

Reducing Communication Cost of Encrypted Data Search with Compressed Bloom Filters

Muhammad Umer*, Tahir Azim*[†] and Zeeshan Pervez[‡]

* National University of Sciences and Technology (NUST), Islamabad, Pakistan

[†] École Polytechnique Fédérale de Lausanne, Switzerland

[‡] University of the West of Scotland, Paisley, Scotland PA1 2BE

Email: 13msitnummer@seecs.edu.pk, tahir.azim@epfl.ch, zeeshan.pervez@uws.ac.uk

Abstract—Consumer data, such as documents and photos, are increasingly being stored in the cloud. To ensure confidentiality, the data is encrypted before being outsourced. However, this makes it difficult to search the encrypted data directly. Recent approaches to enable searching of this encrypted data have relied on trapdoors, locally stored indexes, and homomorphic encryption. However, data communication cost for these techniques is very high and they limit search capabilities either to a small set of pre-defined trapdoors or only allow exact keyword matching.

This paper addresses the problem of high communication cost of bloom-filter based privacy-preserving search for measuring similarity between the search query and the encrypted data. We propose a novel compression algorithm which avoids the need to send the entire encrypted bloom filter index back to the client. This reduces the cost of communicating results to the client by over 95%. Using sliding window bloom filters and homomorphic encryption also enables our system to search over encrypted data using keywords that only partially match the originally stored words. We demonstrate the viability of our system by implementing it on Google Cloud, and our results show that the cost of partially matching search queries on encrypted data scales linearly with the total number of keywords stored on the server.

I. INTRODUCTION

Cloud computing enables its subscribers to use computing services over the Internet and provides tailored computing resources on-demand. Subscribers of cloud storage often outsource their data on public cloud servers to ensure that it is highly available and reliable. As the amount of data stored on the cloud has increased dramatically, the security of this data has become an important concern. Recent cases have included users' private photos and videos getting stolen from compromised cloud servers. Credit card numbers and medical records are other examples of user data whose privacy needs to be ensured.

The straightforward way to achieve data protection and confidentiality is to place data on the cloud after encrypting it using public key encryption. This solves the confidentiality and privacy problem, but encryption hides information within the data and searching becomes difficult. The user has to download all of the data from cloud storage and search it after decrypting it locally. This is obviously an expensive proposition. Therefore we need a mechanism that allows a user to search over encrypted data without revealing private information and without downloading excessive data.

Existing approaches for searching over encrypted data often rely on “trapdoors” [1]–[3]. Trapdoors enable a user to search

the encrypted data for a small set of pre-defined keywords. These approaches restrict the searching capability of a user to a limited number of trapdoors defined during data encryption. More recent work has focused on homomorphic encryption [4] as a solution to this problem. However, to ensure privacy from the cloud service provider (CSP), all the evaluated results are returned from the cloud server in encrypted form without any compression. This causes high communication cost but allows searching for exact keyword matches over the encrypted data.

In this paper, we present a system for performing searches on encrypted data in the cloud. We propose a novel compression technique which allows us to avoid sending the entire encrypted index back to the client. This reduces the cost of communicating search results back to the client by over 95% resulting in faster response times and less budgeting cost to the owner. Besides exact keyword matching, our system also supports similarity-based searching: finding documents with partially matching keywords. Homomorphic encryption enables private search capability, while sliding window bloom filters enable searching for partially matching keywords.

II. RELATED WORK

Although cloud service providers (CSP) allow data to be encrypted for security purposes, encryption leads to two problems. First, searching for keywords within data encrypted using popular encryption algorithms is not possible. Second, if you provide the CSP the decryption keys, you cannot prevent the CSP from accessing the confidential data. Recent work has focused on solving these problems using a variety of techniques, which we broadly classify into two categories.

Using Trapdoor Encryption. Song et al. [1] propose a symmetric key based searchable scheme for encrypted data. Each keyword in the document is encrypted independently using trapdoors. Their approach proposes ways to search encrypted data using both an encrypted index as well as searching the original data itself. Goh et al. [2] propose encrypted search using bloom filters. They generate trapdoors against all the keywords in a file and add them to a bloom filter stored in the cloud. To search, a user computes the trapdoor for a keyword and sends it to the cloud server. The cloud server checks the trapdoor in the bloom filter, and if it exists, returns the corresponding file identifier. Boneh et al [3] present one of the earliest public key based encrypted search algorithm,

enabling authorized users having the private key to search in the index. Yu et al. [5] use Hierarchical Predicate Encryption (HPE) for encrypted search. All of the above schemes support only exact keyword matches and rely on pre-defined trapdoors.

Pal et al. [6] propose an encrypted search scheme using bloom filters and use soundex coding to search similar words. Kuzu et al. [7] carry out similarity-based encrypted search using locality sensitive hashing for similarity score calculation. In general, none of these works attempt to reduce the communication cost of returning results to the client.

Using Homomorphic Encryption. Homomorphic encryption allows computations to be carried out on ciphertext, thus generating an encrypted result, which when decrypted, matches the result of operations performed on the plaintext. An example of homomorphic encryption is the Pascal Paillier cryptosystem [8], which provides two useful properties: (i) the product of two ciphertexts will decrypt to the sum of their corresponding plaintexts and (ii) a ciphertext raised to a constant k will decrypt to the product of the corresponding plaintext and the constant.

Pervez et al. [9] propose an encrypted data search scheme using an inverted index and homomorphic encryption [4]. They use homomorphic encryption to provide end-to-end privacy. Only authorized users can execute queries using their personal proxy re-encryption keys in order to manipulate the index, so a lot of computation is done due to index transformations. They use a trusted third party to rank search results and model queries. While this approach avoids using trapdoors, it still supports only exact keyword matching and has high communication cost. Finally, CryptDB [10] and MONOMI [11] improve the performance of searching over encrypted data by employing customized forms of encryption for specific kinds of data. However, in doing so, they are vulnerable to several kinds of information leakage but with low probability.

III. PROPOSED SYSTEM

We first discuss some of our assumptions about the system and the threat model. Then we describe the main steps involved in the initial index generation and then the actual search process over encrypted data.

A. System Model

Entities involved in the proposed system are the data owner, CSP and end user. A data owner is an entity who wants its confidential data to be stored in cloud storage with the capability of privacy-aware searching. A CSP hosts public cloud storage services for its subscribers on a pay-as-you-use model. The end user is an entity that will perform searches on the encrypted data stored in the cloud. The end user can submit search queries to the CSP, which evaluates the queries and returns results back to the user. However, during query evaluation, the CSP should not be able to learn anything about the query, the stored data or the matching results.

B. Threat Model

The data owner and end user are considered trusted entities while CSP is considered as an untrusted entity because it is

TABLE I
NOTATIONS USED IN MATHEMATICAL AND DESCRIPTIVE DETAILS

Notation	Description
\mathcal{F}	Confidential file that needs to be outsourced
$\mathcal{K}\omega$	List of keywords in a data file \mathcal{F}
$f\omega$	Frequency of a keyword $k\omega$
$\mathcal{E}_h, \mathcal{D}_h$	Homomorphic encryption and decryption functions
$\mathcal{E}_s, \mathcal{D}_s$	Symmetric encryption and decryption functions
σ_{pk}, σ_{sk}	Homomorphic encryption public and private keys
\mathcal{R}_q	Compressed resultant term after bloom filter matching
\mathcal{M}_q	Uncompressed resultant term after bloom filter matching
\mathcal{BF}_i	i 'th encrypted bit of bloom filter \mathcal{BF}
O_b	Number of 1 bits in a bloom filter
I	Index file containing encrypted bloom filter for each keyword, frequency and number of 1s as index entries

maintained by an arbitrary third party. Data communication between the end user and CSP should be considered as untrusted as all requests are routed via the open Internet. The CSP hosts the encrypted data file (\mathcal{F}) and the encrypted index (I). As they are encrypted, the CSP cannot learn any information regarding \mathcal{F} and I . Query evaluation is done by CSP homomorphically: that is, the CSP performs evaluation over encrypted text of I and encrypted user query, so it is neither able to learn anything about the matched results, nor can it relate any subsequent queries from other users.

C. Assumptions and Notations

We ignore the key exchange mechanism between the data owner and the end users, assuming that it happens using secure out-of-band channels. Table I illustrates the notations that we use in order to explain the details of our proposed approach.

D. Index Generation

The data owner generates an encrypted index on the client side to facilitate subsequent searches. At a high level, the index is generated by creating a *sliding window bloom filter* which is then encrypted using the Pascal Paillier homomorphic encryption algorithm [8].

The sliding window bloom filter (SWBF) is a special type of bloom filter for which a window size is defined, and based on that window size, a keyword is sliced and mapped to the bloom filter. For example, suppose we have a word “cloud” and a window size of 2. The word will then be sliced into “cl”, “lo”, “ou”, and “ud” and each of these slices will be independently mapped to the bloom filter. This filter enables us to achieve a partial matching even if the requested keyword does not match completely.

Algorithm 1 describes the indexing procedure. For each file to be uploaded to the server, a separate index file is generated. The index file contains one index entry per keyword. To generate the index file, a sliding window bloom filter is updated for each keyword and the number of 1 bits in the bloom filter are noted. Each bit of the bloom filter is then encrypted using the Pascal Paillier algorithm, and written to the index entry along with frequency and number of 1 bits. At the end of this process, each index entry I_i in an index file I has the structure: $I_i = [\mathcal{BF}_i, f\omega, O_b]$, where O_b is the number of 1's in \mathcal{BF}_i and $\mathcal{BF}_i = \mathcal{BF}_1, \mathcal{BF}_2, \dots, \mathcal{BF}_n$.

Once the encrypted index file has been generated, it is uploaded to the cloud along with the data files. The data files are encrypted using a symmetric encryption algorithm, as they will not be used during the search process.

E. Search Process

To start the search process, the user first generates a query through a similar process as index creation. Each of the user's specified keywords are used to generate sliding window bloom filters, whose bits are then encrypted using the public key σ_{pk} of the Pascal Paillier algorithm. Note that this process takes place on the end user's machine. After encryption, the query (comprising the encrypted bloom filters) is sent to the cloud for terms matching.

On receiving a query, the CSP matches the incoming query with the index entries of the files it is hosting. For a perfect match, all of the corresponding bits of the index and query bloom filters need to match. However, since both the query and the index bloom filters are encrypted, simple matching is not possible. Furthermore, encrypting the same plaintext repeatedly results in completely different ciphertexts, which improves privacy and security but makes matching difficult.

Fortunately, homomorphic encryption provides a way to perform mathematical operations over ciphertext. Since our bloom filter entries are encrypted using Pascal Paillier homomorphic encryption, we can just multiply their ciphertexts, which when decrypted, yields the sum of the plaintexts [8]. The sum of the bits will be either 0, 1 or 2 after decryption. The twos in the decrypted output then represent the matched bits and the ones

represent the unmatched bits. Then, a similarity score can be estimated as the ratio of twos to ones in the decrypted result.

Note that decryption of this result will require users to have the private Pascal Paillier key. Thus, only authorized users who have received the private key from the data owner can actually decrypt the results.

F. Compressing Search Results

Once the encrypted index entries and the query have been multiplied together, we can simply return these products back to the client. The client can then decrypt the products to get the sums and compute similarity scores. While this straightforward approach works and is used by many existing systems (for example, by [9]), it is extremely inefficient. In particular, if there are many documents and keywords in the cloud dataset, this will result in a huge amount of data to be communicated back to the user. Specifically, the product of the query and the encrypted index entry for *every* keyword in the dataset will be returned to the client. This increases both response time over the network as well as the cost of cloud computing for the data owner, since network usage is billed by most CSPs.

Our technique to reduce this communication cost overhead is based on the insight that each returned index entry consists only of 0s, 1s or 2s after decryption. We exploit this property by defining a polynomial \mathcal{P} such that:

$$\mathcal{P} = a^0 \cdot \mathcal{BF}_1 + a^1 \cdot \mathcal{BF}_2 + a^2 \cdot \mathcal{BF}_3 + \dots + a^{n-1} \cdot \mathcal{BF}_n$$

where $a \in \{3, 5, 7, \dots\}$, \mathcal{BF}_i are the Paillier sums of the corresponding bits of an index entry and represent variables for the polynomial \mathcal{P} and $\mathcal{BF}_i \in \{0, 1, 2\}$. We used $a = 3$ as it yields smaller sums and less computation.

The sum of this polynomial \mathcal{P} is the resultant compressed term which we return back to the client: $\mathcal{R}_q = \sum_{i=0}^{n-1} a^i \cdot \mathcal{BF}_{i+1}$. The size of this single compressed term \mathcal{R}_q is equal to the size of the Pascal Paillier key. This multiplication by constants and addition is again possible because our bloom filter entries are encrypted using Pascal Paillier homomorphic encryption. Therefore, if after decryption, we want the product of the plaintext with a constant, we just need to exponentiate the ciphertext by that constant [8].

This means that irrespective of the size of the bloom filter, we will always return a single number back to the client. The client can then decompress it to find out the original Paillier sums using the routine specified in Algorithm 2. In short, this algorithm works correctly because every \mathcal{BF}_i is at most 2 and $a^n > 2 \sum_{j=0}^{n-1} a^j$. This means that taking $\log_a \mathcal{R}_q$ repeatedly will yield the position of the next bloom filter entry that needs to be incremented.

The savings achieved using this algorithm can be evaluated using a simple example. Suppose we use a 32-bit bloom filter and a 64-bit Pascal Paillier key. Then for 1000 keywords, the data owner will create an index file with 1000 entries. The size of the index file will then be $32 \cdot 64 \cdot 1000 / 8 = 256 \text{KB}$. On the other hand, after compression, the number returned per keyword will just have a size equal to the size of the Paillier key, namely 64 bits. So the total data returned will be $64 \cdot 1000 / 8 = 8 \text{KB}$, a saving of over 95%.

Algorithm 1: Index Creation

Input: A collection of text files $C = \langle \mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_n \rangle$

Output: Index files $I = \langle I_1, I_2, \dots, I_n \rangle$ for each file in C

```

1  $\forall \mathcal{F}_i \in \langle \mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_n \rangle;$ 
2 while  $\mathcal{F}_i \in C$  do
3    $\mathcal{K}\omega \leftarrow \text{extractAllKeywords}(\mathcal{F}_i);$ 
4    $\forall k\omega_j \in \mathcal{K}\omega;$ 
5   while  $k\omega_j \in \mathcal{K}\omega$  do
6      $SWBF \leftarrow \text{createSWBF}(k\omega_j);$ 
7      $f_j \leftarrow \text{getKeywordFrequency}(k\omega_j);$ 
8      $O_b \leftarrow \text{getOnes}(\mathcal{BF}_j);$ 
9      $\forall b f_k \in SWBF;$ 
10    while  $b f_k \in SWBF$  do
11       $\mathcal{BF}_k \leftarrow \mathcal{E}_h(\sigma_{pk}, b f_k);$ 
12       $\text{IndexEntry} \leftarrow \text{IndexEntry}, \mathcal{BF}_k;$ 
13       $b f_k \leftarrow \text{getNextBit}(SWBF);$ 
14     $\text{IndexEntry} \leftarrow \text{IndexEntry}, f_j, O_b;$ 
15     $I_j \leftarrow \text{writeToIndexFile}(\text{IndexEntry});$ 
16     $I \leftarrow I, I_j;$ 
17     $k\omega_j \leftarrow \text{getNextKeyword}(\mathcal{K});$ 
18   $\mathcal{F}_i \leftarrow \text{getNextFile}(C);$ 
19 return  $I;$ 
```

Algorithm 2: Index entry decompression

Input: The compressed index sum \mathcal{R}_q

Output: Matched result \mathcal{M}_q comprising a sequence of 0s, 1s and 2s

```
1  $i \leftarrow 0$ ;  
2 while  $\mathcal{R}_q > 0$  do  
3    $i \leftarrow \log_3 \mathcal{R}_q$  ;  
4    $\mathcal{R}_q \leftarrow \mathcal{R}_q - 3^i$  ;  
5    $\mathcal{M}_q[i] \leftarrow \mathcal{M}_q[i] + 1$  ;  
6 return  $\mathcal{M}_q$ ;
```

IV. EVALUATION

We demonstrate the viability of our proposed scheme using Google Cloud. We implement an indexing service, a search service and a client application as standard Java services. We deploy the search service on Google App Engine configured as an F4 class instance with a 2.4 GHz CPU and 512 MB of RAM. We utilize Google Blobstore for hosting index files because it allows us to store and retrieve the complete index files efficiently. We evaluate the results of our approach on a dataset of 150 documents. These documents range in size from 5 MB to 100 MB, containing from 5000 to 129,000 keywords. The keywords are chosen to be at least 8 characters in length and Porter stemming [12] is applied to those keywords. The client is a Lenovo Thinkpad 430 with a 2.6 GHz Intel Core(TM) i5-332M CPU and 8 GB of RAM.

A. Search Performance

For private term matching and data returned, our implementation demonstrates that our compression algorithm results in 95% less data returned to the client compared to traditional approaches, and remains constant even when the size of the bloom filter increases. This is shown in Figure 1.

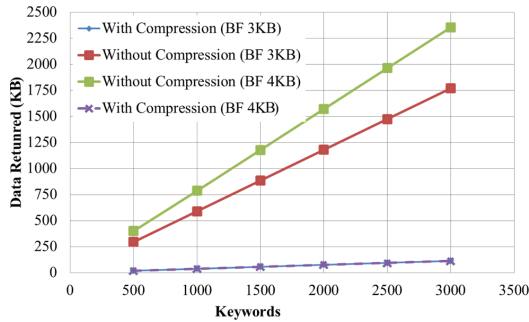


Fig. 1. Data returned by a search query: compressed vs uncompressed results

B. Cost Estimation and Response Times

Finally, we evaluate the dollar cost that CSP will charge to data owners for search queries with partial matching. Our results, shown in Figure 2, demonstrate that a query searching for a single keyword in a dataset having 500-3500 index entries will cost only \$0.000002 to \$0.000002 per 1000 similar queries. Due to the sliding window bloom filter approach, all of the metrics increase only linearly ($O(n)$) with the number of

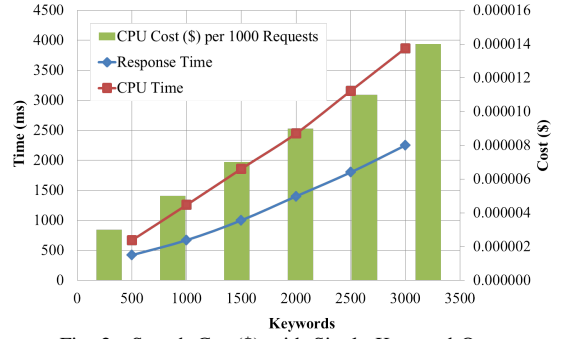


Fig. 2. Search Cost (\$) with Single Keyword Query

keywords on the server. This is substantially better than using a naive algorithm for partial matching, which would result in $O(nm)$ time complexity for n keywords that are on average m characters long. Crucially, privacy is preserved throughout the search process because data is never decrypted at the cloud or by any other untrusted system.

V. CONCLUSION

We presented a privacy-aware similarity-based searching scheme with reduced communication cost over encrypted data residing in an untrusted cloud domain. Using a novel compression algorithm, our system avoids the need to send the bloom filter based index back to the user, reducing communication costs by over 95%. By using homomorphic encryption for index files, a cloud service provider can not learn the contents of the data files, the index files or the search queries. Furthermore, it cannot discern patterns from incoming queries. Moreover, unlike most existing techniques, our proposed system supports similarity-based searching in addition to exact matching. Finally, by not relying on trapdoors, it allows end users to formulate arbitrary queries to search encrypted data.

REFERENCES

- [1] Dawn Xiaodong Song, David Wagner, and Adrian Perrig, "Practical techniques for searches on encrypted data," in *Proceedings of the IEEE Symposium on Security and Privacy, 2000 (S&P 2000)*.
- [2] Eu-Jin Goh et al., "Secure indexes," *IACR Cryptology ePrint Archive*, vol. 2003, pp. 216, 2003.
- [3] Dan Boneh, Giovanni Di Crescenzo, Rafail Ostrovsky, and Giuseppe Persiano, "Public key encryption with keyword search," in *Advances in Cryptology-Eurocrypt 2004*. Springer, 2004, pp. 506–522.
- [4] Craig Gentry et al., "Fully homomorphic encryption using ideal lattices," in *STOC*, 2009, vol. 9, pp. 169–178.
- [5] Shucheng Yu, Cong Wang, Kui Ren, and Wenjing Lou, "Achieving secure, scalable, and fine-grained data access control in cloud computing," in *Infocom, 2010 proceedings IEEE*. Ieee, 2010.
- [6] Sankar K Pal, Puneet Sardana, and Ankita Sardana, "Efficient search on encrypted data using bloom filter," in *Computing for Sustainable Global Development, 2014 International Conference on*. IEEE, 2014.
- [7] Mehmet Kuzu, Mohammad Saiful Islam, and Murat Kantarcioglu, "Efficient similarity search over encrypted data," in *Data Engineering (ICDE), 2012 IEEE 28th International Conference on*. IEEE, 2012.
- [8] Pascal Paillier, "Public key cryptosystems based on composite degree residuosity classes," in *17th international conference on theory and application of cryptographic techniques*. Springer, 1999.
- [9] Zeeshan Pervez, Ammar Ahmad Awan, Asad Masood Khattak, Sungyoung Lee, and Eui-Nam Huh, "Privacy-aware searching with oblivious term matching for cloud storage," *The Journal of Supercomputing*, vol. 63, no. 2, pp. 538–560, 2013.

- [10] Raluca Ada Popa, Catherine Redfield, Nickolai Zeldovich, and Hari Balakrishnan, “Cryptdb: protecting confidentiality with encrypted query processing,” in *Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles*. ACM, 2011.
- [11] Stephen Tu, M Frans Kaashoek, Samuel Madden, and Nickolai Zeldovich, “Processing analytical queries over encrypted data,” in *Proceedings of the VLDB Endowment*. VLDB Endowment, 2013, vol. 6.
- [12] Martin Porter, “The Porter Stemming Algorithm,” <http://tartarus.org/martin/PorterStemmer/>.