

Private Approximate Nearest Neighbor Search with Sublinear Communication

Sacha Servan-Schreiber with Simon Langowski and Srini Devadas



Nearest Neighbor (NN) search



Nearest Neighbor (NN) search



Approximate Nearest Neighbor (ANN) search

(standard relaxation used in practice)







[1]: https://github.com/spotify/annoy













Two-party Computation

Two-party Computation

Fully-homomorphic Encryption

Two-party Computation

Fully-homomorphic Encryption

- Not lightweight:
 - **1 to 5 GB of communication** with databases of 1,000,000 items [1].

Two-party Computation

- Not lightweight:
 - **1 to 5 GB of communication** with databases of 1,000,000 items [1].
- No malicious security:
 - A **malicious** database can deviate from protocol and learn the query.
 - A **malicious** client can deviate from protocol and recover the database.

Fully-homomorphic Encryption

Two-party Computation

- Not lightweight:
 - **1 to 5 GB of communication** with databases of 1,000,000 items [1].
- **No** malicious security:
 - A **malicious** database can deviate from protocol and learn the query.
 - A **malicious** client can deviate from protocol and recover the database.

Fully-homomorphic Encryption

• Is lightweight for the client.

Two-party Computation

- Not lightweight:
 - **1 to 5 GB of communication** with databases of 1,000,000 items [1].
- **No** malicious security:
 - A **malicious** database can deviate from protocol and learn the query.
 - A **malicious** client can deviate from protocol and recover the database.

Fully-homomorphic Encryption

- Is lightweight for the client.
- Is trivially maliciously-secure.

Two-party Computation

- Not lightweight:
 - **1 to 5 GB of communication** with databases of 1,000,000 items [1].
- No malicious security:
 - A **malicious** database can deviate from protocol and learn the query.
 - A **malicious** client can deviate from protocol and recover the database.

Fully-homomorphic Encryption

- Is lightweight for the client.
- Is trivially maliciously-secure.
- **Not** lightweight for the database.
 - Takes **hours** with small databases of e.g., 500 to 2,000 items [2].

Two-party Computation

Fully-homomorphic Encryption

Our goals:

(1) low communication for the client,
(2) concrete efficiency for the database,
(3) privacy for the client and the database,
(4) and malicious security.

The setting: two non-colluding database servers.

Q: Why do we need non-colluding servers?

A: For efficient, symmetric-key cryptography only.¹

¹ See Appendix E of the full version of our paper [1] for a **single-server protocol** that is less concretely efficient but doesn't require any trust assumptions. [1]: Full version of our paper: https://eprint.iacr.org/2021/1157.pdf.



























Servers

- Server A
 - Hold replicas of the database.
 - Do not collude with clients or one another.





Servers

Server A

- Hold replicas of the database.
- Do not collude with clients or one another.





Clients

Will try to learn as much as possible about the database. May collude with other malicious clients.

Guarantees





Guarantees

• Accuracy if both servers follow the protocol.





Guarantees

- Accuracy if both servers follow the protocol.
- **User privacy** even if a server is malicious.





Guarantees

- Accuracy if both servers follow the protocol.
- **User privacy** even if a server is malicious.
- **Database privacy** even if a subset of clients are malicious.





Finding Nearest Neighbors

(non-privately, using Locality-Sensitive Hashing)

Finding the ANN using LSH (non-privately)



Finding the ANN using LSH (non-privately)



Finding the ANN using LSH (non-privately)




















Hashkey	Value(s)











Hashkey	Value(s)			
a3da901f	ID2: (1, 4)			
c26fab1d	ID100: (0,1)			
09ac34fd	ID3: (6,7)			
:	÷			
91ab3cd	ID11: (1,10)			

Step 1: Build an LSH hash table using h.Step 2: Query the hash table.



Hashkey	Value(s)				
a3da901f	ID2: (1, 4)				
c26fab1d	ID100: (0,1)				
09ac34fd	ID3: (6,7)				
:	÷				
91ab3cd	ID11: (1,10)				

Step 1: Build an LSH hash table using h.Step 2: Query the hash table.



Step 1: Build an LSH hash table using h.Step 2: Query the hash table.



Step 2: Query the hash table.

Step 3: Repeat with many different hash tables.

Query (2, 5)	Hashkey	Value(s)	Hashkey	Value(s)	Hashkey	Value(s)
	a3da901f	ID2: (1, 4)	beda11fe			
$ \rightarrow \Pi_1(\bullet) \frown$	c26fab1d	ID100: (0,1)	f12fbe10			
	09ac34fd	ID3: (6,7)	ac33445a			
	:	:	:	:	:	:
•	91ab3cd	ID11: (1,10)	91ab3cd			ID141
Data	Candidate	set: {IC)2: (1,4)			

Step 2: Query the hash table.

Step 3: Repeat with many different hash tables.



Step 2: Query the hash table.

Step 3: Repeat with many different hash tables.



Step 2: Query the hash table.

Step 3: Repeat with many different hash tables.

Step 4: Find closest neighbor in the candidate set.



Hashkey	Value(s)	Hashkey	Value(s)	Hashkey	Value(s)

Data

Candidate set: {ID2: (1,4), **ID5694**: (2,4), ID900: (3,4) }

One-time setup: Construct LSH tables



Step 1: Client uses LSH functions to find the hashkey



Step 1: Client uses LSH functions to find the hashkey



Step 1: Client uses LSH functions to find the hashkey



Q: How can the client retrieve the candidate from each table privately?

A: Using private information retrieval (PIR) [1].

• We use Distributed Point Functions [2] (DPFs) for efficiently querying hash tables in a two-server setting.

Q: How can the client retrieve the candidate from each table privately?

A: Using private information retrieval (PIR) [1].

• We use Distributed Point Functions [2] (DPFs) for efficiently querying hash tables in a two-server setting.

Q: How can the client retrieve the candidate from each table privately?

A: Using private information retrieval (PIR) [1].

• We use Distributed Point Functions [2] (DPFs) for efficiently querying hash tables in a two-server setting.



Q: How can the client retrieve the candidate from each table privately?

A: Using private information retrieval (PIR) [1].

• We use Distributed Point Functions [2] (DPFs) for efficiently querying hash tables in a two-server setting.



Q: How can the client retrieve the candidate from each table privately?

A: Using private information retrieval (PIR) [1].

• We use Distributed Point Functions [2] (DPFs) for efficiently querying hash tables in a two-server setting.



$$[v]_{A} + [v]_{B} = v$$
 (additive secret shares)

Step 1: Client uses LSH functions to find the hashkeyStep 2: Privately query hash tables using PIR



Step 1: Client uses LSH functions to find the hashkeyStep 2: Privately query hash tables using PIR



Step 1: Client uses LSH functions to find the hashkey
Step 2: Privately query hash tables using PIR



69

Step 1: Client uses LSH functions to find the hashkeyStep 2: Privately query hash tables using PIRStep 3: Client recovers the result



Query: (2,5)

Candidate set: $\{ID2: (1,4), ID5694: (2,4), ID900: (3,4), ID101: (6,5)\}$ (union over all tables)

Step 1: Client uses LSH functions to find the hashkeyStep 2: Privately query hash tables using PIRStep 3: Client recovers the result



Query: (2,5)

False positives

Candidate set: {**ID2:** (1,4), **ID5694:** (2,4), **ID900:** (3,4), **ID101:** (6,5) } (union over all tables)

Step 1: Client uses LSH functions to find the hashkeyStep 2: Privately query hash tables using PIRStep 3: Client recovers the result






A strawman protocol

Step 1: Client uses LSH functions to find the hashkeyStep 2: Privately query hash tables using PIRStep 3: Client recovers the result





```
Candidate set: {ID2: (1,4), ID5694: (2,4), ID900: (3,4), ID101: (6,5) }
(union over all tables)
```

Candidate set: {**ID2:** (1,4), **ID5694:** (2,4), **ID900:** (3,4), **ID101:** (6,5)}

Candidate set: {**ID2:** (1,4), **ID5694:** (2,4), **ID900:** (3,4), **ID101:** (6,5)}

I The candidate set leaks a lot about the database to the client!

Candidate set: {**ID2:** (1,4), **ID5694:** (2,4), **ID900:** (3,4), **ID101:** (6,5)}

I The candidate set leaks a lot about the database to the client!

The client learns:

• All near neighbors and their feature vectors.

Candidate set: $\{ID2: (1,4), ID5694: (2,4), ID900: (3,4), ID101: (6,5)\}$

The candidate set leaks a lot about the database to the client!

The client learns:

- All near neighbors and their feature vectors.
- Other feature vectors in the database (false positives).

Candidate set: $\{ID2: (1,4), ID5694: (2,4), ID900: (3,4), ID101: (6,5)\}$

The candidate set leaks a lot about the database to the client!

The client learns:

- All near neighbors and their feature vectors.
- Other feature vectors in the database (false positives).

Baseline leakage



```
Candidate set: \{ID2: (1,4), ID5694: (2,4), ID900: (3,4), ID101: (6,5)\}
```

The candidate set leaks a lot about the database to the client!

The client learns:

- All near neighbors and their feature vectors.
- Other feature vectors in the database (false positives).



Database Privacy with Radix sorting

(hide all feature vectors by pruning without comparisons)

Main idea: use radix sorting for comparison-free sorting

Main idea: use radix sorting for comparison-free sorting

(i.e., let the hash function do the work of sorting by distance)

Main idea: use radix sorting for comparison-free sorting

(i.e., let the hash function do the work of sorting by distance)

Table 1: radius = 0.1



Main idea: use radix sorting for comparison-free sorting

(i.e., let the hash function do the work of sorting by distance)

Table 1: radius = 0.1

Table 2: radius = 0.2



Main idea: use radix sorting for comparison-free sorting

(i.e., let the hash function do the work of sorting by distance)

Table 1: radius = 0.1

Table 2: radius = 0.2

Table 3: radius = 0.3



Main idea: use radix sorting for comparison-free sorting

(i.e., let the hash function do the work of sorting by distance)

 Table 1: radius = 0.1

 Table 2: radius = 0.2

 Table 3: radius = 0.3

 ...

 Table 20: radius = 2.0



Candidates are now sorted by distance from the query!

Candidate Set = $\{0, 0, 0, 1D5694, 1D2, 1D900, 0, 1D101\}$

Candidates are now sorted by distance from the query!

Candidate Set = $\{0, 0, 0, ID5694, ID2, ID900, 0, ID101\}$ Nearest Neighbor

Candidates are now sorted by distance from the query!

Candidate Set =
$$\{0, 0, 0, ID5694, ID2, ID900, 0, ID101\}$$

Nearest Neighbor

We no longer need to include the feature vectors (only the IDs)!

Candidates are now sorted by distance from the query!

Candidate Set =
$$\{0, 0, 0, ID5694, ID2, ID900, 0, ID101\}$$

Nearest Neighbor

We no longer need to include the feature vectors (only the IDs)!

Still leaks many IDs from the database to the client.

More database privacy with Oblivious masking

(hide all candidates except for the nearest neighbor)

After processing the PIR queries, each server has a share of the candidate set:



After processing the PIR queries, each server has a share of the candidate set:



Observation: the nearest neighbor is always preceded by shares of 0s...

$$\left[\text{Candidate Set}\right] = \left\{ [0], [0], [0], [1D5694], [1D2], [0], [1D900], [1D101] \right\}$$

Observation: the nearest neighbor is always preceded by shares of 0s...

$$[Candidate Set] = \left\{ [0], [0], [0], [ID5694], [ID2], [0], [ID900], [ID101] \right\}$$

Algorithm 1: ObliviousMaskingInput: Secret-shared vector $[v] \in \mathbb{F}_p^L$ and randomness rand.Output: Secret-shared vector $[y] = ([y_1], \dots, [y_L]) \in \mathbb{F}_p^L$.Procedure:1: for $i \in \{1, \dots, L\}$:1.1: Sample $r_i \leftarrow$ rand.1.2: Set $[y_i] \leftarrow [v_i] + r_i \cdot \left(\sum_{j=0}^{i-1} [v_j]\right)$.2: Output $y \in \mathbb{F}_p^L$.

Step 1: Servers agree on common randomness $(r_1, r_2, ..., r_n)$ e.g., with a short PRG seed.

$$\left[\text{Candidate Set} \right] = \left\{ [0], [0], [0], [ID5694], [ID2], [0], [ID900], [ID101] \right\}$$

Step 1: Servers agree on common randomness $(r_1, r_2, ..., r_n)$ e.g., with a short PRG seed.

Step 2: Servers copy (and randomize) shares from [Candidate Set] to [Candidate Set].

$$\begin{bmatrix} Candidate Set \end{bmatrix} = \begin{cases} [0], [0], [0], [ID5694], [ID2], [0], [ID900], [ID101] \\ \end{bmatrix}$$
$$\begin{bmatrix} Candidate Set \end{bmatrix} = \begin{cases} \end{cases}$$

Step 1: Servers agree on common randomness $(r_1, r_2, ..., r_n)$ e.g., with a short PRG seed.

Step 2: Servers copy (and randomize) shares from [Candidate Set] to [Candidate Set].

Step 1: Servers agree on common randomness $(r_1, r_2, ..., r_n)$ e.g., with a short PRG seed.

Step 2: Servers copy (and randomize) shares from [Candidate Set] to [Candidate Set].

$$\begin{bmatrix} Candidate Set \end{bmatrix} = \begin{cases} [0], [0], [0], [ID5694], [ID2], [0], [ID900], [ID101] \\ \end{bmatrix}$$
$$\begin{bmatrix} Candidate Set \end{bmatrix} = \begin{cases} [0] \end{cases}$$

Step 1: Servers agree on common randomness $(r_1, r_2, ..., r_n)$ e.g., with a short PRG seed.

Step 2: Servers copy (and randomize) shares from [Candidate Set] to [Candidate Set].

Step 1: Servers agree on common randomness $(r_1, r_2, ..., r_n)$ e.g., with a short PRG seed.

Step 2: Servers copy (and randomize) shares from [Candidate Set] to [Candidate Set].

Step 1: Servers agree on common randomness $(r_1, r_2, ..., r_n)$ e.g., with a short PRG seed.

Step 2: Servers copy (and randomize) shares from [Candidate Set] to [Candidate Set].

Step 1: Servers agree on common randomness $(r_1, r_2, ..., r_n)$ e.g., with a short PRG seed.

Step 2: Servers copy (and randomize) shares from [Candidate Set] to [Candidate Set].

Step 1: Servers agree on common randomness $(r_1, r_2, ..., r_n)$ e.g., with a short PRG seed.

Step 2: Servers copy (and randomize) shares from [Candidate Set] to [Candidate Set].

$$\begin{bmatrix} Candidate Set \end{bmatrix} = \begin{cases} [0], [0], [0], [ID5694], [ID2], [0], [ID900], [ID101] \\ \end{bmatrix} \\ \begin{bmatrix} Candidate Set \end{bmatrix} = \begin{cases} [0], [0], [0], [ID5694], [ID2] + [(0+0+0+ID5694)r_5] \\ \end{bmatrix} \end{cases}$$

Step 1: Servers agree on common randomness $(r_1, r_2, ..., r_n)$ e.g., with a short PRG seed.

Step 2: Servers copy (and randomize) shares from [Candidate Set] to [Candidate Set].

Step 1: Servers agree on common randomness $(r_1, r_2, ..., r_n)$ e.g., with a short PRG seed.

Step 2: Servers copy (and randomize) shares from [Candidate Set] to [Candidate Set].

$$\begin{bmatrix} Candidate Set \end{bmatrix} = \begin{cases} [0], [0], [0], [1D5694], [1D2], [0], [1D900], [1D101] \\ \hline Candidate Set \end{bmatrix} = \begin{cases} [0], [0], [0], [0], [1D5694], [\#$@!#*], [$%@!&], ... \end{cases}$$

All elements *after* the first non-zero element are **random**!

All elements *after* the first non-zero element are **random**!

Result: The client learns nothing beyond the first non-zero element (i.e., the ANN)!

$$\begin{bmatrix} Candidate Set \end{bmatrix} = \begin{cases} [0], [0], [0], [ID5694], [ID2], [0], [ID900], [ID101] \\ \end{bmatrix}$$

$$\begin{bmatrix} Candidate Set \end{bmatrix} = \begin{cases} [0], [0], [0], [0], [ID5694], \\ \end{bmatrix}$$
no leakage
All elements *after* the first non-zero element are **random**!

Result: The client learns \checkmark nothing beyond the first non-zero element (i.e., the ANN)!

The radix bucket of the ANN leaks the approximate distance from the query.



The radix bucket of the ANN leaks the approximate distance from the query.

• Asymptotically optimal leakage!



The radix bucket of the ANN leaks the approximate distance from the query.

- Asymptotically optimal leakage!
- In practice: 2 to 15x more than baseline leakage.



Evaluation

We evaluate on four real world datasets

- **Deep1B** (10 million items, 96 dimensions).
- **SIFT** (1 million items, 128 dimensions).
- **GIST** (1 million items, 916 dimensions).
- MNIST (60,000 items, 784 dimensions).

Efficiency on small datasets (MNIST; 60,000 items)



Server computation for > 95% accuracy (32 core servers):

• 300 milliseconds per query.

Communication: 800 KB between client and both servers.

Efficiency on small datasets (MNIST; 60,000 items)



Server computation for > 95% accuracy (32 core servers):

• 300 milliseconds per query.

Communication: 800 KB between client and both servers.

Efficiency on small datasets (MNIST; 60,000 items)



Server computation for > 95% accuracy (32 core servers):

• 300 milliseconds per query.

Communication: 800 KB between client and both servers.

Efficiency on large datasets (1M to 10M items)



Server computation for > 95% accuracy (32 core servers):

- **1.2 seconds per query** on 1M item datasets.
- **8 seconds per query** on 10M item datasets.

Communication: 1-2MB between client and both servers.

Efficiency on large datasets (1M to 10M items)



Server computation for > 95% accuracy (32 core servers):

Five orders of magnitude less computation compared to FHE-based approaches.

Efficiency on large datasets (1M to 10M items)



1,000 to 3,000X less communication compared to two-party computation!

Communication: 1-2MB between client and both servers.

Thank you!



Sacha Servan-Schreiber

with Simon Langowski and Srini Devadas

Full paper:ia.cr/2021/1157Code:github.com/sachaservan/private-annContact:3s@mit.edu, slangows@mit.edu