14} & x_{15} & x_{16} \end{matrix} \times \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} x_2 \\ x_6 \\ x_{10} \\ x_{14} \end{bmatrix}$$


Communication:  $O(\sqrt{n})$

# Recap: additively homomorphic encryption

Encryption scheme that supports:

- Adding ciphertexts:  $\text{Enc}_k(x) + \text{Enc}_k(y) = \text{Enc}_k(x + y)$
- Multiplication by a constant (by extension):  $c \cdot \text{Enc}_k(x) = \text{Enc}_k(c \cdot x)$


# Recap: single-server PIR scheme



Write database as  $\sqrt{n} \times \sqrt{n}$  matrix

$$\begin{matrix} x_1 & x_2 & x_3 & x_4 \\ x_5 & x_6 & x_7 & x_8 \\ x_9 & x_{10} & x_{11} & x_{12} \\ x_{13} & x_{14} & x_{15} & x_{16} \end{matrix}$$


# Recap: single-server PIR scheme




Write database as  $\sqrt{n} \times \sqrt{n}$  matrix

$$\begin{matrix} x_1 & x_2 & x_3 & x_4 \\ x_5 & x_6 & x_7 & x_8 \\ x_9 & x_{10} & x_{11} & x_{12} \\ x_{13} & x_{14} & x_{15} & x_{16} \end{matrix}$$


# Recap: single-server PIR scheme




# Recap: single-server PIR scheme



# Recap: single-server PIR scheme



# Recap: single-server PIR scheme



# Outline

- 1. SimplePIR**
2. Tiptoe
3. Student presentation: Compass

# Learning with errors assumption

[Regev09]

Parameterized by  $n, m, q$  and error distribution  $\chi$

$$\left( \begin{array}{c} A \\ , \\ A \end{array} \cdot \begin{array}{c} s \\ + \\ e \end{array} \right) \approx \left( \begin{array}{c} A \\ , \\ r \end{array} \right)$$

The diagram illustrates the Learning with Errors (LWE) assumption. It shows two matrices, \$A\$ (represented by a grid), being multiplied by a vector \$s\$ (represented by a blue rectangle) and added to a vector \$e\$ (represented by an orange rectangle). The result is approximately equal to a matrix \$A\$ followed by a random vector \$r\$ (represented by an orange rectangle).

where  $A \in \mathbb{Z}_q^{m \times n}$ ,  $s \in \mathbb{Z}_q^n$ ,  $e \in \chi^m$ , and  $r \in \mathbb{Z}_q^m$  is a random vector

Choosing error distribution: balance between utility and security

- $\chi$  is uniform distribution: secure but not useful
- $\chi$  always outputs 0: useful but not secure
- Standard setting:  $\chi$  is discrete Gaussian distribution restricted to some output range

# Regev encryption

[Regev09]

Parameterized by  $n, m$ , ciphertext modulus  $q$ , plaintext modulus  $p$ , and error distribution  $\chi$

Encrypt message  
 $m \in \mathbb{Z}_p$  with  
secret key  $s$

$$\left( \begin{array}{c} a \\ a^T \\ s \\ e \\ m \end{array} \right) \left( \begin{array}{c} a \\ a^T \\ s \\ e \\ m \end{array} \right) = (a, c)$$

LWE assumption

$$\left( \begin{array}{c} A \\ A \\ s \\ e \end{array} \right) \approx \left( \begin{array}{c} A \\ r \end{array} \right)$$

# Regev encryption

[Regev09]

Parameterized by  $n, m, q$ , plaintext modulus  $p$ , and error distribution  $\chi$

Encrypt message  
 $m \in \mathbb{Z}_p$  with  
secret key  $s$

$$\left( \begin{array}{c} a \\ a^T \\ s \end{array}, \quad a^T \cdot s + e + m \lfloor q/p \rfloor \end{array} \right) = (a, c)$$


Decrypt ciphertext  
 $(a, c)$  with secret  
key  $s$

Round  $c - a^T \cdot s$  to the nearest multiple of  $\lfloor q/p \rfloor$


Succeeds as long as error is “small”, i.e.,  $< \frac{1}{2} \cdot \lfloor q/p \rfloor$

Feature: additively homomorphic –  $\text{Enc}_k(m_1) + \text{Enc}_k(m_2) = \text{Enc}_k(m_1 + m_2)$

# Recall: single-server PIR scheme



# Regev encryption



Message is a vector of values


# Processing a PIR query

$$D \times \left( \begin{array}{c} \theta \\ \theta \\ \theta \\ \theta \\ \theta \end{array} \right) = \left( \begin{array}{c} D \times \text{Random matrix} \\ A \end{array} , D \times \text{Message dependent} \\ b \end{array} \right)$$

The diagram illustrates the processing of a PIR query. On the left, a vector of length  $D$  is shown, consisting of five elements, each marked with a yellow padlock icon, indicating they are encrypted. This vector is multiplied by a  $D \times 5$  matrix. The result is split into two parts: a  $D \times D$  matrix multiplication with a "Random matrix" (labeled "A") and a  $D \times 1$  vector multiplication with a "Message dependent" vector (labeled "b").


# SimplePIR

[Henzinger, Hong, Corrigan-Gibbs, Meiklejohn, Vakuntanathan]




# SimplePIR

Preprocessing:  
Download hint



Online:  
make query



Database **D**

Database **D**

Can re-use hint to  
make many queries

High-throughput queries


# Outline

1. SimplePIR
2. **Tiptoe**
3. Student presentation: Compass

# Web-search queries reveal your sensitive data

|               |   |
|---------------|---|
| Health        | ballet knee problem                     |
| Finances      | job opportunities in west palm beach    |
| Religion      | african american churches in norfolk va |
| Citizenship   | application forms us citizen            |
| Relationships | domestic violence laws in new york city |


# Today: Send query to search engine




# Today: Send query to search engine




# Goal: Search without revealing the query



# Goal: Search without revealing the query



# Goal: Search without revealing the query



- Limitations:**
- does not hide subsequent HTTP(S) requests
  - does not hide *when* the client makes searches
  - does not guarantee **correct** search results

# Challenges



1. Search algorithms are not “crypto friendly”
  - Crypto: Boolean circuits
  - Search: Random access to terabytes of data ☹
2. Computing on encrypted data is expensive
  - Orders of magnitude worse than computation on plaintext
3. The web is large

# Tiptoe: A private web-search engine


[Henzinger, Dauterman, Corrigan-Gibbs, Zeldovich]

- + Server learns no information about client's query  
unlike DuckDuckGo, Google over Tor, ...
- + Searches 364 million web pages with 2.7 seconds of latency  
with 145 core-s of compute, 57 MiB comm., and 300 MiB of client storage
- + Supports image search; extensible to code, audio, video
- Search quality cannot compete with non-private search  
evaluated with standard information-retrieval benchmark (MS-MARCO)


# Architecture



# Architecture




# Architecture



# Tiptoe: Design ideas

1. Reduce private text search to private nearest-neighbor search  
*Key tool:* Semantic embeddings [Osgood57, ...]
2. Reduce private nearest-neighbor search to private matrix multiplication  
*Key tool:* Clustering to reduce communication to  $\sqrt{N}$  (from  $\sqrt{N}$ )
3. Use a (new) fast scheme for private matrix multiplication  
*Key tool:* Fast linearly homomorphic encryption [HHCMV'22]

# 1. Reduce private text search to private nearest-neighbor search using semantic embeddings [Osgood57, ..., MCCD13, BCLT19...]



# Why embeddings are useful

- Reduce a messy problem to a clean one (text search)  
(nearest-neighbor search)
- Can improve/optimize embeddings independently of crypto
- There are open embeddings for text, code, images, videos, audio...
- The embedding has no effect on privacy

# Tiptoe: Design ideas



Reduce private text search to private nearest-neighbor search

Key tool: Semantic embeddings [Osgood57, ...]

2. Reduce private nearest-neighbor search to private matrix multiplication



Key tool: Clustering to reduce communication to  $\sqrt{N}$  (from  $\sqrt{N}$ )

3. Use a (new) fast scheme for private matrix multiplication

Key tool: Fast linearly homomorphic encryption [HHCMV'22]

## 2. Reduce nearest-neighbor search to private matrix multiplication

### Preprocessing batch job




Embedding  
space

## 2. Reduce nearest-neighbor search to private matrix multiplication


### Preprocessing batch job

Cluster  
centroids  
(20 MB)




## 2. Reduce nearest-neighbor search to private matrix multiplication

At query time: Identify “best” cluster




## 2. Reduce nearest-neighbor search to private matrix multiplication

At query time: Fetch scores for docs in “best” cluster



## 2. Reduce nearest-neighbor search to private matrix multiplication

At query time: Fetch scores for docs in “best” cluster




## 2. Reduce nearest-neighbor search to private matrix multiplication

At query time: Fetch scores for docs in “best” cluster

With  $C$  clusters and  
length- $L$  embeddings,  
**Upload:**  $O(C \cdot L)$  bytes

With  $N$  total documents,  
**Download:**  $O(N/C)$  bytes

Balancing upload & download,  
**Total comm.:**  $O(\sqrt{NL})$  bytes




# Tiptoe: Design ideas


-  Reduce private text search to private nearest-neighbor search
  - Key tool: Semantic embeddings [Osgood57, ...]
-  Reduce private nearest-neighbor search to private matrix multiplication
  - Key tool: Clustering to reduce communication to  $\sqrt{N}$  (from  $\sqrt{N}$ )
- 3. Use a (new) fast scheme for private matrix multiplication
  - Key tool: Fast linearly homomorphic encryption [HHCMV'22]

# SimplePIR

Preprocessing:  
Download hint



Online:  
make query



Can re-use hint to  
make many queries

# Avoid storing the hint at the client

Idea: outsource the hint to the server

Most of the decryption algorithm just requires linear operations

- Can perform under encryption at the server

# Tiptoe: Design ideas



Reduce private text search to private nearest-neighbor search

**Key tool:** Semantic embeddings [Osgood57, ...]



2. Reduce private nearest-neighbor search to private matrix multiplication

**Key tool:** Clustering to reduce communication to  $\sqrt{N}$  (from  $\sqrt{N}$ )




Use a (new) fast scheme for private matrix multiplication

**Key tool:** Fast linearly homomorphic encryption [HCMV'22]


# Architecture recap

## Batch job

- Embeds documents
- Clusters embeddings
- Computes the encryption-scheme “hint” for both services



# Architecture recap



# Strengths and weaknesses of Tiptoe from the evaluation


## Strengths

- Evaluates on real webpages (364M from common crawl) and images (400M from LAION-400M)
- Similar MRR@100 score to Coeus, but much better performance
- Costs orders of magnitude lower than client-side index or Coeus
- Clustering makes it possible for communication to scale sub linearly with the number of documents with a comparatively small drop in search quality
- Infrastructure is compatible with text and images

## Weaknesses

- Search quality worse than that of state-of-the-art plaintext search algorithms
- Querying more clusters could help improve search quality, but would substantially increase costs
- Evaluation over hundreds of millions of documents, but web is much larger (shows analytical scaling to tens of billions of documents)
- Compute costs scale linearly with number of documents
- Requires client to store state (embedding function + cluster centroids)

# Architecture recap



## Client download

- Embedding function 270 MB
- Cluster centroids 20 MB

# Implementation

6,500 lines of code (Go and Python)


+ our external private information retrieval library (SimplePIR, 2300 lines)

Search quality: MS Marco doc-rank “dev” dataset, 3.2M pages

Text search: Cleaned common crawl, 364M pages  
msmarco-distilbert-base-tas-b embedding

Image search: LAION-400M data set, 400M images  
CLIP embedding (prepackaged)

# Implementation



# Tiptoe is cheaper than state-of-the-art private search

|                                      | Coeus (SOSP 21) | Tiptoe           | Gain    |
|--------------------------------------|-----------------|------------------|---------|
| Docs searched                        | 5 million       | 364 million      | 72 ×    |
| Client storage                       | -               | 0.3 GiB          | – ∞ ×   |
| Server compute<br>(per million docs) | 2,580 core-s    | 0.4 core-s       | 6,450 × |
| Communication<br>(per million docs)  | 10 MiB          | 0.15 MiB         | 66 ×    |
| AWS cost<br>(per million docs)       | ≤ 1 US cent     | ≤ 0.0008 US cent | 1,250 × |

# Tiptoe is cheaper than state-of-the-art private search

|                                      | Coeus (SOSP 21)   | Tiptoe           | Gain             |
|--------------------------------------|---|------------------|------------------|
| Docs searched                        | 5 million   | 364 million      | 72 ×             |
| Client storage                       | Semantic embeddings are 100 × smaller than a bag-of-words representation. |                  | $-\infty \times$ |
| Server compute<br>(per million docs) | 2,580 core-s  | 0.4 core-s       | 6,450 ×          |
| Communication<br>(per million docs)  | 10 MiB  | 0.15 MiB         | 66 ×             |
| AWS cost<br>(per million docs)       | ≤ 1 US cent   | ≤ 0.0008 US cent | 1,250 ×          |

# Tiptoe is cheaper than state-of-the-art private search

|                                      | Coeus (SOSP 21)   | Tiptoe           | Gain        |
|--------------------------------------|---|------------------|-------------|
| Docs searched                        | 5 million   | 364 million      | 72 ×        |
| Client storage                       | Semantic embeddings are 100 × smaller than a bag-of-words representation.   |                  | $-\infty$ × |
| Server compute<br>(per million docs) | 2.580 core-s  | 0.4 core-s       | 6,450 ×     |
| Communication<br>(per million docs)  | Tiptoe's high-throughput crypto tools are over 10 × faster than prior work. |                  | 66 ×        |
| AWS cost<br>(per million docs)       | ≤ 1 US cent   | ≤ 0.0008 US cent | 1,250 ×     |

# Tiptoe is cheaper than state-of-the-art private search

|                                      | Coeus (SOSP 21)   | Tiptoe      | Gain        |
|--------------------------------------|---|-------------|-------------|
| Docs searched                        | 5 million   | 364 million | 72 ×        |
| Client storage                       | Semantic embeddings are 100 × smaller than a bag-of-words representation.       |             | $-\infty$ × |
| Server compute<br>(per million docs) | 2.580 core-s  | 0.4 core-s  | 6,450 ×     |
| Communication<br>(per million docs)  | Tiptoe's high-throughput crypto tools are over 10 × faster than prior work.     |             | 66 ×        |
| AWS cost<br>(per million docs)       | Clustering lets Tiptoe's communication scale sub-linearly with the corpus size. |             | 1,250 ×     |

# Tiptoe's search executes in seconds

|                    | Text search | Image search |
|--------------------|-------------|--------------|
| Docs searched      | 364 million | 400 million  |
| System config      | 208 vCPUs   | 384 vCPUs    |
| End-to-end latency | 2.7 s       | 3.5 s        |

# Tiptoe's search executes in seconds

|                      | Text search                      | Image search     |                  |
|----------------------|----------------------------------|------------------|------------------|
| Docs searched        | 364 million                      | 400 million      |                  |
| System config        | 208 vCPUs                        | 384 vCPUs        |                  |
| End-to-end latency   | 2.7 s                            | 3.5 s            |                  |
| Comm.                | { ahead of time<br>at query time | 42 MiB<br>15 MiB | 50 MiB<br>21 MiB |
| Client preprocessing | 34 s                             | 35 s             |                  |
| Server compute       | 145 core-s                       | 339 core-s       |                  |

# Tiptoe's search executes in seconds

|                      | Text search                      | Image search     |
|----------------------|----------------------------------|------------------|
| Docs searched        | 364 million                      | 400 million      |
| System config        | 208 vCPUs                        | 384 vCPUs        |
| End-to-end latency   | 2.7 s                            | 3.5 s            |
| Comm.                | { ahead of time<br>at query time | 42 MiB<br>15 MiB |
| Client preprocessing | 34 s                             | 35 s             |
| Server compute       | 145 core-s                       | 339 core-s       |

Extremely parallelizable, DRAM-bandwidth bound

# Tiptoe's search executes in seconds

|                      | Text search      | Image search     |
|----------------------|------------------|------------------|
| Docs searched        | 364 million      | 400 million      |
| System config        | 208 vCPUs        | 384 vCPUs        |
| End-to-end latency   | 2.7 s            | 3.5 s            |
| Comm.                | 42 MiB<br>15 MiB | 50 MiB<br>21 MiB |
| Client preprocessing | 34 s             | 35 s             |
| Server compute       | 145 core-s       | 339 core-s       |

Extremely parallelizable, DRAM-bandwidth bound


Off of the critical path

# Preprocessing

Run on heterogeneous GPU cluster + 80-core machine

|                                      | Text<br>364M docs | Image<br>400M docs |
|--------------------------------------|-------------------|--------------------|
| Embeddings                           | 92 corehours      | 583 corehours      |
| Build clusters                       | 927               | 1,493              |
| Dim. reduction +<br>cluster resizing | 312               | 290                |
| Crypto preproc.                      | 50                | 120                |
| Total time                           | 0.013 coresec/doc | 0.022 coresec/doc  |


# Search quality is “tolerable” (MS MACRO benchmark)



As embeddings improve, Tiptoe search quality will improve.

# Examples: Text search

how long before eagles get feathers



the meaning of haploid cell

The screenshot shows a Britannica article page. The header includes the Britannica logo, a search bar with "Search Britannica...", and navigation links for "Browse" and "Subscribe". The main title of the article is "haploid phase". Below the title, it says "biology". There is a "Share" button with a link. A section titled "LEARN ABOUT THIS TOPIC in these articles:" lists "algae". To the right, there is a thumbnail image of green algae and a text excerpt: "...of chromosomes and is called haploid, whereas in the second stage each cell has two sets of chromosomes and is called diploid. When one haploid gamete fuses with another haploid gamete during fertilization, the resulting combination, with two sets of chromosomes, is called a zygote. Either immediately or at some...". A "READ MORE" link is at the bottom right.

Exact string search is not as good...

# Examples: Image search

close up of a  
pizza sitting on  
top of a plate



zebra standing  
on top of a lush  
green field




shutterstock · 729088339

brown bear  
standing in  
front of a tree



# Apple's Enhanced Visual Search

- Apple's Enhanced Visual Search lets users see if their photos contain particular landmarks
- Uses design techniques drawn from Tiptoe
- For their use case, the cost of scanning over every cluster is too high
- Their solution: add differentially private noise to hide the identity of the queried cluster (can now query a sublinear number of clusters)



# Outline

1. SimplePIR
2. Tiptoe
3. **Student presentation: Compass**

# References

Asi, Hilal, Fabian Boemer, Nicholas Genise, Muhammad Haris Mughees, Tabitha Ogilvie, Rehan Rishi, Kunal Talwar et al. "Scalable private search with wally." *arXiv preprint arXiv:2406.06761* (2024).

Henzinger, Alexandra, Matthew M. Hong, Henry Corrigan-Gibbs, Sarah Meiklejohn, and Vinod Vaikuntanathan. "One server for the price of two: Simple and fast {Single-Server} private information retrieval." In *32nd USENIX Security Symposium (USENIX Security 23)*, pp. 3889-3905. 2023.

Henzinger, Alexandra, Emma Dauterman, Henry Corrigan-Gibbs, and Nickolai Zeldovich. "Private web search with Tiptoe." In *Proceedings of the 29th symposium on operating systems principles*, pp. 396-416. 2023.

Regev, Oded. "On lattices, learning with errors, random linear codes, and cryptography." *Journal of the ACM (JACM)* 56, no. 6 (2009): 1-40.

<https://65610.csail.mit.edu/2025/lec/l05-pke.pdf>

<https://machinelearning.apple.com/research/homomorphic-encryption>