

Achieving Fine-grained Multi-keyword Ranked Search over Encrypted Cloud Data

Guowen Xu[†], Hongwei Li^{† §}

[†] School of Computer Science and Engineering, University of Electronic Science and Technology of China, China

Abstract—With the advancement of Cloud computing, people now store their data on remote Cloud servers for larger computation and storage resources. However, users’ data may contain sensitive information of users and should not be disclosed to the Cloud servers. If users encrypt their data and store the encrypted data in the servers, the search capability supported by the servers will be significantly reduced because the server has no access to the data content. In this paper, we propose a Fine-grained Multi-keyword Ranked Search (FMRS) scheme over encrypted Cloud data. Specifically, we leverage novel techniques to realize multi-keyword ranked search, which supports both mixed “AND”, “OR” and “NO” operations of keywords and ranking according to the preference factor and relevance score. Security analysis indicates that the FMRS scheme can achieve the data confidentiality, the privacy protection of index and trapdoor, and the unlinkability of trapdoor. Extensive experiments using the real-world dataset demonstrate that the FMRS achieves better performance than the existing schemes in terms of functionality and efficiency.

Index Terms—Cloud computing, Multi-keyword search, Privacy-preserving

I. INTRODUCTION

Cloud computing is an emerging data storage and computing service, which is available to the public users over the Internet. It significantly saves the local data storage and computing cost of data owners [1], [2]. However, users’ data may contain sensitive information of users and should not be disclosed to the cloud server. Simply encrypting the data may prevent the server from accessing it, but it also significantly reduces the search capability of the server. In addition, there are other security concerns, like some sensitive information leakage of users.

Recently, the searchable encryption [3]–[5] has been developed as a fundamental approach to enable searching over encrypted cloud data, which proceeds as follows. Firstly, the data owner generates several keywords according to the outsourced documents. These keywords are then encrypted and delivered to the cloud server. When a search user intends to search over the outsourced documents, it can select some relevant keywords and sends the ciphertext of the selected keywords to the cloud server. The cloud server then uses the ciphertext to match the outsourced encrypted keywords, and returns the matching results to the search user at last. Sun et al. [6] propose a multi-keyword text search scheme, which builds the search index based on term frequency and the vector space model with cosine similarity measure to achieve higher search

result accuracy. To improve the search efficiency, Strizhov et al. [7] propose a tree-based Substring Position Searchable Symmetric Encryption (SSP-SSE) to handle substring search queries over encrypted data, which also involves identifying the position of the substring within the document. Li et al. [8] utilize the relevance score and k-nearest neighbor techniques to design an efficient multi-keyword search scheme, which supports the mixed “AND”, “OR” and “NO” operations of keywords. However, these proposed schemes cannot achieve the ranking according to the preference factor and relevance score simultaneously.

In this paper, we propose a Fine-grained Multi-keyword Ranked Search (FMRS) scheme over encrypted cloud data. Specifically, our original contributions can be summarized as follows:

- We utilize novel techniques to develop the multi-keyword ranked search which can achieve mixed “AND”, “OR” and “NO” operations of keywords while supporting fine-grained ranking according to the preference factor and relevance score.
- Security analysis indicates that the proposed FMRS scheme can achieve the data confidentiality, the privacy protection of index and trapdoor, and the unlinkability of trapdoor. Through extensive experiments using the real-world dataset, we show the performance of the FMRS is better than other schemes [9], [6], [10] and [8] in terms of functionality and efficiency.

The remainder of this paper is organized as follows. In Section III, we describe the preliminaries of the proposed schemes. In Section II, we outline the system model, threat model and security requirements. We present the developed scheme in Section IV. Then we carry out the security analysis and performance evaluation in Section V and Section VI, respectively. Finally, Section VIII concludes the paper.

II. SYSTEM MODEL, THREAT MODEL AND SECURITY REQUIREMENTS

A. System Model

As shown in Fig. 1, we consider a system consists of three entities.

- *Data owner*: The data owner outsources her data to the cloud server which supports reliable data access. In order to protect the privacy of data, data is usually encrypted

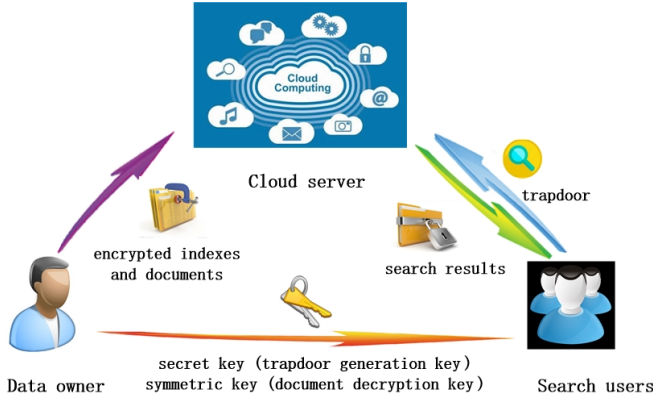


Fig. 1: System model

through symmetric encryption before being outsourced. To improve the search efficiency, the data owner needs to generate a number of keywords for each outsourced document, and create the corresponding index by the keywords and a secret key. Finally, the data owner sends the encrypted documents as well as the corresponding indexes to the cloud server, and sends the symmetric key and secret key to search user.

- *Cloud server*: The cloud server is an intermediate entity which stores the encrypted documents and corresponding indexes that are received from the data owner, and provides data access and search services to given search users. When the cloud server receives a trapdoor from the search user, it would return a collection of matching documents based on certain operation.
- *Search user*: The cloud server supports reliable data access to the corresponding search users. Generally speaking, a search user queries the documents in the cloud server with following two steps. Firstly, the search user uses the secret key and symmetric key to generate a trapdoor which is sent to the cloud server, then she receives the matching document collection from the cloud server and decrypts them with the symmetric key.

B. Threat Model and Security Requirements

In our threat models, the cloud server is generally considered “honest-but-curious”, which is the same as most related works on secure cloud data search [6], [9], [10]. Specifically, on one hand, the cloud server is honest and honestly follows the designated protocol specification and provides appropriate services. However, the cloud server is also curious, and could be “curious” to infer and analyze data (including index) stored in the cloud server so as to learn additional information. Based on this situation, we consider two threat models depending on the information available to the cloud server.

- **Known Ciphertext Model**: In this model, the cloud server only knows the encrypted document collection \mathcal{C} and the corresponding index collection \mathcal{I} , both of which are outsourced from the data owner.

- **Known Background Model**: In this stronger model, the cloud server is supposed to possess more knowledge than what can be accessed in the known ciphertext model, such as the correlation relationship of trapdoors and the related statistics of other information, i.e., cloud server can obtain a large amount of statistical information through a known database which bears the similar nature to the targeting dataset.

Based on the above threat model, we define the security requirements as follows:

- *Confidentiality of documents*: The documents of the data owner are stored in the cloud server. Due to the privacy of data, data is usually encrypted through symmetric encryption before being outsourced, so the contents of documents should not be identifiable except by the data owner and the authorized search users.
- *Privacy protection of index and trapdoor*: As discussed in Section II-A, for the convenience of users to search the encrypted cloud data, the index and the trapdoor should be created by the documents’ keywords and the search keywords, respectively. If the cloud server knows some information of index or trapdoor, she can deduce any association between keywords and encrypted documents. Therefore, the contents of index and trapdoor privacy should be ensured and cannot be identified by the cloud server.
- *Unlinkability of trapdoor*: Unlinkability of trapdoor, i.e. the cloud server can not get any keyword information according to the trapdoors. Even if some documents stored in the cloud server are searched many times, the cloud server should be unable to determine whether two trapdoors are generated from the same search request. Otherwise, the cloud server can deduce relationship of trapdoors which threatens the privacy of keywords. Therefore, for the same keywords, trapdoors should be generated randomly, rather than deterministic.

III. PRELIMINARIES

A. Notation

- \mathcal{F} —the document collection to be outsourced, denoted as a set of N documents $\mathcal{F} = (F_1, F_2, \dots, F_N)$.
- \mathcal{C} —the encrypted document collection according to \mathcal{F} , denoted as a set of N documents $\mathcal{C} = (C_1, C_2, \dots, C_N)$.
- \mathcal{FTD} —the identity collection of encrypted documents \mathcal{C} , denoted as $\mathcal{FTD} = (FID_1, FID_2, \dots, FID_N)$.
- \mathcal{W} —the keyword dictionary, including m keywords, denoted as $\mathcal{W} = (w_1, w_2, \dots, w_m)$.
- \mathcal{I} —the index stored in the cloud server, which is built from the keywords of each document, denoted as $\mathcal{I} = (I_1, I_2, \dots, I_N)$.
- \mathcal{W} —the query keyword set generated by a search user, which is a subset of \mathcal{W} .
- $T_{\mathcal{W}}$ —the trapdoor for keyword set \mathcal{W} .
- \mathcal{FTD} —the identity collection of documents returned to the search user.

B. Secure kNN Computation

We leverage the work of Wong et al. [11]. Wong et al. propose a secure k -nearest neighbor (kNN) scheme which can encrypt two vectors and calculate their Euclidean distance secretly. Firstly, the secret key (S, M_1, M_2) needs to be generated by data owner. The role of binary vector S is to split plaintext vector into two random vectors, which can change the original value of plaintext vector. Then the M_1 and M_2 are used to encrypt the split vectors. Detailed introduction for kNN can refer to the literature [11].

C. Relevance Score

The relevance score between a keyword and a document indicates the frequency of the keyword in the document. Without loss of generality, we adopt a widely used expression in [12] to evaluate the relevance score as

$$Score(\widetilde{W}, F_j) = \sum_{w \in \widetilde{W}} \frac{1}{|F_j|} \cdot (1 + \ln f_{j,w}) \cdot \ln(1 + \frac{N}{f_w}) \quad (1)$$

where $f_{j,w}$ represents the frequency of keyword w in document F_j ; f_w represents the number of documents that contain the keyword w ; N represents the number of files in the collection; and $|F_j|$ represents the length of F_j , which obtained by counting the number of indexed keywords.

D. Reference Factor

The preference factors of keywords indicate the importance of keywords in the search keyword set that are personally defined by the search user. As a search user, she can set the weight of keywords according to her own preferences. In this paper, for the convenience of deriving, the search user randomly chooses a super-increasing sequence $(a_1 > 0, a_2, \dots, a_l)$ (i.e., $\sum_{i=1}^{j-1} a_i \cdot D < a_j (j = 2, 3, \dots, l)$), where a_i is the preference factor of keyword w_i .

IV. PROPOSED SCHEME

In this section, we propose a fine-grained multi-keyword ranked search scheme (FMRS), which can support both mixed “AND”, “OR” and “NO” operations of keywords and ranking according to the preference factor and relevance score.

1) *Initialization*: The data owner randomly generates the secret key $k = (S, M_1, M_2)$, where S is a $(m+1)$ -dimensional binary vector, M_1 and M_2 are two $(m+1) \times (m+1)$ invertible matrices, respectively, m is the number of keywords in W . Then the data owner sends (k, sk) to search user through a secure channel, where sk is the symmetric key used to encrypt documents which are outsourced to the cloud server.

2) *Index building*: The data owner firstly utilizes symmetric encryption algorithm to encrypt the document collection (F_1, F_2, \dots, F_N) with the symmetric key sk , and the document collection are encrypted as $C_j (j = 1, 2, \dots, N)$, then the data owner generates an m -dimensional binary vector P_i for every document C_i , where the values of vector $P_j (j = 1, 2, \dots, N)$ determined by the $TF \times IDF$ weighting technique [11], and the j th item of P represents the relevance score of keyword w_j in the document.

The data owner extends the P to a $(m+1)$ -dimension vector P' , where $P'[m+1] = 1$. The data owner splits P' into two $(m+1)$ -dimension vectors (P_a, P_b) using the key S . i.e, if $S[j] = 0$, $P_a[i] = P_b[i] = P'[i]$, otherwise $P'[i] = P_a[i] + P_b[i]$, (the value of $P'[i]$ will be randomly split into $P_a[i]$ and P_b). Therefore, the index of encrypted document C_j can be denoted as $I_j = (P_a M_1, P_b M_2)$. Finally, the data owner sends $C_j \parallel FID_j \parallel I_j (j = 1, 2, \dots, N)$ to the cloud server.

3) *Trapdoor generation*: The search user firstly generates the keyword set \widetilde{W} , then an m -dimensional binary vector Q is generated where each bit $Q[j]$ indicates the preference factors of w_j . The preference factors of keywords indicate the importance of keywords in the search keywords set that are personally defined by the search user[2]. Q is firstly extended to $(m+1)$ dimension vector Q' , where $Q'[m+1]$ is set to $-s$ (the value of $-s$ will be defined in the following schemes in detail), then Q' is scaled by a random number $r \neq 0$ to generate $Q'' = r \cdot Q'$. After applying the same splitting and encryption processes as above, the tarpdoor $T_{\widetilde{W}}$ is generated as $(M_1^{-1} q_a, M_2^{-1} q_b)$. Finally, the search user sends $T_{\widetilde{W}}$ to the cloud server.

4) *Query*: With the index $I_j (j = 1, 2, \dots, N)$ and tarpdoor $T_{\widetilde{W}}$ the final query result is as follows

$$\begin{aligned} R_j &= I_j \cdot T_{\widetilde{W}} = (P_a M_1, P_b M_2) \cdot (M_1^{-1} q_a, M_2^{-1} q_b) \\ &= P_a \cdot q_a + P_b \cdot q_b = P' \cdot Q'' \\ &= r \cdot (P \cdot Q - s) \end{aligned} \quad (2)$$

A. Model analysis

Compared with the traditional model, we replace the values of $P[i]$ and $Q[i]$ by the relevance scores and the preference factors of keywords, respectively. For ease of calculation, the score is rounded up, i.e, $score(w_i, F_j) = \lceil 10 * score(w_i, F_j) \rceil$, and we assume that the score is not greater than D , i.e. $score(w_i, F_j) < D$. And we also assume that the keyword sets of the the “OR” “AND” and “NO” operations are $(w'_1, w'_2, \dots, w'_l)$, $(w''_1, w''_2, \dots, w''_l)$ and $(w'''_1, w'''_2, \dots, w'''_l)$, respectively, the “OR”, “AND” and “NO” operations denoted by \vee, \wedge and \neg , respectively. Here we assume that the “NO” keywords have maximum weight, “AND” second, “OR” minimum. Thus the corresponding rule can be represented as $(w'_1 \vee w'_2 \vee \dots \vee w'_l) \wedge (w''_1 \wedge w''_2 \wedge \dots \wedge w''_l) \wedge (\neg w'''_1 \wedge \neg w'''_2 \wedge \dots \wedge \neg w'''_l)$ by the ascending order of keyword weight, For “OR”, “AND” and “NO” operations, the search user chooses a super-increasing sequence $(a_1, a_2, \dots, a_{l_1}), (a_{l_1+1}, a_{l_1+2}, \dots, a_{l_1+l_2})$ and $(a_{l_1+l_2+1}, a_{l_1+l_2+2}, \dots, a_{l_1+l_2+l_3})$ $(\sum_{i=1}^{j-1} a_i \cdot D < a_j (j = 2, 3, \dots, N))$ to achieve searching with keyword weight, respectively, and we assume that $l_1 + l_2 + l_3 = N$. So according to the search keyword set $(w'_1, w'_2, \dots, w'_l, w''_1, w''_2, \dots, w''_l, w'''_1, w'''_2, \dots, w'''_l)$, the corresponding values in Q are set as $(a_1, \dots, a_{l_1}, a_{l_1+1}, \dots, a_{l_1+l_2}, a_{l_1+l_2+1}, \dots, a_{l_1+l_2+l_3})$. Other values in Q are set as 0.

We demonstrate that our model(FMRS) can achieve logical search operation as following(for the convenience of deduction, we still use w_i to replace w'_i, w''_i and w'''_i in the following paragraphs):

Step 1: A search user needs to be sure of the keywords to be searched for. In FMRS, firstly, the values of $score(w_i, F_j)$ which belongs to“NO” keyword set will be set as 1, then the search results is:

$$R_j = r \cdot (P \cdot Q - s) = r \cdot \left(\sum_{i=1}^N score(w_i, F_j) a_i - s \right) \quad (3)$$

(here s is equal to the minimum value of the “NO” keywords weight, i.e., $s = a_{l_1+l_2+1}$)

Step 2: Check whether R_j is less than 1 by computing following equation (4)

$$\begin{aligned} R_j &= r \cdot (P \cdot Q - s) \\ &= r \cdot \left(\sum_{i=1}^N score(w_i, F_j) \cdot a_i - s \right) \\ &= r \left(\sum_{i=1}^{l_1+l_2} score(w_i, F_j) a_i + \sum_{i=l_1+l_2+1}^N score(w_i, F_j) a_i - s \right) \end{aligned} \quad (4)$$

We know that $\sum_{i=1}^{j-1} a_i \cdot D < a_j (j = 2, 3 \dots N)$, if all the keywords in the “NO” keyword set are not in the keyword sets of $F_j (j = 2, 3 \dots N)$, we can infer $\sum_{i=l_1+l_2+1}^N score(w_i, F_j) \cdot a_i = 0$, therefore

$$\begin{aligned} R_j &= r \left(\sum_{i=1}^{l_1+l_2} score(w_i, F_j) \cdot a_i - s \right) \\ &= r \left(\sum_{i=1}^{l_1+l_2} score(w_i, F_j) \cdot a_i - a_{l_1+l_2+1} \right) \\ &< \sum_{i=1}^{l_1+l_2} D \cdot a_i - a_{l_1+l_2+1} \\ &< 0 \end{aligned} \quad (5)$$

In the same way, If there is a keyword belongs to “NO” keyword set and which is in the keyword sets of $F_j (j = 2, 3 \dots N)$, there must be $R_j > 0$. So, if $R_j > 0$, we choose a new R_j and return to **Step 2**. Otherwise, we go to the next step.

Step 3: Use R_j to mod $(-r \cdot a_{l_1+l_2+1}, r \cdot a_{l_1+l_2}, \dots, r \cdot a_{l_1+1})$ in turn, then check whether the quotient is over or equal to 1 each time. Besides, the remainder can't be zero.

For the first time, $R_j = r \left(\sum_{i=1}^{l_1+l_2} score(w_i, F_j) \cdot a_i - s \right) = r \cdot \sum_{i=1}^{l_1+l_2} score(w_i, F_j) \cdot a_i - r \cdot s$. Then $R_j \bmod -r \cdot a_{l_1+l_2+1} = r \cdot \sum_{i=1}^{l_1+l_2} score(w_i, F_j) \cdot a_i$. Obviously, quotient is 1, then we

go to the next step. (the purpose of this step is to eliminate the effect of s).

For the second time , the value of R_j is equal to the remainder operated last time, which is $r \cdot \sum_{i=1}^{l_1+l_2} score(w_i, F_j) \cdot a_i$, then we mod $r \cdot a_{l_1+l_2}$, here we know

$$\begin{aligned} R_j &= r \cdot \sum_{i=1}^{l_1+l_2} score(w_i, F_j) \cdot a_i \\ &= r \cdot \sum_{i=1}^{l_1+l_2-1} score(w_i, F_j) \cdot a_i + r \cdot score(w_{l_1+l_2}, F_j) \cdot a_{l_1+l_2} \end{aligned} \quad (6)$$

If the keyword $w_{l_1+l_2}$ (i.e., $w''_{l_1+l_2}$) is not in the keyword sets of $F_j (j = 2, 3 \dots N)$, then

$$\begin{aligned} R_j &= r \cdot \sum_{i=1}^{l_1+l_2-1} score(w_i, F_j) \cdot a_i \\ &< r \cdot \sum_{i=1}^{l_1+l_2-1} D \cdot a_i \\ &< r \cdot a_{l_1+l_2} \end{aligned} \quad (7)$$

The quotient of $R_j \bmod r \cdot a_{l_1+l_2} = 0$, we should also choose a new R_j and return to **Step 2**. On the contrary, if the keyword $w_{l_1+l_2}$ (i.e., $w''_{l_1+l_2}$) is in the keyword sets of F_j , $R_j \bmod r \cdot a_{l_1+l_2} \geq 1$, we go to the next step.

In a similar way, using the R_j to mod $(-r \cdot a_{l_1+l_2-1}, r \cdot a_{l_1+l_2-2}, \dots, r \cdot a_{l_1+1})$ in turn, if all the keywords in the “AND” keyword set are in the keyword sets of $F_j (j = 2, 3 \dots N)$, the quotient is ≥ 1 in each time. Besides, the final remainder can't be zero, i.e., the remainder of R_j mods $a_{l_1+1} \neq 0$. (The remainder is not 0 guarantee that at least one “OR” keyword in the document)can satisfy the above matching rule with “OR”, “AND”and “NO”.

Step 4: For all of the documents that conform to the above query scheme. Return K documents with the highest scores by equation (2).

V. SECURITY ANALYSIS

In this section, we analyze the security of the proposed model. In particular, we focus on how to achieve confidentiality of documents, privacy protection of index and trapdoors, and unlinkability of trapdoors of our proposed model. Other security features are not the key issues of our scheme.

A. Confidentiality of Documents

In FMRS, the documents of data owner are stored in the cloud server. For the consideration of privacy, data is generally encrypted by symmetric encryption (e.g., AES) before outsourcing. In addition, the secret key sk is generated by the data owner and is sent to the search user by a secure channel. AES was proved to be secure in [13]. Any entity that is unable to obtain the information or content of a document without the secret key sk . Therefore, the confidentiality of encrypted documents can be achieved.

B. Privacy Protection of Index and Trapdoor

As discussed in Section IV, for the convenience of searching, the index and the trapdoor should be created by the documents' keywords and the search keywords, respectively. All the index $I_j = (p_a M_1, p_b M_2)$ and the trapdoor $T_{\mathcal{W}} = (M_1^{-1} q_a, M_2^{-1} q_b)$ are ciphertexts of vectors (P, Q) . The secret key is $K = (S, M_1, M_2)$ that generated by data owner in our model, where S functions as a splitting indicator that divides P and Q into (p_a, p_b) and (q_a, q_b) respectively, then, we use two invertible matrices M_1 and M_2 to encrypt (p_a, p_b) and (q_a, q_b) . The security of the encryption algorithm has been proved under the known ciphertext model [11]. Thus, the contents of index and trapdoor cannot be identified. Therefore, privacy protection of index and trapdoor is supported in FMRS.

C. Unlinkability of Trapdoor

Some documents stored in the cloud server may be frequently retrieved. Unlinkability refers to the case where the cloud server can not obtain keyword information from the trapdoors. Once unlinkability of trapdoor is broken, the cloud server can deduce relationship of trapdoors, and threaten the privacy of keywords. Therefore, for the same keywords, trapdoor should be generated randomly, rather than deterministic. If the trapdoor generation algorithm is deterministic, the cloud server may reveal the relationship of keywords, even though the cloud server cannot decrypt the trapdoors. We consider whether the trapdoor $T_{\mathcal{W}} = (M_1^{-1} q_a, M_2^{-1} q_b)$ can be linked to the keywords, we prove FMRS can achieve the unlinkability of trapdoors in a strong threat model, i.e., known background model [9].

Known Background Model: In this model, the cloud server is supposed to possess more knowledge than what can be accessed in the known ciphertext model. i.e., the cloud server can obtain a large amount of statistical information through another database which has a similar nature to the targeted dataset. In our model, the trapdoor is made up of two parts. The values of $a_i (i = 1, 2, \dots, N)$ are the super-increasing sequence randomly selected by the search user (assume there are α possible sequences). And the $(m + 1)$ dimension is $-s$ defined by the search user, where the value of s is equal to the minimum value of the "NO" keyword weights. i.e., $s = a_{l_1+l_2+1}$. Assuming that the number of different $a_{l_1+l_2+1}$ is represented as β . Further, $Q'' = r \cdot Q'$, Q' is used to multiply a positive random number r , assuming that all the possible values of r is 2^{η_r} (if the search user chooses η_r -bit r). Finally, Q'' is split into (q_a, q_b) by the splitting indicator S . Specifically, if $S[i] = 0 (i = 1, 2, \dots, m + 1)$, the value of $Q''[i]$ will be randomly split into $q_a[i]$ and $q_b[i]$, assuming the number of '0' in S is μ , and each dimension with q_a and q_b is η_q bits. Note that η_s, η_r, μ and η_q are independent of each other. Then in our model, we calculate the probabilities of two trapdoors which are the same as follows:

$$P_2 = \frac{1}{\beta \cdot 2^{\eta_r} \cdot (2^{\eta_q})^\mu} = \frac{1}{\beta \cdot 2^{\eta_r + \mu \eta_q}} \quad (8)$$

Therefore, the larger β, η_r, μ and η_q can achieve the stronger security, we choose 1024-bit r , then the probability $P_1 < 1/2^{1024}$. Thus, the probabilities of two trapdoors which are the same is negligible. In summary, we present the comparison results of security level in Table I, where (FMRS) represents our model. Clearly, all the schemes can achieve confidentiality of documents and privacy protection of index and trapdoor, but the OPE schemes [14] cannot achieve the unlinkability of trapdoor very well because of the similarity relevance mentioned in [10]. Comparison with our model, the scheme [8] can not return precise results because the relevance scores is not utilized.

TABLE I: Comparison of Security Level

	[14]	[6], [9], [10]	[8]	FMRS
Confidentiality	✓	✓	✓	✓
Privacy protection	✓	✓	✓	✓
Unlinkability		✓	✓	✓

VI. PERFORMANCE EVALUATION

In this section, we analyze the performance of the model by using the method of simulation and comparison with the existing models [6], [9], [10]. We randomly select a certain number of data through a real database 1990-2003 [15], and conduct real-world experiments on an Intel Core i5 2.6 GHz system.

A. Functionality

We compare the function of [9], [6], [10] and our scheme in Table II, where (FMRS) represent our model.

MRSE [9] can achieve multi-keyword search and return the coordinate matching results by using secure kNN computation scheme. And [6] and [10] consider the relevance scores of keywords. To a certain extent, it can improve the accuracy of the search results when compared with [9]. [8] utilizes the relevance score and k-nearest neighbor techniques to design an efficient multi-keyword search scheme which supports the mixed "AND", "OR" and "NO" operations of keywords (here [8]_I represents the model one of [8], [8]_II represents the model two of [8]). However, the proposed scheme cannot achieve the ranking according to the preference factor and relevance score simultaneously. Note that if the values of all relevance scores and preference factors of keywords as the same, our model degrades to MRSE and the coordinate matching can be achieved. Besides, a series of logical operations of "OR", "AND", and "NO" can be realized according to the different choices of the search user.

B. Efficiency

1) *Computation overhead:* In order to provide a comprehensive analysis of computation overhead, we discuss it from following phases.

Index building. Note that the *Index building* phase of FMRS, which contains the relevance score computing. Considering the cost of calculating the relevance score, it is

TABLE II: Comparison of Functionalities

	[9]	[10]	[8]_I	[8]_II	FMRS
Multi-keyword search	✓	✓	✓	✓	✓
Coordinate matching	✓	✓	✓	✓	✓
Relevance score		✓	✓		✓
Preference factor			✓		✓
AND OR NO operations				✓	✓

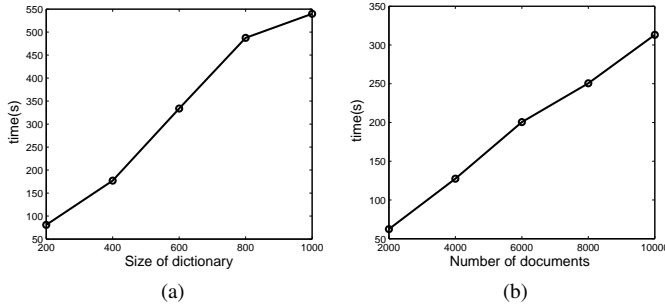


Fig. 2: Time for building index. (a) For the different size of dictionary with the same number of documents, $N=6000$. (b) For the different number of documents with the same size of dictionary, $|\mathcal{W}| = 400$.

negligible in comparison with the cost of index building, we do not distinguish them. Moreover, because the index instruction mainly involves with the two multiplications of a $(m+1) \times (m+1)$ invertible matrix and a $(m+1)$ - dimension splitting vector. as shown in Fig. 2, we can see the time of building index is significantly associated with the number of documents and dictionaries.

Trapdoor generation. In *Trapdoor generation* phase, our model randomly generates a super increasing sequence and a weight sequence, respectively, which is same as [8]. As shown in Fig. 3, the time of generating trapdoors is also significantly associated with the number of dictionaries, instead of the number of query keywords. It is partly because even if we do not want to search some keywords, we will still need to set values for the corresponding elements. With the increase in the number of the keywords in the dictionary, the time cost rises.

Query. The computation overhead in *Query* phase, as shown in Fig. 4, is significantly associated with the size of dictionary and the number of documents, instead of the number of query keywords. Note that, in *Trapdoor generation* and *Query* phases, the computation overheads are irrelative to the number of query keywords. Thus our schemes are more efficient compared with some multiple-keyword search schemes [16], [17], as their cost is linear with the number of query keywords. Besides, comparing with [8], our model returns more precise results because of using the relevance scores.

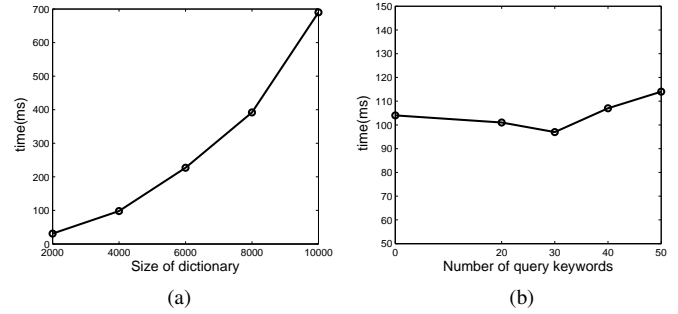


Fig. 3: Time for generating trapdoor. (a) For the different size of dictionary with the same number of query keywords, $|\mathcal{W}|=20$. (b) For the different number of query keywords with the same size of dictionary, $N = 4000$.

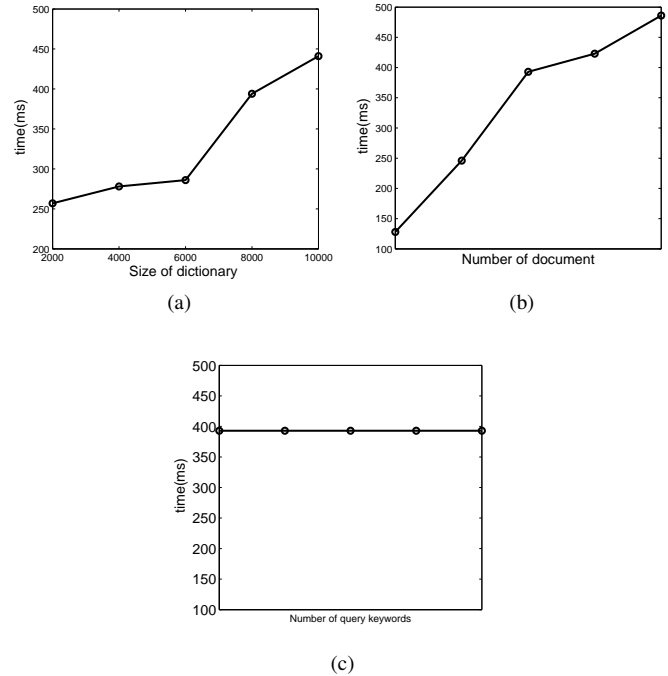


Fig. 4: Time for query. (a) For the different size of dictionary with the same number of documents and number of search keywords, $N = 6000$, $|\mathcal{W}| = 20$. (b) For the different number of documents with the same size of dictionary and number of search keywords, $|\mathcal{W}| = 8000$, $|\mathcal{W}| = 20$. (c) For the different number of search keyword with the same size of dictionary and number of documents, $N = 6000$, $|\mathcal{W}| = 8000$.

VII. RELATED WORK

Searchable encryption has been recently developed as a fundamental approach to enable searching over encrypted cloud data. Wang et al. [14] propose a ranked keyword search scheme which considers the relevance scores of keywords. However, the above schemes cannot efficiently support multi-keyword search which is widely used to provide the better experience to the search user. Later, Cao et al. [18] propose a searchable scheme which solves the problem of privacy-preserving multi-keyword query over encrypted graph-structured data in cloud computing, it establishes a set of strict privacy requirements for such a secure cloud data utilization system to become a reality. Recently, Li et al. [8] utilize the relevance score and k-nearest neighbor techniques to design an efficient multi-keyword search scheme which supports the mixed “AND”, “OR” and “NO” operations of keywords. However, the proposed scheme cannot achieve the ranking according to the preference factor and relevance score simultaneously. In this paper, we propose a fine-grained multi-keyword ranked search scheme (FMRS), which can support both mixed “AND”, “OR” and “NO” operations of keywords and ranking according to the preference factor and relevance score.

VIII. CONCLUSION

In this paper, we propose a Fine-grained Multi-keyword Ranked Search (FMRS) scheme over the encrypted cloud data. Specifically, we develop the multi-keyword ranked search to support both mixed “AND”, “OR” and “NO” operations of keywords and ranking according to the preference factor and relevance score. Security analysis indicates that FMRS scheme can preserve confidentiality of documents, privacy protection of index and trapdoor and unlinkability of trapdoor. Real-world experiments demonstrate that FMRS can achieve better performance in terms of functionality and efficiency compared to the existing proposals.

IX. ACKNOWLEDGEMENT

This work is supported by the National Natural Science Foundation of China under Grants 61472065, U1233108, U1333127, and 61272525, the International Science and Technology Cooperation and Exchange Program of Sichuan Province, China under Grant 2014HH0029, China Post-doctoral Science Foundation funded project under Grants 2014M552336 and 2015T80972, and State Key Laboratory of Information Security foundation Open Foundation under Grant 2015-MS-02.

REFERENCES

- [1] C. Wang, N. Cao, K. Ren, and W. Lou, “Enabling secure and efficient ranked keyword search over outsourced cloud data,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 23, no. 8, pp. 1467–1479, 2012.
- [2] M. Yu, K. Yang, L. Wei, and J. Sun, “Practical private information retrieval supporting keyword search in the cloud,” in *2014 Sixth International Conference on Wireless Communications and Signal Processing (WCSP)*, Oct 2014, pp. 1–6.

- [3] D. X. Song, D. Wagner, and A. Perrig, “Practical techniques for searches on encrypted data,” in *Proceedings of S&P*. IEEE, 2000, pp. 44–55.
- [4] R. Li, Z. Xu, W. Kang, K. C. Yow, and C.-Z. Xu, “Efficient multi-keyword ranked query over encrypted data in cloud computing,” *Future Generation Computer Systems*, vol. 30, pp. 179–190, 2014.
- [5] H. Li, D. Liu, Y. Dai, T. H. Luan, and X. Shen, “Enabling efficient multi-keyword ranked search over encrypted cloud data through blind storage,” *IEEE Transactions on Emerging Topics in Computing*, 2014, DOI:10.1109/TETC.2014.2371239.
- [6] W. Sun, B. Wang, N. Cao, M. Li, W. Lou, Y. T. Hou, and H. Li, “Verifiable privacy-preserving multi-keyword text search in the cloud supporting similarity-based ranking,” *IEEE Transactions on Parallel and Distributed Systems*, vol. DOI: 10.1109/TPDS.2013.282, 2013.
- [7] M. Strizhov and I. Ray, “Substring position search over encrypted cloud data using tree-based index,” in *2015 IEEE International Conference on Cloud Engineering (IC2E)*, March 2015, pp. 165–174.
- [8] H. Li, Y. Yang, T. H. Luan, X. Liang, L. Zhou, and X. S. Shen, “Enabling fine-grained multi-keyword search supporting classified sub-dictionaries over encrypted cloud data,” *IEEE Transactions on Dependable and Secure Computing*, 2015, DOI: 10.1109/TDSC.2015.2406704.
- [9] N. Cao, C. Wang, M. Li, K. Ren, and W. Lou, “Privacy-preserving multi-keyword ranked search over encrypted cloud data,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 25, no. 1, pp. 222–233, 2014.
- [10] J. Yu, P. Lu, Y. Zhu, G. Xue, and M. Li, “Towards secure multi-keyword top-k retrieval over encrypted cloud data,” *IEEE Transactions on Dependable and Secure Computing*, vol. 10, no. 4, pp. 239–250, 2013.
- [11] W. K. Wong, D. W.-I. Cheung, B. Kao, and N. Mamoulis, “Secure knn computation on encrypted databases,” in *Proceedings of SIGMOD International Conference on Management of data*. ACM, 2009, pp. 139–152.
- [12] J. Zobel and A. Moffat, “Exploring the similarity space,” in *ACM SIGIR Forum*, vol. 32, no. 1. ACM, 1998, pp. 18–34.
- [13] N. Ferguson, R. Schroepel, and D. Whiting, “A simple algebraic representation of rijndael,” in *Selected Areas in Cryptography*. Springer, 2001, pp. 103–111.
- [14] C. Wang, N. Cao, J. Li, K. Ren, and W. Lou, “Secure ranked keyword search over encrypted cloud data,” in *Proceedings of ICDCS*. IEEE, 2010, pp. 253–262.
- [15] <http://kdd.ics.uci.edu/databases/nsfabs/nsfawards.html>.
- [16] P. Golle, J. Staddon, and B. Waters, “Secure conjunctive keyword search over encrypted data,” in *Applied Cryptography and Network Security*. Springer, 2004, pp. 31–45.
- [17] D. Boneh and B. Waters, “Conjunctive, subset, and range queries on encrypted data,” in *Theory of cryptography*. Springer, 2007, pp. 535–554.
- [18] N. Cao, Z. Yang, C. Wang, K. Ren, and W. Lou, “Privacy-preserving query over encrypted graph-structured data in cloud computing,” in *Proceedings of ICDCS*. IEEE, 2011, pp. 393–402.