

Don't be Dense: Efficient Keyword PIR for Sparse Databases

Sarvar Patel, Joon Young Seo, Kevin Yeo

Published in USENIX

Michael Rodel, Nir Segal and Hallel Weinberg

this project is submitted as part of completing your duties in the Secure Computation class instructed by Dr. Adi Akavia.

Abstract

The paper "Don't be Dense: Efficient Keyword PIR for Sparse Databases" [1] introduces a novel keyword private information retrieval (PIR) scheme called SparsePIR, tailored specifically for sparse databases. In sparse databases, a significant portion of the entries are empty or contain no meaningful data.

SparsePIR stands out by achieving a response overhead that is only half of what state-of-the-art keyword PIR schemes typically require. Moreover, it offers this improved efficiency without necessitating the client to store any long-term data. This feature makes SparsePIR a highly promising solution for various applications, especially those dealing with large and sparse databases, such as search engines and advertising platforms.

We will present a review of this paper and discuss its results.

1 Introduction

PIR. Private Information Retrieval (PIR) is a cryptographic protocol that allows a user to retrieve information from a database without revealing which specific information they are interested in.

The database in the server is represented as an array.

The protocol divided to four algorithms (*Init*, *Query*, *Answer*, *Decrypt*):

- Init algorithm initializes the keys for communication between the server and the client.
- Query algorithm make a query in the client and sends it encrypted to the server.
- Answer algorithm receives the database and the encrypted query and send a response to the client.
- Decrypt algorithm receives the response and decrypt it.

2 Results summary and significance

	Client Storage	Encoding Size	Response Overhead
Client Mapping	n	n	$1 \times$
Cuckoo Hashing	$O(1)$	$(2 + \epsilon)n$	$2 \times$
Constant- Weight	$O(1)$	n	$2 - 9 \times$
SparsePIR	$O(1)$	$(1 + \epsilon)n$	$1 \times$
SparsePIR^g	$O(1)$	$(1 + \epsilon)n$	$1 \times$
SparsePIR^c	$< 11.2KB$	$(1.03)n$	$1 \times$

Figure 1: One-round keyword PIR comparison. Response overhead is compared to state-of-the-art standard PIR, Spiral, that query over dense n -entry arrays

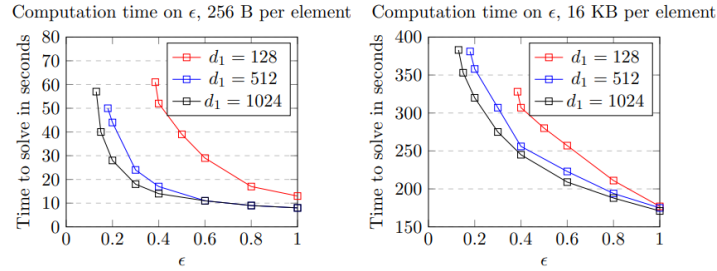


Figure 2: Computation time to solve the SparsePIR linear system vs ϵ on 2^{20} elements for dimension sizes $d_1 = 128, 512, 1024$. The graphs suggest that the ϵ lower bounds for $d_1 = 128, 512, 1024$ are roughly 0.38, 0.18, 0.13 respectively. We chose band parameter $w \in [50, 60]$.

Database	Onion CH-PIR	Onion SparsePIR	Onion SparsePIR ^g	Onion SparsePIR ^c	Spiral CH-PIR	Spiral SparsePIR	Spiral SparsePIR ^g	Spiral SparsePIR ^c
2²⁰ × 256 B								
Query Size	63 KB	63 KB	63 KB	63 KB	14 KB	14 KB	14 KB	14 KB
Response Size	254 KB	127 KB	127 KB	127 KB	42 KB	21 KB	21 KB	21 KB
Computation	3.63 s	3.04 s	3.10 s	3.05	1.41 s	1.44 s	1.45 s	1.42 s
Rate	0.001	0.002	0.002	0.002	0.006	0.012	0.012	0.012
Server Cost	\$0.000034	\$0.000027	\$0.000029	\$0.000028	\$0.0000091	\$0.0000074	\$0.0000074	\$0.0000073
2¹⁴ × 30 KB								
Query Size	63 KB	63 KB	63 KB	63 KB	14 KB	14 KB	14 KB	14 KB
Response Size	254 KB	127 KB	127 KB	127 KB	172 KB	86 KB	86 KB	86 KB
Computation	32.25 s	41.91 s	32.24 s	32.28 s	10.02 s	11.57 s	10.21 s	10.18 s
Rate	0.118	0.236	0.236	0.236	0.174	0.349	0.349	0.349
Server Cost	\$0.00015	\$0.00017	\$0.00014	0.00014	\$0.000054	\$0.000052	\$0.000047	\$0.000047
2¹⁴ × 100 KB								
Query Size	63 KB	63 KB	63 KB	63 KB	14 KB	14 KB	14 KB	14 KB
Response Size	1016 KB	508 KB	508 KB	508 KB	484 KB	242 KB	242 KB	242 KB
Computation	14.43 s	17.32 s	15.14 s	15.10 s	4.93 s	5.91 s	5.11 s	5.17 s
Rate	0.098	0.197	0.197	0.197	0.207	0.413	0.413	0.413
Server Cost	\$0.00014	\$0.00011	\$0.00010	\$0.00010	\$0.000061	\$0.000044	\$0.000041	\$0.000041

Figure 3: Comparison of cuckoo hashing (CH) and SparsePIR frameworks for various PIR protocols.

The schemes reduce the response size by at least 2x compared to prior keyword PIR constructions. In addition, keyword PIR may be built with identical communication and computation costs as standard PIR.

3 Techniques overview

3.1 Key-Value Pairs Representation

A more efficient representation of a database often involves using a key-value representation with hashing. This approach is preferred because the traditional array representation can be inefficient when dealing with sparse data, where most of the array elements are empty.

By employing a hash table, the data can be organized in a way that reduces sparsity compared to the array representation. Hashing allows for faster access to specific values based on their keys, which can significantly improve data retrieval performance.

For some hash key K , we denote $rep(K, k_i, v_i)$ as the hash evaluation of k_i concatenated with v_i (k_i, v_i are the client's key and value). Note that $rep(K, k_i, v_i) = F(K, k_i) || v_i$ where $F(K, k_i)$ is the hashed index.

3.2 Fully Homomorphic Encryption

All the computation are under fully homomorphic encryption scheme.

Given plaintext pt , the encrypted pt is $Enc(pt) = (r, r \cdot s + e + pt)$, where r is a uniformly random element, s is the secret key in FHE scheme and e is a polynomial noise where each coefficient is typically small.

To decrypt ciphertext (ct_1, ct_2) , one will compute $ct_2 - (ct_1 \cdot s) = e + pt$.

In the paper, FHE encryption was implemented in the communication and the Answer algorithm on the server side, ensuring that the server will not know the client's query.

3.3 Recursion

SparsePIR is compatible with the PIR optimization known as recursion.

Recursion allows a client to obtain multiple entries from the database in a single query. This can be done by recursively splitting the query into smaller queries, and then using the results of the smaller queries to reconstruct the results of the original query.

Represent the database, n -entry array, as a hypercube of dimensions $d_1 \times \dots \times d_z$.

The product of the dimensions is $d_1 \cdot \dots \cdot d_z \geq n$.

To query an entry, the client will generate z indicator vectors $v_1 \in \{0, 1\}^{d_1}, \dots, v_z \in \{0, 1\}^{d_z}$, where each v_i has zeros everywhere except for a one indicating the location of the entry with respect to the dimension d_i .

Benefits: Efficiency and Scalability.

3.4 Encoding

The process of encoding the database involves generating random binary vectors $v_k \in \{0, 1\}^m$ ($m = d_1$) for each entry in the database. These vectors are then encoded using $rep(K_r, k_i, v_i)$. The resulting encoded database, represented as matrix E , will be as follows:

$$\begin{pmatrix} v_{k_1}^T \\ v_{k_2}^T \\ \vdots \\ v_{k_n}^T \end{pmatrix} \cdot E = M \cdot E = y = \begin{pmatrix} rep(K_r, k_1, v_1) \\ rep(K_r, k_2, v_2) \\ \vdots \\ rep(K_r, k_n, v_n) \end{pmatrix}$$

Where, $v_{k_i}^T$ denotes the transpose of the randomly generated binary vector v_{k_i} for the i -th entry and y represents the database.

By following this encoding process, the original entries of the database are transformed into an encoded form.

3.5 Partition

The paper proposes a novel partitioning technique for SparsePIR.

The partitioning technique is based on the idea of dividing the database into a number of partitions, such that each partition contains a small number of database entries. This allows the client to generate a smaller random vector, which can significantly reduce the communication cost of SparsePIR.

The partitioning technique is implemented as follows:

First, the database is sorted by the number of keywords that each entry contains.

Then, the database is divided into a number of partitions, such that each partition contains a roughly equal number of entries. The number of partitions is chosen such that the average number of keywords per entry in each partition is small.

When the client generates a query, it first determines the partition that contains the entry that it is querying. Then, the client generates a random vector that is only used to select the shares from the partition that contains the desired entry. This allows the client to generate a smaller random vector, which can significantly reduce the communication cost of SparsePIR.

Algorithm 1 Partition

- 1: Create $b = (1 + \epsilon)n/d_1$ partitions where each part will be size d_1 .
- 2: Given input of Database $D = \{(k_1, v_1), \dots, (k_n, v_n)\}$ that needs to be encoded, assign the n key-value pairs to the b partitions uniformly at random, ensuring that in each part $n/b = d_1/(1 + \epsilon) < d_1$ pairs will be allocated.
- 3: As a result, the problem is effectively reduced to efficiently encoding databases of $O(d_1)$.
- 4: Consider the i -th part $P_i = \{(k_1, v_1), \dots, (k_{|P_i|}, v_{|P_i|})\}$:
- 5: Generate M_i that is a $|P_i| \times d_1$ matrix where each entry is uniformly chosen from $\{0, 1\}$ (in particular, the i -th row is generated randomly using K_2 and key k_i).
- 6: Compute $y_i^t = [rep(K_r, k_i, v_i), \dots, rep(K_r, k_{|P_i|}, v_{|P_i|})]$.
- 7: Finally, Solve the linear system $M_i \cdot e_i = y_i$ and obtain the encoding e_i for part P_i .
- 8: To build the final encoded database E , put each of e_1, \dots, e_b as column vectors in E to construct a $d_1 \times b$ matrix.
- 9: To query for a key k , the client will compute $i = F_1(K_1, k)$ to determine the partition associated with k .
- 10: The client randomly generates the associated random vector using hash key K_2 and k , denoted as v_{k_1} .
- 11: The server applied v_{k_1} to the first dimension of E :

$$[v_{k_1} \cdot e_1 \quad \dots \quad v_{k_1} \cdot e_b]$$

- 12: The only entry needed to be retrieved is $v_{k_1} \cdot e_i$ (given k is assigned to the i -th partition).
-

Benefits: Reduced communication cost, improved scalability and simplified implementation.

4 PIR

The underlying PIR scheme will represent any m -element database as a hypercube with dimensions $d_1 \times \dots \times d_z$. PIR is constructed from four algorithm:

- $(ck, sk) \leftarrow \text{Init}(1^\lambda)$: The Init algorithm produces client key ck and server key sk .
- $(st, req) \leftarrow \text{Query}(ck, k)$: The Query algorithm receives z vectors of length $d_1 \times \dots \times d_z$ that it will homomorphically encrypt and upload to the server. For standard PIR, these z vectors are indicator vectors representing the query index in each dimension.
- $resp \leftarrow \text{Answer}(sk, D, req)$: The Answer algorithm receives the database, in a form of two-dimensional matrix of size $d_1 \times \lceil m/d_1 \rceil$, and homomorphic encryptions of vectors v_1, \dots, v_n . The Answer algorithm will perform the standard PIR algorithm of applying v_1 to E to obtain a $\lceil m/d_1 \rceil$ vector, arrange the vector into a $d_2 \times \lceil m/d_1 \cdot d_2 \rceil$ matrix and apply v_2 to obtain a vector of size $\lceil m/d_1 \cdot d_2 \rceil$. Repeat this process for all z dimensions.
- $v \leftarrow \text{Decrypt}(ck, st, resp)$: The Decrypt algorithm receives the server's response and produces a decrypted answer.

5 SparsePIR

5.1 The Contribution

Algorithm 2 SparsePIR.Init Algorithm

Input: 1^λ : security parameter.
Output: (ck, sk) : client and server key.
 $(ck, sk) \leftarrow \Pi_{\text{PIR}}.\text{Init}(1^\lambda)$
return (ck, sk)

Algorithm 3 SparsePIR.Encode Algorithm

Input: $D = \{(k_1, v_1), \dots, (k_n, v_n)\}$: database.
Output: $(prms, E)$: parameters and encoding.
 Sample hash keys K_1, K_2 and K_r .
 $b \leftarrow (1 + \epsilon) \cdot n/d_1$
 $m \leftarrow (1 + \epsilon) \cdot n$
 $P_1 \leftarrow \emptyset, \dots, P_b \leftarrow \emptyset$
for $i = 1, \dots, n$ **do**
 $j \leftarrow F_1(K_1, k_i)$
 $P_{j+1} \leftarrow P_{j+1} \cup \{(k_i, v_i)\}$
for $j = 1, \dots, b$ **do**
 $e_i \leftarrow \text{GenerateEncode}(K_2, K_r, P_j)$
 $prms \leftarrow (K_1, K_2, K_r)$
 $E \leftarrow [e_1, \dots, e_b]$
return $(prms, E)$

Algorithm 4 GenerateEncode Algorithm

Input: (K_2, K_r, P) : hash keys and a part.
Output: e : an encoding of the part.
 $M \leftarrow []$ as empty array
 $y \leftarrow []$ as empty array
for $(k, v) \text{ in } P$ **do**
 Append $\text{RandVector}(K_2, k)^T$ to M as row
 Append $\text{rep}(K_r, k, v)$ to y
 $e \leftarrow \text{SolveLinearSystem}(M, y)$
return e

Algorithm 5 RandVector Algorithm

Input: (K_2, k) : the hash key and database key.
Output: v_k : a randomly generated vector.
 $v_k \leftarrow [0]^{d_1}$
for $i = 1, \dots, d_1$ **do**
 $v_k[i] \leftarrow F_2(K_2, k || i)$
return v_k

Algorithm 6 SolveLinearSystem Algorithm

Input: (M, y) : band matrix and values to solve for.
Output: e : a solution to the linear system $M \cdot e = y$.
Sort rows of M and y according to first non-zeros entry of M
Execute Gaussian elimination to get e
return e

Algorithm 7 SparsePIR.Query Algorithm

Input: $(prms, ck, k)$: parameters and the query key.
Output: (st, req) : temporary state and request.
 $b \leftarrow (1 + \epsilon) \cdot n/d$
Parse $prms = (K_1, K_2, K_r)$
 $v_1 \leftarrow \text{RandVector}(K_2, k)$
 $j \leftarrow F_1(K_1, k)$
for $i = 2, \dots, z$ **do**
 $j_i \leftarrow \lfloor j / \lceil b/d_i \rceil \rfloor$
 $v_i \leftarrow [0]_i^d$
 $v_i[j_i + 1] \leftarrow 1$
 $b \leftarrow \lfloor b/d_i \rfloor$
 $j \leftarrow j \bmod b$
 $(st_{\text{PIR}}, req) \leftarrow \Pi_{\text{PIR}}.\text{Query}(ck, v_1, \dots, v_z)$
return $((st_{\text{PIR}}, k), req)$

Algorithm 8 SparsePIR.Answer Algorithm

Input: $(prms, sk, E, req)$: parameters, server key, encoded databases and the request.
Output: $resp$: the response to the request.
 $resp \leftarrow \Pi_{\text{PIR}}.\text{Answer}(sk, E, req)$
return $resp$

Algorithm 9 SparsePIR.Decrypt Algorithm

Input: $(prms, ck, st, resp)$: parameters, client key, temporary state and response.

Output: v : output value.

Parse $prms$ as (K_1, K_2, K_r)

Parse st as (st_{PIR}, k)

$x \leftarrow \Pi_{PIR}.Decrypt(ck, st_{PIR}, resp)$

Parse x as (id, v)

if $id = F(K_r, k)$ **then**

return v

else

return \perp

5.2 Strengths Of SparsePIR

SparsePIR's main advantage lies in its efficiency. It surpasses state-of-the-art keyword PIR schemes by achieving only half the response overhead, making it highly promising for applications dealing with large and sparse databases.

Furthermore, the simplicity of SparsePIR is noteworthy. It utilizes a straightforward encoding scheme, which makes it easily comprehensible and implementable. Consequently, SparsePIR becomes an excellent option for applications prioritizing simplicity over the utmost security.

5.3 Weaknesses Of SparsePIR

Firstly, SparsePIR introduces a slight increase in server-side computation time due to cryptographic operations and encoding/decoding requirements.

However, while SparsePIR offers simplicity and efficiency, it may provide reduced security guarantees compared to more complex PIR schemes for the following reasons:

1. SparsePIR only supports simple queries, limiting its query capabilities. If a user can only issue queries for a single item, an attacker could exploit this to gain insights into the user's data through a series of queries.
2. SparsePIR is optimized for smaller databases compared to other PIR schemes. If the database is small, attackers may attempt to brute-force the encryption scheme to access the user's data.

Moreover, when a user interacts with a single server instead of multiple servers, there is a potential privacy concern. The server can observe the user's repeated requests over time, allowing them to infer the user's interests. This could lead to tracking of the user's activities.

Additionally, SparsePIR is most effective in scenarios with large and sparse databases, but it may not perform as well when dealing with dense databases or databases with frequent updates.

5.4 Open Problems

SparsePIR is not suitable for multiple servers, and it also lacks support for complex queries in its protocol.

References

- [1] Sarvar Patel, Joon Young Seo, and Kevin Yeo. Don't be dense: Efficient keyword pir for sparse databases. USENIX, 2023. <https://www.usenix.org/conference/usenixsecurity23/presentation/patel>.