

# Secure and Efficient Multi-keyword Ranked Search over Encrypted Cloud Data

Xingming Sun, Lu Zhou, Zhangjie Fu, Zhihua Xia and Jiangang Shu

College of Computer and Software & Jiangsu Engineering Center of Network Monitoring,  
Nanjing University of Information Science and Technology, Nanjing 210044, CHINA

**Abstract:** With the development of cloud computing, the sensitive information of outsourced data is at risk of unauthorized accesses. To protect data privacy, the sensitive data should be encrypted by the data owner before outsourcing, which makes the traditional and efficient plaintext keyword search technique useless. Hence, it is an especially important thing to explore secure encrypted cloud data search service. In this paper, we propose a practically efficient and flexible searchable encrypted scheme which supports multi-keyword ranked search. To support multi-keyword search and result relevance ranking, we adopt Vector Space Model (VSM) to build the searchable index to achieve accurate search result. To improve search efficiency, we design a tree-based index structure. We propose a secure search scheme to meet the privacy requirements in the threat model. Finally, experiments on real-world dataset show that our scheme is efficient.

**Keywords:** Multi-keyword search, ranked search, encrypted cloud data, security

## 1 Introduction

Cloud Computing is a new but increasingly mature model of enterprise IT infrastructure that provides on-demand high quality applications and services from a shared pool of configuration computing resources [1]. The cloud customers, individuals or enterprises, can outsource their local complex data system into the cloud to avoid the costs of building and maintaining a private storage infrastructure. However, some problems may be caused in this circumstance since the Cloud Service Provider (CSP) possesses full control of the outsourced data. Unauthorized operation on the outsourced data may exist on account of curiosity or profit. To protect the privacy of sensitive information, sensitive data (e.g., emails, photo albums, personal health records, financial records, etc.) should be encrypted by the data owner before outsourcing [2], which makes the traditional and efficient plaintext keyword search technique useless. The simple and awkward method of downloading all the data and decrypting locally is obviously impractical. So, two aspects should be concentrated on to explore privacy-preserving effective search service. Firstly, ranked search, which can enable data users to find the most relevant information quickly, is a very important issue. The number of documents outsourced to the cloud is so large that the cloud should have the ability to perform search result ranking to meet the demand for

effective data retrieval [3]. Secondly, multi-keyword search is also very important to improve search result accuracy as single keyword search often return coarse search results.

In this paper, we propose a practically efficient and flexible searchable encrypted scheme. To address multi-keyword search and result ranking, we use Vector Space Model (VSM) [6] to build document index. To improve search efficiency, we use a tree-based index structure which is a balanced binary tree (see the details in section 2.5). We construct the searchable index tree based on the document index vectors. Our encryption scheme can meet the privacy requirements in the threat model. Our contributions are summarized as follows:

- (1) We study the problem of multi-keyword ranked search over encrypted cloud data while supporting strict privacy requirements.
- (2) With the index tree designed in this paper, the search time complexity close to  $O(r \log m)$  ( $r$  is the number of documents including the search keywords and  $m$  is the whole number of documents in the dataset).
- (3) We make security analysis for our scheme which proves privacy guarantees. And experiments on the real-world dataset show that proposed scheme is indeed efficient.

## 2 Problem formulation

### 2.1 System model

Like in [4], we consider there are three different entities in the system model: the data owner, the user, and the cloud server respectively. The data owner encrypts document collection  $DC$  in the form of  $C$  before outsourcing it to the cloud in order to protect the sensitive data from unauthorized entities. And for the purpose of searching interested data, the data owner will also generate an encrypted searchable index  $I$  based on a set of distinct keywords  $W$  extracted from  $DC$ . In the search stage, the system will generate an encrypted search trapdoor based on the keywords entered by the user (has been authorized by data owner). Given the trapdoor, the cloud server will search the index  $I$  and then return the ranked search results to the user. As the search results are well ranked by the cloud server, the user can send a parameter  $k$  together with search query to get top- $k$  most relevant documents. As the issue of key distribution is out of the scope of this paper, we assume that data users have been authorized by data owner.

### 2.2 Threat model

In our system model, we consider that the cloud server is “honest-but-curious”

adopted by most previous searchable encryption schemes [3-5]. That is to say, the cloud server honestly implements the protocol and correctly returns the search results, but it is also curious to infer sensitive information during execute protocol. In the known ciphertext model, only the encrypted dataset  $C$ , the encrypted search query and the searchable index  $I$  are available to the cloud server.

### 2.3. Notations

The main notations used in this paper are showed as follows:

- $DC$  – the plaintext document collection, expressed as a set of  $m$  documents  $DC = \{d \mid d_1, d_2, \dots, d_m\}$ .
- $C$  – the encrypted document collection for  $DC$  stored in the cloud server, expressed as  $C = \{c \mid c_1, c_2, \dots, c_m\}$ .
- $W$  – the dictionary, including  $n$  keywords extracted from  $DC$ , expressed as  $W = \{w \mid w_1, w_2, \dots, w_n\}$ .
- $\tilde{W}$  – a subset of  $W$ , representing the keywords in a search request, expressed as  $\tilde{W} = \{w_{i1}, w_{i2}, \dots, w_{it}\}$ .
- $I$  – the searchable index tree generated from the whole document set  $DC$ . Each leaf node in the index tree is associated with a document in  $DC$ .
- $D_d$  – the index vector of document  $d$  for all the keywords in  $W$ .
- $Q$  – the query vector for the keyword set  $\tilde{W}$ .
- $\tilde{D}_d$  – the encrypted index vector for  $D_d$ .
- $\tilde{Q}$  – the encrypted query vector for  $Q$ .
- $f(key, \cdot)$ ,  $g(key, \cdot)$  – pseudorandom function (PRF), defined as:  $\{0,1\}^* \times key \rightarrow \{0,1\}^k$ .
- $Enc(key, \cdot)$ ,  $Dec(key, \cdot)$  – symmetric encryption/decryption function.
- $T_{\tilde{W}}$  – the encrypted form of  $\tilde{W}$ .

### 2.4. Preliminaries

**Keywords hash table:** A static hash table for all the keywords in  $W$ , denoted as  $\lambda$ . There are  $n$  entries in  $\lambda$  and each entity is a tuple  $(key, value)$ , in which the  $key$  is from a domain of exponential size, i.e., from  $\{0,1\}^k$  representing a keyword in  $W$ , and  $value$  is a boolean value which has been encrypted. For a

key  $x \in \{0,1\}^k$ , the corresponding value is denoted as  $\lambda[x]$ .

**Similarity function:** We employ the similarity evaluation function for cosine measure from [7]. Each document  $d$  in the dataset is corresponding to an  $n$ -dimension index vector  $D_d$ , and each dimension of  $D_d$ , denoted as  $D_d[i]$ , is related to a keyword  $w_i$  in  $W$ . If document  $d$  contains keyword  $w_i$ ,  $D_d[i]$  stores the normalized TF weight of  $w_i$  within document  $d$ , otherwise  $D_d[i] = 0$ . For a search request  $\tilde{W}$ , an  $n$ -dimension query vector  $Q$  is also generated. It is similar with document index vector that each dimension of  $Q$  is related to a keyword in  $W$ . And if  $\tilde{W}$  contains keyword  $w_i$ ,  $Q[i]$  stores the normalized IDF weight of  $w_i$ , otherwise  $Q[i] = 0$ . The notations used in our similarity evaluation function are showed as follows:

$$SC(Q, D_d) = \frac{\sum_{i=1}^n w_{q,i} \cdot w_{d,i}}{\sqrt{\sum_{i=1}^n (w_{q,i})^2} \cdot \sqrt{\sum_{i=1}^n (w_{d,i})^2}}$$

(1)

where  $w_{d,i}$  represents the TF weight of  $w_i$  within document  $d$ ,  $w_{q,i}$  represents the IDF weight of keyword  $w_i$ . We use functions  $w_{d,i} = 1 + \ln(f_{d,i})$  and  $w_{q,i} = \ln(1 + m / f_i)$  to compute the value of TF and IDF weight respectively, where  $f_{d,i}$  represents the TF of keyword  $w_i$  within the document  $d$ ,  $f_i$  represents the number of documents containing the keyword  $w_i$ ,  $m$  represents the total number of documents in the document collection,  $n$  represents the total number of keywords in the keyword dictionary. And hence, the vector  $Q$  and  $D_d$  are both unit vectors.

## 2.5. Searchable index tree

Our searchable index is a balanced binary tree. In order to make it easy to understand, the document index vector in each leaf node only stores TF information rather than normalized TF weight. Given the document collection  $DC = \{d \mid d_1, d_2, \dots, d_m\}$ , we can build the index tree  $I$ . The data structure is built using the procedure, denoted as *buildIndex* ( $DC$ ), showed as follows: (1) For each document  $d_j$  in  $DC$ , we generate a leaf node where stores document identifier  $j$  and index vector  $D_{dj}$ . (2) Then we build the tree following a postorder traversal with all leaf nodes generated

in (1). Each internal node  $u$  of the index tree stores an  $n$ -bit vector  $D_u$ . If  $D_u[i] = 1$ , then there is at least one path from  $u$  to a leaf node of which corresponding document contains keyword  $w_i$ . (3) In this step, we introduce how to generate vector  $D_u$  in each internal node  $u$ . Let  $v$  and  $w$  be the left child and right child of internal node  $u$  respectively, then  $D_u[i] = 1$  if  $D_v[i] \neq 0$  or  $D_w[i] \neq 0$ , otherwise  $D_u[i] = 0$ . Note that when the node  $v$  ( $w$ ) is a leaf node and stores identifier  $j$ ,  $D_v = D_{dj}$  ( $D_w = D_{dj}$ ).

**Tree-based search algorithm:** The sequential search process for keywords in a search request  $\tilde{W}$  conducts as follows: the procedure starts from the root node and when arrives at an internal node  $u$ , if at least a keyword  $w_k$  ( $k$  is the order number of  $w$  in  $W$ ) in  $\tilde{W}$  leads to  $D_u[k] = 1$ , it continues to search both subtrees of  $u$ , otherwise stops searching in the subtree  $t_u$  ( $t_u$  denotes the tree whose root is  $u$ ) because none of leaf node in  $t_u$  contains keyword in search query. When arrives at a leaf node, the process computes the cosine value between the index vector stored in the leaf node and the query vector as the similarity score. We denote the number of documents that contain the keyword in the search query as  $r$ . In the sequential search, the procedure will traverse as many paths as  $r$ . So, the search complexity is  $O(r \log m)$  as the height of a balanced binary tree with  $m$  leaf nodes is  $\log m + 1$ .

### 3 Secure Index Scheme

In order to achieve accurate multi-keyword ranked search, we propose to adopt VSM and cosine measure to evaluate similarity scores. By using the cryptographic methods similar to the techniques adopted in [4, 8], the document index vector and query are both well protected. The scheme is described as follows:

**Setup:** In this phase, we initialize our system. The data owner generates the secret key  $SK$ . The  $SK$  includes: 1) a  $n$ -bit vector  $S$  which is randomly generated; 2) two  $n \times n$  also randomly generated invertible matrices  $\{M_1, M_2\}$ ; 3) two randomly picked key  $sk_1$  and  $sk_2$ . Hence,  $SK$  is in the form of a 5-tuple as  $\{S, M_1, M_2, sk_1, sk_2\}$ .

**GenIndex**( $DC, SK$ ): The data owner calls procedure *buildIndex* ( $DC$ ) that defined in section 2.5. Then, every document index vector  $D_d$  is split into two random vectors denoted as  $\{D'_d, D''_d\}$ . The splitting procedure is expressed

as follow: take  $S$  as the splitting indicator, if the  $j$ -th bit of  $S$  is 0,  $D'_d[j]$  and  $D''_d[j]$  are set as the same as  $D_d[j]$ ; if the  $j$ -th bit of  $S$  is 1,  $D'_d[j]$  and  $D''_d[j]$  are set randomly so long as their sum is equal to  $D_d[j]$ . So, the encrypted index vector  $\tilde{D}_d$  is denoted as  $\tilde{D}_d = \{M_1^T D'_d, M_2^T D''_d\}$ . Store  $\tilde{D}_d$  at the leaf node that stores correspondent  $D_d$  and delete  $D_d$ . For each internal node  $u$  in the index tree, a static hash table  $\lambda$  is generated. There are  $n$  tuples  $(key, value)$  in  $\lambda$ , and for every  $i = 1, 2, \dots, n$ , set  $\lambda_u[f(sk_1, w_i)] = Enc(g(sk_2, w_i), D_u[i])$ . Store  $\lambda_u$  in internal node  $u$  and delete  $D_u$ . Finally, the encrypted searchable index tree  $I$  is generated.

**GenQuery** ( $\tilde{W}, SK$ ): With the interested keywords in  $\tilde{W}$ , the  $n$ -dimension query vector  $Q$  is generated. Each dimension of  $Q$  is a normalized IDF weight of corresponding keyword. Specifically, if  $i$ -th keyword of  $W$  is in  $\tilde{W}$ ,  $Q[i] = w_{q,i}$ , otherwise  $Q[i] = 0$ . Next,  $Q$  is also split into two random vectors as  $\{Q', Q''\}$  using the similar splitting procedure used for document index vector. The difference is that if the  $j$ -th bit of  $S$  is 0,  $Q'[j]$  and  $Q''[j]$  are set randomly so long as their sum is equal to  $Q[j]$ ; if the  $j$ -th bit of  $S$  is 1,  $Q'[j]$  and  $Q''[j]$  are set as the same as  $Q[j]$ . Then, the encrypted query vector  $\tilde{Q}$  is in the form of  $\{M_1^{-1} Q', M_2^{-1} Q''\}$ . Next,  $T_{\tilde{W}} = \{f(sk_1, w_{i1}), g(sk_2, w_{i1}), f(sk_1, w_{i2}), g(sk_2, w_{i2}), \dots, f(sk_1, w_{in}), g(sk_2, w_{in})\}$  is produced by encrypting each item in  $\tilde{W}$ . Finally, the  $\{T_{\tilde{W}}, \tilde{Q}\}$  is sent to the cloud server.

**Search** ( $I, \tilde{Q}, T_{\tilde{W}}, k$ ): The cloud server follow the search algorithm expressed in section 2.5. Let  $u$  be an internal node in  $I$ , and let  $a_i = \lambda_u[f(sk_1, w_i)]$  for each item  $f(sk_1, w_i)$  in  $T_{\tilde{W}}$ . If exist at least one  $a_i$  satisfies  $Dec(g(sk_2, w_i), a_i) = 1$ , the procedure continue to search all children of  $u$ . When arrive at a leaf node, the procedure obtains the encrypted document vector  $\tilde{D}_d$  and compute the similarity of  $\tilde{D}_d$  and  $\tilde{Q}$  using Eq(2). Finally, the cloud server will rank the searched documents based on their similarity scores.

$$\begin{aligned}
 SC(\tilde{D}_d, \tilde{Q}) &= \{M_1^T D'_d, M_2^T D''_d\} \cdot \{M_1^{-1} Q', M_2^{-1} Q''\} \\
 &= D'_d \cdot Q' + D''_d \cdot Q'' \\
 &= D_d \cdot Q
 \end{aligned} \tag{2}$$

## 4 Performance Analysis

In this section, we estimate the overall performance of our proposed scheme by implementing the secure search system using C# language on a Windows7 server with Intel(R) Core(TM)2 Quad CPU 2.83GHz. The document set is built from the real-world data set: Request for comments database (RFC) [9], which includes about 6500 publications.

### 4.1. Index tree construction

It is obvious that the time cost of the index tree construction is mainly affected by the number of documents in the dataset and keywords in the dictionary. For each internal node in the searchable index tree, the major computation is the encryption of the hash table, the time cost of which is proportional to the number of keywords in the dictionary. And for each leaf node in the index tree, the main computation is the encryption of the document index vector, which mainly depends on the time cost for two multiplications of a  $H \times H$  matrix and an  $H$ -dimension vector where  $H$  is  $n$  in our scheme. And the whole number of nodes in the index tree is related to the number of documents in the dataset. Fig. 1(a) shows that, given the same dictionary with 4000 keywords, the time cost for building the index tree is nearly linear with the number of documents in the dataset. Fig. 1(b) indicates that the time cost for constructing index tree is proportional to the number of keywords in the dictionary with the same size of dataset. Although the time cost for constructing index tree is not an ignorable overhead for the data owner, it is a one-time operation before data outsourcing.

### 4.2. Query generation

The time cost for generating the search query is greatly affected by the size of keywords dictionary. Two multiplications of a matrix and a split query vector are conducted in every query generation phase. The dimensionality of matrix ( $n \times n$ -dimension) and query vector ( $n$ -dimension) depends on the size of keywords dictionary. Fig. 2 shows the relationship between the time of generating search query and the number of keywords in dictionary.

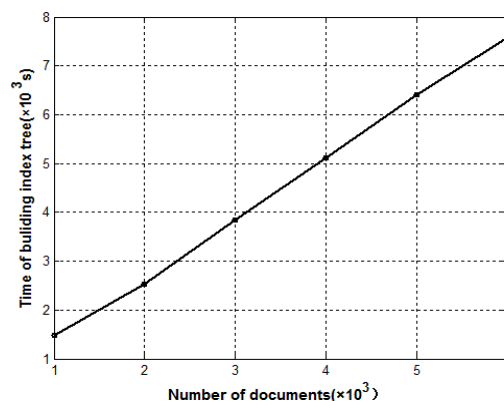
### 4.3. Search efficient

The search process, which is implemented by the cloud server, is composed by computing the similarity scores of relevant documents and result ranking based on these scores. Fig. 3(a) shows the search time for our scheme with different size of dataset. Let  $r$  represent the number of documents including the search keywords. From Fig. 3(a), we can know that the search time is mainly depends on the number of documents in the dataset when  $r$  is fixed. Fig. 3(b) shows the relationship between search time and  $r$  with same dataset, in which search time is almost linear to  $r$ .

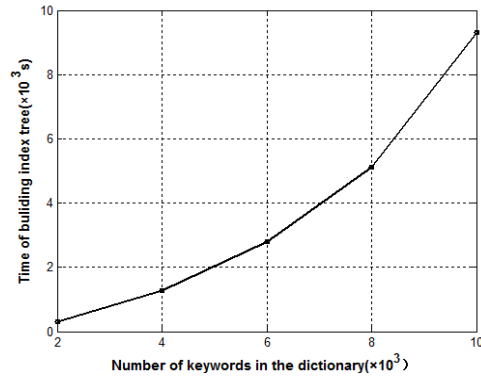
## 5 Conclusion and Future Work

In this paper, we propose secure search scheme supporting multi-keyword ranked search over encrypted cloud data. We make contributions mainly in two aspects: similarity ranked search for more accurate search result and tree-based searchable index for more efficient searching. In term of accuracy, we adopt the vector space model combined with cosine measure to evaluate the similarity between search request and document and acquire accurate search result instead of undifferentiated result. For the efficiency aspect, we propose a tree-based index structure. We propose a secure scheme to meet privacy requirements in the threat model. Finally, we analyze the performance of our scheme in detail by the experiment on real-world dataset. But, there still exist some problems, such as dynamic update for searchable index. We will do more research in the future.

**Acknowledgements.** This work is supported by the NSFC (61232016, 61173141, 61173142, 61173136, 61103215, 61373132, 61373133), GYHY201206033, 201301030, 2013DFG12860, BC2013012 and PAPD fund.



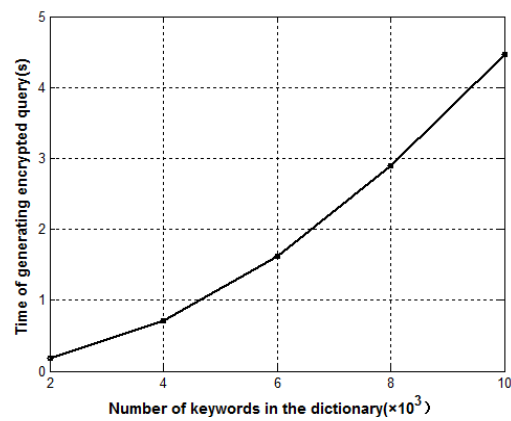




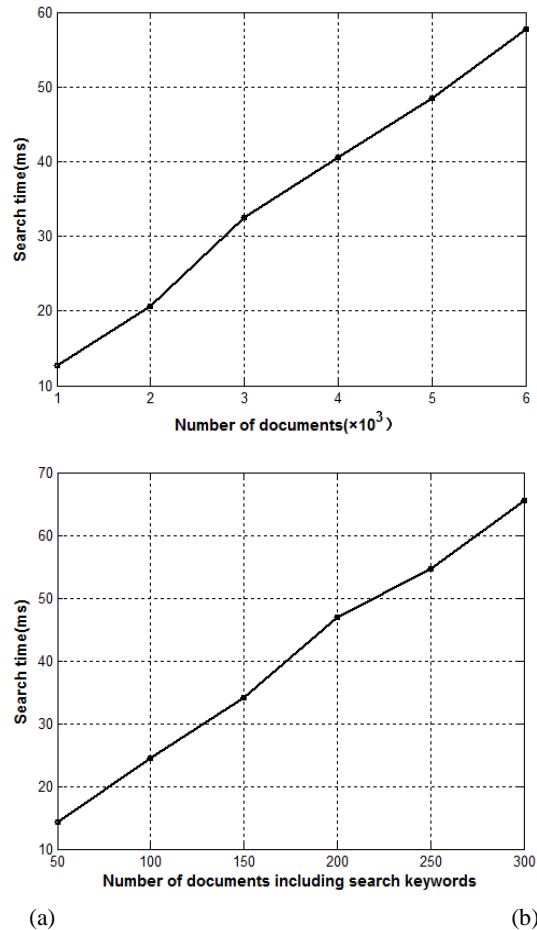
(a)

(b)

**Fig.1.** Time cost for building index tree. (a) For the different number of documents in the dataset with the same dictionary,  $n=4000$ . (b) For the different number of keywords in the dictionary with the same dataset,  $m=1000$ .



**Fig.2.** Time cost of generating encrypted query.



**Fig.3.** Search efficiency. (a) For the different size of dataset with the same number of documents including search keywords,  $r=90$ . (b) For the different number of documents including the search keywords with the same size of dataset,  $m=2000$ .

## Reference

1. L. M. Vaquero, L. Roderio-Merino, J. Caceres and M. Lindner. A break in the clouds: towards a cloud definition [J]. ACM SIGCOMM Computer Communication Review, 2009, 39(1): 50–55.
2. S. Kamara and K. Lauter. Cryptographic cloud storage[C]//Financial Cryptography and Data Security, Springer Berlin Heidelberg Publishing, 2010: 136-149.
3. C. Wang, N. Cao, K. Ren and W. Lou. Enabling secure and efficient ranked keyword search over outsourced cloud data [J]. IEEE Transactions on Parallel and Distributed Systems, 2012, 23(8): 1467–1479.
4. N. Cao, C. Wang, M. Li, K. Ren and W. Lou. Privacy-preserving multi-keyword ranked

- search over encrypted cloud data [J]. IEEE Transactions on Parallel and Distributed Systems, 2014, 25(1): 222-233.
5. W. Sun, B. Wang, N. Cao, M. Li, W. Lou, YT. Hou and H.L. Privacy-preserving multi-keyword text search in the cloud supporting similarity-based ranking[C]//Proceedings of the 8th ACM SIGSAC symposium on Information, computer and communications security, 2013: 71-82.
  6. I. H. Witten, A. Moffat and T. C. Bell. Managing gigabytes: Compressing and indexing documents and images [M]. San Francisco: Morgan Kaufmann Publishing, 1999.
  7. D.A.Grossman and O. Frieder. Information retrieval: Algorithms and heuristics[M]. Springer Publishing, 2004.
  8. W. K. Wong, D. W. Cheung, B. Kao and N. Mamoulis. Secure knn computation on encrypted databases[C]//Proceedings of SIGMOD, 2009: 139–152.
  9. RFC (Request For Comments Database) [DB/OL]. <http://www.ietf.org/rfc.html>.