

Multi-keyword Similarity Search Over Encrypted Cloud Data*

Mikhail Strizhov and Indrajit Ray

Colorado State University
Fort Collins, CO, 80523, USA
{strizhov, indrajit}@CS.ColoState.EDU

Abstract. Searchable encryption allows one to upload encrypted documents on a remote *honest-but-curious* server and query that data at the server itself without requiring the documents to be decrypted prior to searching. In this work, we propose a novel secure and efficient multi-keyword similarity searchable encryption (MKSIm) that returns the matching data items in a ranked ordered manner. Unlike all previous schemes, our search complexity is sublinear to the total number of documents that contain the queried set of keywords. Our analysis demonstrates that proposed scheme is proved to be secure against adaptive chosen-keyword attacks. We show that our approach is highly efficient and ready to be deployed in the real-world cloud storage systems.

Keywords: searchable encryption, secure cloud outsourcing, multi-keyword ranked search, homomorphic encryption

1 Introduction

Cloud computing enables new types of services where the computational and network resources are available online through the Internet. One of the most popular services of cloud computing is data outsourcing. For reasons of cost and convenience, public as well as private organizations can now outsource their large amounts of data to the cloud and enjoy the benefits of remote storage. At the same time, confidentiality of remotely stored data on untrusted cloud server is a big concern. In order to reduce these concerns, sensitive data, such as, personal health records, emails, income tax and financial reports, etc. are usually outsourced in encrypted form using well-known cryptographic techniques. Although encrypted data storage protects remote data from unauthorized access, it complicates some basic, yet essential data utilization services such as plaintext keyword search. A simple solution of downloading the data, decrypting and searching locally is clearly inefficient since storing data in the cloud is meaningless unless it can be easily searched and utilized. Thus, cloud services should enable efficient search on encrypted data to provide the benefits of a first-class cloud computing environment.

* This work was partially supported by the U.S. National Science Foundation under Grant No. 0905232

Researchers have investigated this problem quite extensively in the context of encrypted documents [1–6, 9, 10, 12, 13, 16, 20, 23, 25, 26]. Solutions generally involve building an encrypted searchable index such that its content is hidden from the remote server yet allowing the corresponding documents to be searched. These solutions differ from each other mostly in terms of whether they allow single keyword search or multi-keyword search and the types of techniques they use to build the trapdoor function. A few of them, most notably [4, 5, 9, 25], allow the notion of similarity search. The similarity search problem consists of a collection of data items that are characterized by some features, a query that specifies a value for a particular feature, and a similarity metric to measure the relevance between the query and the data items. However, these techniques either do not allow searching on multiple keywords and ranking the retrieved document in terms of similarity scores or are very computationally intensive. Moreover, none of these schemes are protected against *adaptive* adversaries [12]. Taking into account large volumes of data available today, there is need in efficient methods to perform secure similarity search over encrypted data outsourced into the cloud. In this work, we propose a novel secure and efficient multi-keyword similarity searchable encryption scheme that returns the matching data items in a ranked order.

Our contributions can be summarized as follows:

1. We present a secure searchable encryption scheme that allows multi-keyword query over an encrypted document corpus and retrieves the relevant documents ranked based on a similarity score.
2. We construct the searchable encryption scheme that is CKA2-secure in the random oracle model [12, 17]. Our scheme achieves semantic security against *adaptive* adversaries that choose their search queries as a function of previously obtained trapdoors and search outcomes.
3. We present a construction that achieves the *optimal* search time. Unlike all previous schemes that are glued to the linear search complexity, our search is sublinear to the total number of documents that contain the queried set of keywords. We show that this type of searchable encryption scheme can be extremely efficient.

The rest of the paper is organized as follows: Section 2 gives an outline of the most recent related work. In Section 3 we discuss the threat model for our problem space, give an overview of the system model, notations and preliminaries and introduce the building blocks used in our solution. In Section 4, we provide the framework of the proposed searchable encryption scheme and define the necessary security terms and requirements. The security analysis and comparison to other existing schemes is given in Section 5. Finally, Section 6 is devoted for the concluding remarks.

2 Related Work

Searchable encryption has been an active research area and many quality works have been published [1–6, 9–13, 16, 20–23, 25, 26]. Traditional searchable encryption schemes usually build an encrypted searchable index such that its content is

hidden to the server, however it still allows performing document searching with given search query. Song *et al.* [23] were the first to investigate the techniques for keyword search over encrypted and outsourced data. The authors begin with idea to store a set of plaintext documents on data storage server such as mail servers and file servers in encrypted form to reduce security and privacy risks. The work presents a cryptographic scheme that enables indexed search on encrypted data without leaking any sensitive information to the untrusted remote server. Goh [16] developed a per-file Bloom filter-based secure index, which reduce the searching cost proportional to the number of files in collection. Recent work by Moataz *et al.* [21] proposed boolean symmetric searchable encryption scheme. Here, the scheme is based on the orthogonalization of the keywords according to the Gram-Schmidt process. Orencik's solution [22] proposed privacy-preserving multi-keyword search method that utilizes minhash functions. Boneh *et al.* [6] developed the first searchable encryption using the asymmetric settings, where anyone with the public key can write to the data stored remotely, but the users with private key execute search queries. The other asymmetric solution was provided by Di Crescenzo *et al.* in [11], where the authors propose a public-key encryption scheme with keyword search based on a variant of the quadratic residuosity problem.

All secure index based schemes presented so far, are limited in their usage since they support only exact matching in the context of keyword search. Wang *et al.* [25] studied the problem of secure ranked keyword search over encrypted cloud data. The authors explored the statistical measure approach that embeds the relevance score of each document during the establishment of searchable index before outsourcing the encrypted document collection. The authors propose a single keyword searchable encryption scheme using ranking criteria based on keyword frequency that retrieves the best matching documents. Cao *et al.* [9] presented a multi-keyword ranked search scheme, where they used the principle of "coordinate matching" that captures the similarity between a multi-keyword search query and data documents. However, their index structure is uses a binary representation of document terms and thus the ranked search does not differentiate documents with higher number of repeated terms than documents with lower number of repeated terms.

3 Background and Building Blocks

3.1 System & Threat Model

Consider a cloud data hosting service that involves three entities: data owner, cloud server and data user. The data owner may be an individual or an enterprise, who wishes to outsource a collection of documents $D = (D_1, D_2, \dots, D_n)$ in encrypted form $C = (C_1, C_2, \dots, C_n)$ to the cloud server and still preserve the search functionality on outsourced data. $C_i = E_S[D_i]$ is the encrypted version of the document D_i computed using a semantically secure encryption scheme E with a secret key S . To enable multi-keyword ranked search capability, the data owner constructs searchable index I that is built on m distinct keywords $K = (k_1, k_2, \dots, k_m)$ extracted from the original dataset D . Both I and C are outsourced to the cloud server. To securely search the document collection for one or more keywords $\bar{K} \in K$, the authorized

data user uses *search trapdoor* (distributed by the data owner) that generates the search request to the cloud server. Once the cloud server receives such request, it performs a search based on the stored index I and returns a ranked list of encrypted documents $L \subseteq C$ to the data user. The data user then uses the secret key S , securely obtained from the data owner, to decrypt received documents L to original view.

We assume a *honest-but-curious* model for the cloud server. The cloud server is *honest*, that is, it is always available to the data user and it correctly follows the designated protocol specification, and it provides all services that are expected. The *curious* cloud server may try to perform some additional analysis to breach the confidentiality of the stored data. In the rest of the paper, the cloud server and the adversary are the same entity. That way, the adversary has access to the same set of information as the cloud server. For this work, we are not concerned about the cloud server being able to link a query to a specific user; nor are we concerned about any denial-of-service attacks.

3.2 Notations and Preliminaries

Let $D = (D_1, D_2, \dots, D_n)$ be a set of documents and $K = (k_1, k_2, \dots, k_m)$ be the dictionary consisting of unique keywords in all documents in D , where $\forall i \in [1, m]$ $k_i \in \{0, 1\}^*$. $C = \{C_1, C_2, \dots, C_n\}$ is an encrypted document collection stored in the cloud server. I_i is a searchable index associated with the corresponding encrypted document C_i . If A is an algorithm then $a \leftarrow A(\dots)$ represents the result of applying the algorithm A to given arguments. Let R be an operational ring, we write vectors in bold, e.g. $\mathbf{v} \in R$. The notation $\mathbf{v}[i]$ refers to the i -th coefficient of \mathbf{v} . We denote the dot product of $\mathbf{u}, \mathbf{v} \in R$ as $\mathbf{u} \otimes \mathbf{v} = \sum_{i=1} \mathbf{u}[i] \cdot \mathbf{v}[i] \in R$. We use $\lfloor x \rfloor$ to indicate rounding x to the nearest integer, and $\lfloor x \rfloor, \lceil x \rceil$ (for $x \geq 0$) to indicate rounding down or up.

Cryptographic Notations. A private-key encryption scheme is a set of three polynomial-time algorithms $SKE = (\text{Gen}, \text{Enc}, \text{Dec})$ such that Gen is a probabilistic algorithm that takes a security parameter k and returns a secret key K_{secret} ; Enc is a probabilistic algorithm that takes a key K_{secret} and a message m , and outputs a ciphertext ξ ; Dec is a deterministic algorithm that takes a secret key K_{secret} and a ciphertext ξ , and outputs m if K_{secret} is the valid secret key. We say that SKE is CPA-secure if the ciphertexts it outputs do not reveal any partial information about the original plaintext to an adversary that can adaptively query an encryption oracle. We also make use of pseudo-random function (PRF), which is a polynomial-time computable function that cannot be distinguished from random functions by any probabilistic polynomial-time adversary. We refer the reader [12, 17, 19] for formal definitions of *semantic security*, CPA-security and PRFs.

We now review definitions related to homomorphic encryption. Our definitions are based on Gentry's works [14] and [15], but we slightly relax the definition of decryption correctness, to allow a negligible probability of error.

Definition 1. Homomorphic encryption: A homomorphic encryption scheme Hom consist of four algorithms:

- KeyGen : Given security parameter λ , returns a secret key sk and a public key pk .

- Enc: Given the plaintext $m \in \{0, 1\}$ and public key pk , returns ciphertext ϕ .
- Dec: Given the ciphertext ϕ and secret key sk , returns the plaintext m .
- Eval: Given public key pk , a t -input circuit C (consisting of addition and multiplication gates modulo 2), and a tuple of ciphertext $(\phi_1, \phi_2, \dots, \phi_t)$ (corresponding to the t input bits of C), returns a ciphertext ϕ .

Hom is correct for a family \mathbb{C} of circuits with $\leq t = \text{Poly}(\lambda)$ input bits if for any $C \in \mathbb{C}$ and input bits $(m_i)_{i \leq t}$, the following holds with overwhelming probability over the randomness of KeyGen and Enc:

$$\text{Dec}(sk, \text{Eval}(pk, C, (\phi_1, \phi_2, \dots, \phi_t))) = C(m_1, m_2, \dots, m_t), \quad (1)$$

where $(sk, pk) = \text{KeyGen}(\lambda)$ and $\phi_i = \text{Enc}(pk, m_i)$ for $i = 1, \dots, t$.

Hom is compact if for any circuit C with $\leq t = \text{Poly}(\lambda)$ input bits, the bit-size of the ciphertext $\text{Eval}(pk, C, (\phi_1, \phi_2, \dots, \phi_t))$ is bounded by a fixed polynomial $b(\lambda)$.

3.3 Brakerski’s Homomorphic Cryptosystem

Brakerski *et al.* [7, 8] recently proposed encryption schemes that efficiently support low-degree homomorphic computations needed for our constructions. We use the ring-LWE-based variant of Brakerski’s homomorphic cryptosystem [7] that operates over polynomial rings modulo a cyclotomic polynomial. One appealing property of Brakerski’s homomorphic cryptosystem is to use the “batching mode” where a single ciphertext represent a vector of encrypted values and single homomorphic operation on two such ciphertexts applies the homomorphic operation component-wise to the entire vector. In such way, for the cost of a single homomorphic operation we get the compute on a entire vector of encrypted plaintexts. Let $\Delta_m(x)$ be the m^{th} cyclotomic polynomial and let $R = \mathbb{Z}[x]/\Delta_m(x)$ denote the operational ring. Moreover, let $R_p = \mathbb{Z}_p[x]/\Delta_m(x)$ be our plaintext space and let $R_q = \mathbb{Z}_q[x]/\Delta_m(x)$ be our ciphertext space for some $q > p$. Here, the secret keys are vectors of elements in R_q and the plaintext comprises of elements of R_p . We present the necessary basics of ring-LWE-based homomorphic encryption scheme: key-generation, encryption and decryption mechanisms:

Hom.KeyGen: We sample a polynomial $\bar{\mathbf{s}} \in R_q$ from a Hamming weight distribution $\mathbb{HWT}(h)$. Here, h specifies the number of nonzero coefficients in $\bar{\mathbf{s}}$ and each nonzero element in $\bar{\mathbf{s}}$ is ± 1 with equal probability. The vector $\mathbf{s} = (1, \bar{\mathbf{s}}) \in R_q^2$ is the secret key. The public key is generated by sampling uniformly a polynomial \mathbf{A} from R_q and \mathbf{e} from a noise distribution. The noise distribution is set to be the zero-mean discrete Gaussian distribution $\mathbb{DG}(\sigma^2)$ with variance σ^2 . The public key is $\mathbf{P} = [\mathbf{b} \parallel -\mathbf{A}] \in R_q^2$, where $\mathbf{b} = [\mathbf{A} \cdot \bar{\mathbf{s}} + \mathbf{e}]_q$.

Hom.Enc(\mathbf{P}, \mathbf{m}): The algorithm that encrypts the message $\mathbf{m} \in R_p$ with public key \mathbf{P} . We sample a polynomial \mathbf{r} with coefficients varying in $\{0, \pm 1\}$ and $\mathbf{e} = (\mathbf{e}_0, \mathbf{e}_1)$ from the noise distribution $\mathbb{DG}^2(\sigma^2)$. The encryption of message \mathbf{m} using public key \mathbf{P} is a vector of ciphertexts $\mathbf{c} = \left[\mathbf{P}^T \mathbf{r} + \left\lfloor \frac{q}{p} \right\rfloor \mathbf{m} \right]_q$, where $\lfloor \cdot \rfloor$ is the floor function and $[\cdot]_q$ denotes reduction modulo q .

Hom.Dec(S, c): The algorithm that decrypts a ciphertext \mathbf{c} using the secret key \mathbf{s} as follows: $\mathbf{m} = \left[\left[p \cdot \frac{[(\mathbf{c}, \mathbf{s})]_q}{q} \right] \right]_p$.

Addition: $\mathbf{c}_{\text{sum}} = \mathbf{c}_1 + \mathbf{c}_2$. E_{sum} is the result noise in \mathbf{c}_{sum} and $E_{\text{sum}} \leq E_1 + E_2$, where E_1 and E_2 denote noise for \mathbf{c}_1 and \mathbf{c}_2 .

Multiplication: $\mathbf{c}_{\text{prod}} = \left[\frac{p}{q} \cdot (\mathbf{c}_1 \otimes \mathbf{c}_2) \right]$. E_{prod} here denote the maximum noise in \mathbf{c}_{prod} and $E_{\text{prod}} \leq \delta_1 + \delta_2 + \delta_3$, where $(\delta_i)_{1 \leq i \leq 3}$ is the noise inflicted by the key switching process. We refer the reader to [7, 8] for more detailed discussion on the properties of Brakerski's homomorphic cryptosystem.

3.4 Term Frequency-Inverse Document Frequency

One of the main problems of information retrieval is to determine the relevance of documents with respect to the user information needs. The most commonly used technique to represent the relevance score in the information retrieval community is Term Frequency-Inverse Documents Frequency measure [18, 27, 28]. It is computed based on two independent measures - the Term Frequency and the Inverse Document Frequency. The Term Frequency (TF) is a statistical measure that represents the frequency of repeated terms in a documents. The TF value calculates the number of times a given term appears within a single document. The Inverse Documents Frequency (IDF) is a measure of a term's importance across the whole document collection. It is defined as the logarithm of the ratio of the number of documents in a collection to the number of documents containing a given term.

We adopt the following equation for TF-IDF measure from [27]:

$$\text{Score} = \log(f_{i,j} + 1) \times \log \left(1 + \frac{n}{\sum_{k=1}^n \chi(f_{i,k})} \right) \quad (2)$$

where $f_{i,j}$ specifies the Term Frequency of term j in document D_i , n is the total number of documents in the corpus and $\sum_{k=1}^n \chi(f_{j,k})$ denotes IDF value for term j among the entire document collection D .

To provide the ranked results to user's queries we choose to use the dot product as similarity metric. We use vector space model [18], where the documents and search query are represented as high dimensional vectors. The similarity metric is measured by applying the dot product between each document vector and the search query vector as follows: $\text{dotprod}(D_i, Q) = D_i \otimes Q$, where D_i is a vector that represents i -th document and Q is a search query vector.

4 Multi-keyword Similarity Searchable Encryption (MKSIm)

Recall that we are targeting the following scenario: the data owner creates secure searchable index and sends it along with encrypted data files to the cloud server. The

index is constructed in such a way that it provides enough information to perform the search on the outsourced data, but does not give away any information about the original data. Once the server receives the index and encrypted document files, it performs a search on the index and retrieves the most relevant documents according to data user's query.

4.1 Algorithm Definitions

Definition 2. *Multi-keyword Similarity Searchable Encryption (MKSIm) scheme over the set of documents D consists of five polynomial-time algorithms, as follows:*

1. $S_1, S_2, PK, SK \leftarrow \text{Gen}(1^s)$: a probabilistic algorithm that is run by the data owner to setup the scheme. The algorithm outputs secret keys S_1, S_2, PK and SK .
2. $(I, C) \leftarrow \text{BuildIndex}(S_1, S_2, PK, D, K)$: a probabilistic algorithm run by the data owner that takes as input a collection of documents D , the keyword dictionary K and secret keys S_1, S_2 and PK , and outputs the collection of encrypted documents $C = \{C_1, C_2, \dots, C_n\}$ and the searchable index I .
3. $\Omega \leftarrow \text{MakeQuery}(S_2, PK, K, \bar{K})$: a (possibly probabilistic) algorithm run by the data user that takes an input a secret keys S_2 and PK , keyword dictionary K and multiple keywords of interest \bar{K} , and outputs the search query Ω .
4. $L \leftarrow \text{Evaluate}(I, \Omega)$: a deterministic algorithm run by the cloud server. The algorithm inputs a search query Ω and the secure index I . It outputs the sequence of identifiers $L \subseteq C$ matching the search query.
5. $D_i \leftarrow \text{Decrypt}(S_1, SK, C_i)$: a deterministic algorithm that takes as input secret keys S_1 and SK , and a ciphertext C_i and outputs a document D_i .

An index-based MKSIm scheme is correct if $\forall s \in \mathbb{N}, \forall S_1, S_2, PK, SK$ generated by $\text{Gen}(1^s), \forall D, \forall (I, C)$ output by $\text{BuildIndex}(S_1, S_2, PK, D, K), \forall \bar{K} \in K, \forall 1 \leq i \leq n$,

$$\text{Evaluate}(I, \text{MakeQuery}(S_2, PK, K, \bar{K})) = D(\bar{K}) \bigwedge \text{Decrypt}(S_1, SK, C_i) = D_i \quad (3)$$

4.2 MKSIm Security Model

Definition 3. History: Let \bar{K} be a collection of keywords of interest, D be a collection of documents and $\Omega = \{\Omega_1, \Omega_2, \dots, \Omega_q\}$ be a vector of q search queries. The history of is defined as $H(D, \Omega)$.

Definition 4. Access Pattern: Let \bar{K} be a collection of keywords of interest, D be a collection of documents and $\Omega = \{\Omega_1, \Omega_2, \dots, \Omega_q\}$ be a vector of q search queries. The access pattern induced by a q -query history $H(D, \Omega)$, is defined as follows: $\alpha(H) = (D(\Omega_1), D(\Omega_2), \dots, D(\Omega_q))$.

Definition 5. Search Pattern: Let \bar{K} be a collection of keywords of interest, D be a collection of documents and $\Omega = \{\Omega_1, \Omega_2, \dots, \Omega_q\}$ be a vector of q search queries, The search pattern of the history $H(D, \Omega)$ is a symmetric binary matrix $\sigma(H)$ such that for $1 \leq i, j \leq q$, the element in the i^{th} row and j^{th} column is 1 if $\Omega_i = \Omega_j$, and 0 otherwise.

Definition 6. Trace: Let \bar{K} be a collection of keywords of interest, D be a collection of documents and $\Omega = \{\Omega_1, \Omega_2, \dots, \Omega_q\}$ be a vector of q search queries. The trace induced by the history $H(D, \Omega)$ is a sequence $\tau(H) = (|D_1|, |D_2|, \dots, |D_n|, \alpha(H), \sigma(H))$ comprised of the lengths of document collection $D = (D_1, D_2, \dots, D_n)$, and the access and search patters induced by the history $H(D, \Omega)$.

Definition 7. Randomized query: Let $\Omega_{1 \leq i \leq q}$ be a sequence of q generated search queries with the same set of keywords \bar{K} . We say that the scheme has (q, ϵ) -randomized query if: $\forall i, j \in 1, q \Pr(\Omega_i = \Omega_j) < \epsilon$.

Definition 8. Adaptive Semantic Security: Let $MKSim = (\text{Gen}, \text{BuildIndex}, \text{MakeQuery}, \text{Evaluate}, \text{Decrypt})$ be an index-based similarity searchable encryption scheme, s be the security parameter, and $A = (A_0, \dots, A_q)$ be an adversary such that $q \in \mathbb{N}$ and $\mathbb{S} = (\mathbb{S}_0, \dots, \mathbb{S}_q)$ be a simulator. Consider the following probabilistic experiments $\text{Real}_{MKSim, A}^*(s)$ and $\text{Sim}_{MKSim, A, \mathbb{S}}^*(s)$:

$\text{Real}_{MKSim, A}^*(s)$:	$\text{Sim}_{MKSim, A, \mathbb{S}}^*(s)$:
$(D, st_A) \leftarrow A_0(1^s)$	$(D, st_A) \leftarrow A_0(1^s)$
$S_1, S_2 \leftarrow \text{Gen}(1^s)$	$(I, C, st_S) \leftarrow S_0(\tau(D))$
$(I, C) \leftarrow \text{BuildIndex}(S_1, S_2, D, K)$	$(\bar{K}_1, st_A) \leftarrow A_1(st_A, I, C)$
$(\bar{K}_1, st_A) \leftarrow A_1(st_A, I, C)$	$(\Omega_1, st_S) \leftarrow S_1(st_S, \tau(D, \bar{K}_1))$
$\Omega_1 \leftarrow \text{MakeQuery}(S_2, K, \bar{K}_1)$	for $2 \leq i \leq q$
for $2 \leq i \leq q$	$(\bar{K}_i, st_A) \leftarrow A_i(st_A, I, C, \Omega_1, \dots, \Omega_{i-1})$
$(\bar{K}_i, st_A) \leftarrow A_i(st_A, I, C, \Omega_1, \dots, \Omega_{i-1})$	$(\Omega_i, st_S) \leftarrow S_i(st_S, \tau(D, \bar{K}_1, \dots, \bar{K}_i))$
$\Omega_i \leftarrow \text{MakeQuery}(S_2, K, \bar{K}_i)$	let $\Omega = \{\Omega_1, \Omega_2, \dots, \Omega_q\}$
let $\Omega = \{\Omega_1, \Omega_2, \dots, \Omega_q\}$	output $o = (I, C, \Omega)$ and st_A
output $o = (I, C, \Omega)$ and st_A	

where st_A is the state of adversary, st_S is the state of the simulator. We say that $MKSim$ is adaptively semantically secure if for all polynomial-size adversaries $A = (A_0, \dots, A_q)$ such that q is polynomial in s , there exist a non-uniform polynomial-size simulator $\mathbb{S} = (\mathbb{S}_0, \dots, \mathbb{S}_q)$ such that for all polynomial-size \mathbb{R} :

$$|\Pr[\mathbb{R}(o, st_A) = 1] - \Pr[\mathbb{R}(\bar{o}, st_A)]| \leq \epsilon, \quad (4)$$

where $(o, st_A) \leftarrow \text{Real}_{MKSim, A}^*(s)$, $(\bar{o}, st_A) \leftarrow \text{Sim}_{MKSim, A, \mathbb{S}}^*(s)$ and the probabilities are over the coins of Gen and BuildIndex and MakeQuery .

4.3 MKSim Construction

Our searchable scheme is based on SSE-2 inverted index data construction previously introduced in [12]. We enhance SSE-2 scheme with addition of TF-IDF statistical measurement and dot product for ranked search. We show that our construction is very efficient and it achieves the same semantic security guarantees as SSE-2 scheme. Fig. 1 shows an outline of MKSim scheme.

Our searchable index consist of two main algorithms: building the lookup filter T , based on SSE-2 construction and building the TF-IDF table Φ , based on TF-IDF

word importance. We first couple the document collection D with the dictionary K to produce the lookup filter T . For each word k we add an entry in T , there value is the document identifier with the instance of word k . Note, for a given word k and the set of documents that contains the word k , we derive a label for k with j^{th} document identifier. For example, if word k is a “colorado” and there is only one document with this word, then $k||j$ is “colorado1”. We represent the family of k with matching j^{th} documents as follows: $F_k = \{k||j : 1 \leq j \leq |D(k)|\}$, where $|D(k)|$ represents the list of matching documents. For instance, if the word k = “state” exists in a set of four documents, then family F_k is {“state1”, “state2”, “state3”, “state4”}. In our construction, searching for word k becomes equal of searching for all labels in a form of $k||j$ in the family F_k .

We guard the unique number of words in each document by adopting the idea of padding the lookup filter T such that the identifier of each document appears in the same number of entries. To protect the keyword content in the table T , we use the pseudo-random permutation π with secret parameter S_2 such as $\{0, 1\}^{S_2} \times \{0, 1\}^{l+\log_2(n+\max)} \rightarrow \{0, 1\}^{l+\log_2(n+\max)}$, where \max denote the maximum number of distinct keywords in the largest document in D , n is the number of documents in D and each keyword is represented using at most l bits. Our lookup table T is $(\{0, 1\}^{l+\log_2(n+\max)} \times \{0, 1\}^{\log_2(n)} \times \{p\})$, where $p = \max * n$.

In our second step, for each distinct keyword k_j in a document D_i , we calculate the TF-IDF value using equation (2). We then construct a table Φ where each row corresponds to the document id and each column is a keyword in the dictionary K . Each cell element of table Φ contains the TF-IDF value of a keyword k_j . Unfortunately, outsourcing the table elements to the cloud leaks some important information. It is well known fact [24] that an adversary (in our case, the cloud server) may know some of keywords and their TF distributions. Using this information, an adversary can infer the keyword index or even the document content. Based on this observation, we decided to improve the security of our solution. Our table includes the set of dummy keywords Z that are added to the keyword dictionary K . This gives us the randomness that hides the original keyword distribution of TF-IDF values. Finally, we use the “batching mode” encryption of Brakerski’s homomorphic cryptosystem to protect the values of TF-IDF table Φ . We apply Brakerski’s $\text{Hom.Enc}()$ with secret PK on each row of TF-IDF table Φ .

Once both the lookup filter T and TF-IDF table Φ are constructed, we use secure symmetric encryption scheme SKE to encrypt each document D_i with secret key S_1 to form C_i . We outsource searchable index $I = (T, \Phi)$ and encryption collection $C = \{C_1, C_2, \dots, C_n\}$ to the cloud server.

Now the collection is available for selective retrieval from the cloud server. To search for keywords of interest \bar{K} , the data user uses the trapdoor to output the search query Ω to the cloud server. The trapdoor utilizes the pseudo-random permutation π with secret parameter S_2 and $\text{Hom.Enc}()$ with secret PK to form the search query Ω . The server then locates the document identifiers id in the filter table T that matches the keywords of interest. For each encrypted document C_j , where $1 \leq j \leq id$, the cloud server executes the dot product between the search query and the values in TF-IDF table Φ to form the set of polynomials $\{\text{score}_{C_1}, \text{score}_{C_2}, \dots, \text{score}_{C_{id}}\}$.

$\text{Gen}(1^s)$: generate $S_1 \leftarrow \text{SKE.Gen}(1^s)$ and $S_2 \xleftarrow{R} \{0, 1\}^s$. Sample $\mathbf{SK}, \mathbf{PK} \leftarrow \text{Hom.KeyGen}()$.
Output S_1, S_2, \mathbf{SK} and \mathbf{PK} .
 $\text{BuildIndex}(S_1, S_2, \mathbf{PK}, D, K)$:
Initialization:

1. Scan document corpus D , extract $k_1, k_2, \dots, k_p \leftarrow$ from D_i .
2. Construct dictionary K with dummy Z .
3. For each $k \in K$, build $D(k)$ (i.e., the sequence of documents with k).

Build lookup filter T :

1. for each $k_i \in K$:
 - for $1 \leq j \leq |D(k_i)|$:
 - value = $\text{id}(D_{i,j})$, where $\text{id}(D_{i,j})$ is the j^{th} identifier in $D(k_i)$.
 - set $T[\pi_{S_2}(k_i||j)] = \text{value}$.
2. let $\bar{p} = \sum_{k_i \in K} |D(k_i)|$. If $\bar{p} < p$, assign value = $\text{id}(D)$ for all $D_i \in D$ for exactly max entries, set the address to random values.

Build TF-IDF table Φ :

1. for each $D_i \in D$
 - for each $k_j \in K$
 - $z_j \leftarrow \text{TFIDF}(k_j)$ (i.e., calculate TF-IDF value for keyword k_j).
 - $\Phi_i = \text{Hom.Enc}(\mathbf{PK}, \mathbf{z})$.

Output:

1. for each $D_i \in D$, let $C_i \leftarrow \text{SKE.Enc}(S_1, D_i)$.
2. output (I, C) , where $I = (T, \Phi)$ and $C = \{C_1, C_2, \dots, C_n\}$.

 $\text{MakeQuery}(S_2, \mathbf{PK}, K, \bar{K})$:

1. for each $k_i \in \bar{K} \in K$, $t_i = (\pi_{S_2}(k_i||1), \dots, \pi_{S_2}(k_i||n))$, $\mathbf{b}_i \leftarrow \text{TFIDF}(k_i)$.
2. $\mathbf{x} = \text{Hom.Enc}(\mathbf{PK}, \mathbf{b})$.
3. output: $\Omega = (t, \mathbf{x})$.

 $\text{Evaluate}(I, \Omega)$:

1. $\text{id} \leftarrow T[t]$.
2. for all $1 \leq j \leq \text{id}$, $\text{score}_{C_j} \leftarrow \Phi_{C_j} \otimes \mathbf{x}$.
3. output $\{\text{score}_{C_1}, \text{score}_{C_2}, \dots, \text{score}_{C_{\text{id}}}\}$.

 $\text{Decrypt}(S_1, \mathbf{SK}, C_i)$:

1. $\text{score}_{D_i} = \text{Hom.Dec}(\mathbf{SK}, \text{score}_{C_i})$.
2. $(\text{score}_{D_i})_m \xleftarrow{\text{top-}m} \{\text{score}_{D_1}, \text{score}_{D_2}, \dots, \text{score}_{D_{\text{id}}}\}$.
3. output $(D_i)_m \leftarrow \text{SKE.Dec}(S_1, (C_i)_m)$.

Fig. 1. MKSim Construction

The data user decrypts these polynomials using $\text{Hom.Dec}()$ with secret \mathbf{SK} and then he retrieves top- m documents with highest output scores.

5 Security Analysis and Complexity Results

5.1 Security

Theorem 1. *If $\text{MKSim} = (\text{Gen}, \text{BuildIndex}, \text{MakeQuery}, \text{Evaluate})$ is a index-based searchable encryption scheme, then it is adaptively semantically secure.*

Proof. We are going to describe a polynomial-size simulator $\mathbb{S} = \{\mathbb{S}_0, \mathbb{S}_1, \dots, \mathbb{S}_q\}$ such that for all polynomial-size adversaries $A = \{A_1, A_2, \dots, A_q\}$, the outputs of $\text{Real}_{\text{MKSIm}, A}^*(s)$ and $\text{Sim}_{\text{MKSIm}, A, S}^*(s)$ are computationally indistinguishable. Consider the simulator $\mathbb{S} = \{\mathbb{S}_0, \mathbb{S}_1, \dots, \mathbb{S}_q\}$ that adaptively generates the output $o^* = (I^*, C^*, \Omega^*)$ as follows:

- $\mathbb{S}_0(1^s, \tau(D))$: the simulator has a knowledge of history H that includes the number and the size of the documents. \mathbb{S}_0 start with generating $C_i^* \xleftarrow{R} \{0, 1\}^{|\mathcal{D}_i|}$ where $i \in \{1, 2, \dots, n\}$ and index $I^* = (T^*, \Phi^*)$. Here $T^* \xleftarrow{R} \{0, 1\}^{l + \log_2(n + \max)}$ is a lookup filter, and $\Phi^* \xleftarrow{R} \{0, 1\}^K$ is TF-IDF table. \mathbb{S}_0 now includes I^* in st_S and outputs (I^*, C^*, st_S) . Since st_A does not include secrets S_1 and PK , I^* is indistinguishable from the real index. Otherwise \mathbb{S}_0 can distinguish between the output of pseudo-random permutation π and a random values of size $l + \log_2(n + \max)$ and K . At the same time, st_A does not include secret S_2 , thus the output C^* is indistinguishable from the real ciphertext
- $\mathbb{S}_1(st_S, \tau(D, \Omega_1))$: now the simulator \mathbb{S}_1 has a knowledge of all document identifiers corresponding to the search query. However the search query does not disclose its structure and the content. Recall that $\mathcal{D}(\Omega_1)$ is the set of all matching document identifiers. For all $1 \leq j \leq |\mathcal{D}(\Omega_1)|$, the simulator first makes an association between each document identifier $\text{id}(\mathcal{D}_j)$ and a generated search query such that $(\mathcal{D}(\Omega_1)_i, I_i^*)$ are pairwise distinct. \mathbb{S}_1 then creates $\Omega_1^* = (t^*, \chi^*)$, where $t^* \xleftarrow{R} (\text{id}(\mathcal{D}_1), \dots, \text{id}(|\mathcal{D}(\Omega_1)|))$ and $\chi^* \xleftarrow{R} \{0, 1\}^K$. \mathbb{S}_1 stores the association between Ω_1^* and Ω_1 in st_S , and outputs (Ω_1^*, st_S) . Since st_A does not include secret S_1 , the output t_1^* is indistinguishable from the real generated query t_1 , otherwise one could distinguish between the output of π and a random string of size $l + \log_2(n + \max)$. Similarly, since st_A does not include secret PK , the output χ^* is indistinguishable from the real χ , otherwise one could distinguish between the output of $\text{Hom.Enc}()$ and a random string of a size K . Thus, Ω_1^* is indistinguishable from Ω_1 .
- $\mathbb{S}_i(st_S, \tau(D, \Omega_1, \Omega_2, \dots, \Omega_q))$ for $2 \leq i \leq q$: first \mathbb{S}_i checks if the search query Ω_i was executed before, that is, if it appeared in the trace $\sigma[i, j] = 1$, where $1 \leq j \leq i-1$. If $\sigma[i, j] = 0$, the search query has not appeared before and \mathbb{S}_i generates the search query as \mathbb{S}_1 . If $\sigma[i, j] = 1$, then \mathbb{S}_i retrieves previously searched query, and constructs Ω_i^* . \mathbb{S}_i outputs (Ω_i^*, st_S) , where Ω_i^* is indistinguishable from real Ω_i . The final output $\Omega = (\Omega_1^*, \dots, \Omega_q^*)$ is indistinguishable from generated query $(\Omega = \Omega_1, \dots, \Omega_q)$, and outputs of experiments $\text{Sim}_{\text{MKSIm}, A, S}^*(s)$ and $\text{Real}_{\text{MKSIm}, A}^*(s)$ are indistinguishable.

□

Theorem 2. *Injection of dummy keywords provides randomized search queries.*

Proof. Let us consider two search queries Ω_1 and Ω_2 , both constructed from the same keyword set K and a randomly chosen set of dummy keywords Z_1 and Z_2 from a dictionary \mathcal{Z} of a size $n = |\mathcal{Z}|$. We are aiming to prove that: $\forall i, j \quad i \neq j \quad \Pr(\Omega_i = \Omega_j) < \epsilon$ where ϵ tends to zero as n increases.

We first estimate the probability $\Pr(k)$ that the intersection \mathbf{Z} of two sets $Z_1, Z_2 \subseteq \mathcal{X}$ is equal to some value k . We have:

$$\Pr(k) = \frac{\text{\#of ways of choosing } Z_1, Z_2 \text{ with } |\mathbf{Z}| = k}{\text{\# of ways of choosing } Z_1, Z_2} \quad (5)$$

Note, there are $\binom{n}{k}$ choices for \mathbf{Z} . If Z_1 has size k , then there is one choice for Z_1 , and we can choose Z_2 arbitrarily from 2^{n-k} possibilities. If Z_1 has size $k+1$, then there are $\binom{n-k}{1}$ choices for Z_1 and 2^{n-k-1} choices for Z_2 . Let $m = n - k$, then there are $\sum_{j=0}^m 2^{m-j} \binom{m}{j} = 3^m$ possible choices for Z_1, Z_2 with intersection \mathbf{Z} . Thus, there are $3^{n-k} \binom{n}{k}$ possible ways of choosing the subsets. Since there are 4^n ways of choosing any two subsets of \mathcal{X} , we have the following: $\Pr(k) = \frac{3^{n-k} \binom{n}{k}}{4^n}$. We now evaluate $\Pr(k)$ with input $k \rightarrow 0$: $\lim_{k \rightarrow 0} \Pr(k) = \lim_{k \rightarrow 0} \frac{3^{n-k} \binom{n}{k}}{4^n} = \left(\frac{3}{4}\right)^n$. As n increases, $\lim_{k \rightarrow 0} \Pr(k) \rightarrow 0$ and hence Theorem 2 is preserved. \square

5.2 Complexity

We compare MKSim scheme with previous searchable encryption solutions in Table 1. Our comparison is based on few simple metrics: matching technique, query randomization, security notions and search complexity. We use security notations from [12]. Note that all previous work able to achieve the linear search complexity within the total number of documents in the collection. In contrast, the search in our solution is proportional to the number of files that contain a certain set of keywords.

Table 1. Comparison of several searchable encryption schemes. n is the document collection, $|n|$ denotes its bit length, $\#n$ is the number of files in the collection n , $\#n_{\bar{k}}$ is the number of files that contain keywords of interest \bar{k} .

Scheme	Matching	Query randomization	Security	Search complexity
Song <i>et al.</i> [23]	Exact	no	CPA	$O(n)$
Goh <i>et al.</i> [16]	Exact	no	CKA1	$O(\#n)$
Cao <i>et al.</i> [9]	Similarity	yes	CKA1	$O(\#n)$
Moataz <i>et al.</i> [21]	Exact	yes	CKA2	$O(\#n)$
Curtmola <i>et al.</i> [12]	Exact	no	CKA2	$O(\#n_{\bar{k}})$
Orencik <i>et al.</i> [22]	Exact	no	CKA2	$O(\#n)$
MKSim	Similarity	yes	CKA2	$O(\#n_{\bar{k}})$

6 Conclusion

Searchable encryption is a technique that enables secure searches over encrypted data stored on remote servers. We define and solve the problem of multi-keyword

ranked search over encrypted cloud data. In particular, we present an efficient similarity searchable encryption scheme that supports multi-keyword semantics. Our solution is based two building blocks: Term Frequency - Inverse Document Frequency (TF-IDF) measurement and ring-LWE-based variant of homomorphic cryptosystem. We use the dot product to quantitatively evaluate similarity measure and rank the outsourced documents with their importance to the search query. We show that our scheme is adaptive semantically secure against adversaries and able to achieve optimal sublinear search time. As future work, we plan to optimize the index construction algorithm and continue to research on usable and secure mechanisms for the effective utilization over outsourced cloud data.

References

1. M. Abdalla, M. Bellare, D. Catalano, E. Kiltz, T. Kohno, T. Lange, J. Malone-Lee, G. Neven, P. Paillier, and H. Shi. Searchable encryption revisited: consistency properties, relation to anonymous ibe, and extensions. *Journal of Cryptology*, 21(3):350–391, 2008.
2. R. Agrawal, J. Kiernan, R. Srikant, and Y. Xu. Order-preserving encryption for numeric data. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, Paris, France, June 2004.
3. M. Bellare, A. Boldyreva, and A. O’Neill. Deterministic and efficiently searchable encryption. In *Proceedings of the 27th Annual International Cryptology Conference*, Santa Barbara, CA, USA, Aug. 2007.
4. M. Blanton. Achieving full security in privacy-preserving data mining. In *Proceedings of the 3rd IEEE International Conference on Privacy, Security, Risk and Trust*, Boston, MA, USA, Oct. 2011.
5. M. Blanton, M. J. Atallah, K. B. Frikken, and Q. M. Malluhi. Secure and efficient outsourcing of sequence comparisons. In *Proceedings of the 17th European Symposium on Research in Computer Security*, Pisa, Italy, Sept. 2012.
6. D. Boneh and B. Waters. Conjunctive, subset, and range queries on encrypted data. In *Proceedings of the 4th IACR Theory of Cryptography Conference*, Amsterdam, The Netherlands, Feb. 2007.
7. Z. Brakerski. Fully homomorphic encryption without modulus switching from classical gapsvp. In R. Safavi-Naini and R. Canetti, editors, *Advances in Cryptology - CRYPTO 2012*, volume 7417 of *Lecture Notes in Computer Science*, pages 868–886. Springer Berlin Heidelberg, 2012.
8. Z. Brakerski, C. Gentry, and V. Vaikuntanathan. Fully homomorphic encryption without bootstrapping. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, Cambridge, MA, USA, 2012.
9. N. Cao, C. Wang, M. Li, K. Ren, and W. Lou. Privacy-preserving multi-keyword ranked search over encrypted cloud data. In *Proceedings of the 30th IEEE International Conference on Computer Communications*, Shanghai, China, Apr. 2011.
10. Y.-C. Chang and M. Mitzenmacher. Privacy preserving keyword searches on remote encrypted data. In *Proceedings of the 3rd International Conference on Applied Cryptography and Network Security*, New York, NY, USA, June 2005.
11. G. D. Crescenzo and V. Saraswat. Public key encryption with searchable keywords based on jacobi symbols. In K. Srinathan, C. P. Rangan, and M. Yung, editors, *Proceedings of the 8th International Conference on Cryptology in India*, volume 4859 of *Lecture Notes in Computer Science*, Chennai, India, December 2007. Springer.

12. R. Curtmola, J. Garay, S. Kamara, and R. Ostrovsky. Searchable symmetric encryption: Improved definitions and efficient constructions. In *Proceedings of the 13th ACM Conference on Computer and Communications Security*, Alexandria, VA, October 2006.
13. S. D. C. di Vimercati, S. Foresti, S. Jajodia, S. Paraboschi, and P. Samarati. Private data indexes for selective access to outsourced data. In *Proceedings of the 10th Annual ACM Workshop on Privacy in the Electronic Society*, Chicago, IL, USA, Oct. 2011.
14. C. Gentry. *A fully homomorphic encryption scheme*. PhD thesis, Stanford University, 2009. crypto.stanford.edu/craig.
15. C. Gentry. Fully homomorphic encryption using ideal lattices. In *Proceedings of the 41st Annual ACM Symposium on Theory of Computing*, Bethesda, MD, 2009.
16. E.-J. Goh. Secure indexes. Cryptology ePrint Archive, Report 2003/216, 2003. <http://eprint.iacr.org/2003/216/>.
17. O. Goldreich. *The foundations of cryptography - volume 2, basic applications*. Cambridge University Press, 2004.
18. D. Grossman and O. Frieder. *Information retrieval: algorithms and heuristics*. Kluwer international series on information retrieval. Springer, 2004.
19. J. Katz and Y. Lindell. *Introduction to Modern Cryptography (Chapman & Hall/Crc Cryptography and Network Security Series)*. Chapman & Hall/CRC, 2007.
20. M. Kuzu, M. S. Islam, and M. Kantarcioglu. Efficient similarity search over encrypted data. In *Proceedings of the 28th IEEE International Conference on Data Engineering*, Washington, DC, USA, Apr. 2012.
21. T. Moataz and A. Shikfa. Boolean symmetric searchable encryption. In *Proceedings of the 8th ACM SIGSAC Symposium on Information, Computer and Communications Security*, Hangzhou, China, May 2013.
22. C. Orencik, M. Kantarcioglu, and E. Savas. A practical and secure multi-keyword search method over encrypted cloud data. In *Proceedings of the 6th IEEE International Conference on Cloud Computing*, Santa Clara, CA, USA, June 2013.
23. D. X. Song, D. Wagner, and A. Perrig. Practical techniques for searches on encrypted data. In *Proceedings of the 2000 IEEE Symposium on Security and Privacy*, Berkeley, CA, USA, May 2000.
24. A. Swaminathan, Y. Mao, G.-M. Su, H. Gou, A. L. Varna, S. He, M. Wu, and D. W. Oard. Confidentiality-preserving rank-ordered search. In *Proceedings of the ACM Workshop on Storage Security and Survivability*, Alexandria, VA, 2007.
25. 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*, 23(8):1467–1479, 2012.
26. C. Wang, K. Ren, S. Yu, and K. M. R. Urs. Achieving usable and privacy-assured similarity search over outsourced cloud data. In *Proceedings of the 31st Annual IEEE International Conference on Computer Communications*, Orlando, FL, USA, Mar. 2012.
27. I. H. Witten, A. Moffat, and T. C. Bell. *Managing gigabytes : compressing and indexing documents and images*. Morgan Kaufmann, San Francisco, CA, 2 edition, 1999.
28. J. Zobel and A. Moffat. Exploring the similarity space. *SIGIR FORUM*, 32:18–34, 1998.