

Efficient Retrieval over Documents Encrypted by Attributes in Cloud Computing

Na Wang, Junsong Fu, Bharat K. Bhargava, *Fellow, IEEE*, Jiwen Zeng

Abstract—Secure document storage and retrieval is one of the hottest research directions in cloud computing. Though many searchable encryption schemes have been proposed, few of them support efficient retrieval over the documents which are encrypted based on their attributes. In this paper, a hierarchical attribute-based encryption scheme is first designed for a document collection. A set of documents can be encrypted together if they share an integrated access structure. Compared with the ciphertext-policy attribute-based encryption (CP-ABE) schemes, both the ciphertext storage space and time costs of encryption/decryption are saved. Then, an index structure named attribute-based retrieval features (ARF) tree is constructed for the document collection based on the TF-IDF model and the documents' attributes. A depth-first search algorithm for the ARF tree is designed to improve the search efficiency which can be further improved by parallel computing. Except for the document collections, our scheme can be also applied to other datasets by modifying the ARF tree slightly. A thorough analysis and a series of experiments are performed to illustrate the security and efficiency of the proposed scheme.

Index Terms—Cloud computing, document retrieval, file hierarchy, attribute-based encryption.

I. INTRODUCTION

MORE and more people and enterprises are motivated to outsource their local document management systems to the cloud which is a promising information technique (IT) to process the explosive expanding of data [1]. Cloud computing can collect and reorganize a huge amount of IT resources and apparently, the cloud servers can provide more secure, flexible, various, economic and personalized services compared with the local servers. Despite the advantages of cloud services, leaking the sensitive information, such as personal information, company financial data and government documents, to the public is a big threat to the data owners. In addition, to make full use of the data on the cloud, the data users need to access them flexibly and efficiently. Consequently, a huge challenge of outsourcing the data to the cloud is how to protect the confidentiality of the data properly while maintaining their searchability.

An intuitive approach is encrypting the documents first and then outsourcing the encrypted documents to the cloud. A large

number of searchable document encryption schemes have been proposed in the literatures, including single keyword Boolean search schemes [2]–[6], single keyword ranked search schemes [7]–[9] and multi-keyword Boolean search schemes [10]–[14]. However, all these schemes cannot support effective, flexible and efficient document search because of their simple functionalities. Privacy-preserving multi-keyword ranked document search schemes [15]–[18] are more promising and practical. However, all the documents in these schemes are organized by a coarse-grained access control mechanism, i.e., each authorized data user can access all the encrypted documents. As an example, the whole *IEEE Xplore Digital Library* can be accessed by all the authorized organizations (e.g., the universities) at present and this cannot satisfy the data owners and users in the future.

In this paper, a new situation is considered. A data user may want to access part of the library (e.g., computers and data related papers) and intuitively she wants to pay less money compared with the data users who want to access the whole library. In other words, in the document collection, each document can be accessed only by a set of specific data users. In this case, we need to design a fine-grained access control mechanism for the documents and it is more reasonable compared with the present method.

To make the data users able to access part of *IEEE Xplore Digital Library* on demands, a possible approach is encrypting the documents through attribute-based encryption (ABE) schemes [19], [20] before outsourcing them to the cloud. Meanwhile, the authorized data users are assigned with a set of attributes. A data user can decrypt a file if and only if her attributes match the file's attributes. Recently, ciphertext-policy attribute-based encryption (CP-ABE) [21]–[27] is a hot research area and it can provide fine-grained, one to many and flexible access control. In these schemes, each document is encrypted individually and their encryption efficiency can be improved by employing hierarchical attribute-based encryption schemes [28]–[31]. However, these schemes cannot be employed directly to solve our problem properly. First, most existing schemes focus on encrypting a single access tree. However, it is impossible that all the documents in *IEEE Xplore Digital Library* share a single access tree and how to construct a set of optimized access trees for the document collection is a huge challenge. Second, in most existing schemes, when the documents are mapped to a set of shared access trees, the data users need to store a large number of secret keys which will be analyzed in Section IV.B. Apparently, this is a heavy burden for the data users especially for an extremely large document collection and how

This research is supported by National Natural Science Foundation of China (Grant No.11261060).

N. Wang and J.W. Zeng are with the School of Mathematical Science, Xiamen University, Xiamen, 361005, China, e-mail: (wangna@stu.xmu.edu.cn, jwzeng@xmu.edu.cn).

J.S. Fu is with the School of Electronic and Information Engineering, Beijing Jiaotong University, Beijing, 100044, China, e-mail: (12120067@bjtu.edu.cn).

B.K. Bhargava is with the Department of Computer Science, Purdue University, West Lafayette, IN 47906, USA, e-mail: (bbshail@purdue.edu).

to decrease the amount of secret keys for the data users is another challenge. Except for access control, document search efficiency is also a challenge for a large document collection. To our knowledge, most existing schemes cannot support time-efficient retrieval over the documents which are organized under attribute-based access control mechanism.

To support the previously discussed service, we first design an algorithm to generate hierarchical access trees for the document collection. The proposed algorithm employs the greedy strategy to build the access trees incrementally and each access tree grows by continuously splitting the nodes in the tree. Then we design a ciphertext-policy attribute-based hierarchical document collection encryption scheme called CP-ABHE. In the proposed scheme, a set of documents can share a same integrated access tree and be encrypted together rather than being encrypted individually. In this way, both the ciphertext storage space and time costs of the encryption/decryption are saved. The security of the proposed scheme is proved theoretically and its effectiveness is also evaluated by simulation.

To support accurate and efficient document search over the encrypted documents, a complicated index structure is then constructed for the document collection. We first map the documents to document vectors based on the TF-IDF model and, in addition, the attributes of the documents are also taken into consideration. The similarity function between the document vectors is carefully designed and the vectors are organized based on their relative similarities in the ARF tree. Specifically, the similar vectors compose micro clusters which are then aggregated with each other to generate macro clusters until all the vectors belong to one cluster. The ARF vectors of the nodes in the tree are used to describe the inherent properties of the clusters represented by the nodes. At last, a depth-first search algorithm for the ARF tree is designed to guarantee both the search efficiency and accuracy.

The main contributions of this paper are summarized as follows:

- A practical hierarchical attribute-based document collection encryption scheme is proposed in which the documents are organized and controlled based on attributes. The proposed scheme can greatly decrease the storage and computing burdens.
- We map the documents to vectors in which both the keywords and associated attributes are considered. The ARF tree is proposed to organize the document vectors and support time-efficient document retrieval. In addition, a depth-first search algorithm is designed.
- A thorough simulation is performed to illustrate the security, efficiency and effectiveness of our scheme. Specifically, the proposed encryption scheme performs very well in both time and storage efficiency. In addition, our scheme also provides efficient and accurate document retrieval method.

The rest of this paper is organized as follows: The related work is provided in Section II. In Section III, we state the problem and present some preliminary techniques. The hierarchical attribute-based document encryption scheme is designed in Section IV and we present the time efficient document retrieval approach based on the ARF tree in Section V. The security

and efficiency analysis of our scheme is provided in Section VI and we further evaluate the performance of the proposed approach by experiments in Section VII. At last, Section VIII concludes this paper.

II. RELATED WORK

Our approach is mainly related with two research fields of cloud computing, i.e., ciphertext-policy attribute-based document encryption and encrypted document retrieval. The related work in these two fields is provided in the following.

Since Sahai et al. proposed the identity-based encryption (IBE) scheme [19], many ABE schemes [20], [32]–[34] have been proposed in which CP-ABE schemes [21]–[26], [35] are very promising because of their flexibility and scalability. In these CP-ABE schemes, the documents with different access structures need to be encrypted individually. To improve the encryption/decryption efficiency and scalability hierarchical attribute-based encryption has been widely researched [28]–[31] in which a set of documents may share a common access structure and can be encrypted together. Wang et al. propose a hierarchical attribute-based encryption scheme named FH-CP-ABE [28] and have proved its security theoretically. An advantage of the scheme is that the data users can decrypt all the authorized documents by computing the secret key once. Therefore, both the time costs of encryption and decryption are saved. Wang et al. design a scheme named HABE [29] with the traits of high performance, fine-grained access control, scalability and full delegation. HABE is a combination of hierarchical identity-based encryption and CP-ABE. Wan et al. propose hierarchical attribute-set-based encryption scheme (HASBE) [30] by extending ciphertext-policy attribute-set-based encryption (ASBE) with a hierarchical structure of the data users. The HASBE scheme can be seamlessly incorporated with a hierarchical structure of system users by applying a delegation algorithm to ASBE. Deng et al. extend ABE to CP-HABE [31] to support hierarchically distributing and delegating the secret keys which can be used in large organizations. Guo et al. [36] propose a resilient-leakage hierarchical attribute-based encryption scheme to defend against the auxiliary input leakage attack and the security of the scheme is detailedly analyzed.

In addition to encrypting the documents, we also attempt to search the encrypted document efficiently and accurately. Consequently, multi-keywords ranked document retrieval over encrypted document collections is also strongly related with our scheme. In [17], Cao et al. first propose a basic privacy-preserving multi-keyword ranked search scheme based on secure kNN algorithm [37]. A set of strict privacy requirements are established and then two schemes are proposed to improve the security and search experience. However, an apparent drawback of this scheme is that the search efficiency is linear with the cardinality of the document collection and consequently, it cannot be used to process extremely large document databases. Xia et al. [18] design a keyword balanced binary (KBB) tree to organize the document vectors and propose a “Greedy Depth-First Search” algorithm to improve the search efficiency. Moreover, the index tree can be updated

dynamically with an acceptable communication burden. However, the document vectors are chaotically organized in the tree and the search efficiency can be further improved. Chen et al. [15] take the relationships of documents into consideration and a hierarchical-clustering-based index structure is designed to improve the search efficiency. In addition, a verification scheme is also integrated into their scheme to guarantee the correctness of the results. Though the index structure can obtain sub-linear search efficiency, it cannot return the accurate search results. Fu et al. [16] present a personalized multi-keyword ranked search scheme in which an interest model of the data users is integrated into the document retrieval system to support personalized search and improve users' search experience. Specifically, the interest model of a data user is built based on her search history with the help of WordNet [38] in order to depict her behaviors in fine grit level. However, this scheme cannot support dynamic update operations, because the document vectors are constructed based on the statistical information of all the documents in the collection. In addition, though a MDB-tree is employed to improve the search efficiency, the effectiveness of the tree is hard to predict. Li et al. [39] propose a new attribute-based encryption scheme (KSF-OABE) which can implement keyword search function. Though the design goal of KSF-OABE is some similar with our scheme, it cannot hierarchically encrypt a document collection and support efficient multi-keyword document retrieval.

III. PROBLEM STATEMENT AND PRELIMINARIES

In this section, we state the problem and provide the related preliminary techniques. For convenience, some notations are first defined as follows:

- \mathcal{F} — The plaintext document collection of N files, denoted as $\mathcal{F} = \{F_1, F_2, \dots, F_N\}$. Each document is treated as a sequence of keywords. Note that, each file $F_i (1 \leq i \leq N)$ has a unique identifier $f_i (1 \leq i \leq N)$ in the whole document collection.

- \mathcal{A} — The attribute dictionary, denoted as $\mathcal{A} = \{A_1, A_2, \dots, A_n\}$. Each document and data user is associated with a set of attributes in \mathcal{A} .

- \mathcal{C} — The ciphertext of \mathcal{F} . In this paper, \mathcal{F} is symmetrically encrypted by content secret keys $ck = \{ck_1, ck_2, \dots, ck_N\}$, i.e., $C_i = E_{ck_i}(F_i), i = 1, 2, \dots, N$ and all the ciphertexts of the files compose \mathcal{C} .

- \mathcal{I} — The index structure of \mathcal{F} . Each document is first mapped to a document vector and the vectors are organized in an ARF tree.

- \mathcal{W} — The keyword dictionary, denoted as $\mathcal{W} = \{w_1, w_2, \dots, w_m\}$, which is used to generate the document vectors and query vectors.

- \mathcal{W}_Q — A subset of \mathcal{W} , representing the keywords in a query.

- \mathcal{Q} — The document query request of a data user. Each query contains multiple keywords \mathcal{W}_Q which are employed to describe the interested documents. In addition, the attributes of the data user are also added into \mathcal{Q} to check the legality of a document. We say that a document has legal attributes if the

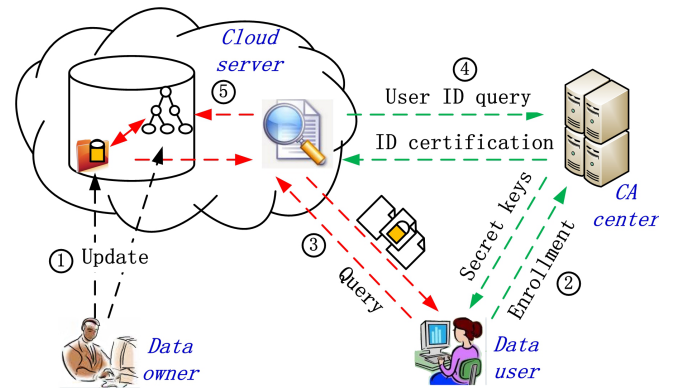


Fig. 1. System model

document's attribute set is a subset of the data user's attribute set and this will be discussed in Section IV.A.

- \mathcal{SR} — The result of a search request, i.e., a set of encrypted documents which are the top- k relevant documents to the request under the constraint of a data user's attributes.

A. System Model and Design Goals

In this paper, we attempt to design a fine-grained access control mechanism for the encrypted documents which also support efficient document search. The search result of a query is defined as the top- k relevant encrypted documents with legal attributes. The process of executing a document query is presented in Fig. 1 and it is mainly composed of five stages:

- ① The data owner is responsible for collecting and pre-processing the documents, and then obtains a set of high quality files \mathcal{F} . He sets the attributes for each document and then hierarchically encrypts the document collection based on attributes. In addition, an index vector is extracted from each document based on the document's content and attributes. An index structure \mathcal{I} is constructed based on the index vectors of the documents. At last, both the encrypted documents \mathcal{C} and encrypted index structure are sent to the cloud server. The cloud server is responsible for storing the encrypted documents and executing document search based on the index structure.

- ② When a data user wants to search a set of interested documents, she first needs to register herself as an authorized data user at the certificate authority (CA) center. Then, if possible, several attributes selected from \mathcal{A} are assigned to the data user by CA and a corresponding secret key associated with these attributes is sent to the data user. At last, the data user can send a query request \mathcal{Q} to the cloud server.

- ③④⑤ Once a query is received from a data user, the cloud server first communicates with the CA to check the legality of the data user and her attributes. If the data user is authorized, the cloud server searches the index structure to obtain the search result \mathcal{SR} . Then the corresponding encrypted documents are extracted from the encrypted document collection \mathcal{C} and sent to the data user. At last, the data user decrypts the documents by her secret key. Note that, the legality checking functionality is optional which can be employed to improve the security level of the whole system. With legality checking, the data users who didn't register themselves in the CA center

cannot search the interested documents through the cloud server. However, the security of the system doesn't greatly decrease without this functionality and it can be explained by the fact that the illegal data users cannot decrypt the documents returned by the cloud server because they don't have the secret keys.

In this paper, we assume that the CA center and the cloud server are trustable. Specifically, the CA center can distribute proper attributes to the data users and the cloud server can execute all the instructions honestly. We further assume that the data users are greedy and attempt to obtain as many plaintext files as possible. The data users try to collude with other users to decrypt the encrypted documents. We mainly restrict our attention to the process of encryption, document search and decryption, and the design goals of our scheme are presented as follows:

- *Flexibility*. The documents can be encrypted and decrypted flexibly based on their attributes. In general, we hope that the proposed scheme can get logarithmic encryption and decryption time efficiency.
- *Compactability*. For a data user with an attribute set, she needs to store only one secret key and the key can be used to decrypt all the documents that have legal attributes.
- *Accuracy*. The search results are accurate according to the data users' search request.
- *Efficiency*. Our scheme aims to achieve logarithmic search efficiency over the encrypted files in general and at least sub-linear search efficiency in the worst case.

B. Document/Query Vector

In this paper, the vector of a document is composed of two parts including a normalized content vector and an attribute vector. To build the content vector, each document is treated as a stream of keywords and we use the normalized term frequency (TF) vector to quantize the documents [40]. For each keyword w_i in keyword dictionary \mathcal{W} , we denote the number of times that this keyword appears in the document F_j by f_{j,w_i} and the TF value of keyword w_i in F_j is defined as $TF'_{j,w_i} = \ln(1 + f_{j,w_i})$. We construct the content vector of F_j as $(TF'_{j,w_1}, TF'_{j,w_2}, \dots, TF'_{j,w_m})$ and further normalize this vector by

$$TF_{j,w_i} = \frac{TF'_{j,w_i}}{\sqrt{\sum_{w_k \in \mathcal{W}} (TF'_{j,w_k})^2}}, \quad i = 1, 2, \dots, m \quad (1)$$

At last, the normalized content vector for F_j is denoted as $V_j = (TF_{j,w_1}, TF_{j,w_2}, \dots, TF_{j,w_m})$. The inverse document frequency (IDF) value of the keyword w_i is defined as $IDF_{w_i} = \ln(\frac{N}{N_{w_i}})$ where N is the number of documents in the whole collection and N_{w_i} is the number of documents that contain the keyword w_i . Further, the query vector of a query \mathcal{Q} is represented as $V_Q = (q_1, q_2, \dots, q_m)$ where q_i is 0, if $w_i \notin \mathcal{W}_Q$; and q_i is IDF_{w_i} , if $w_i \in \mathcal{W}_Q$.

The attribute vector of F_j is denoted as $V'_j = (V'_j[1], V'_j[2], \dots, V'_j[n])$ which is constructed based on the attribute dictionary $\mathcal{A} = \{A_1, A_2, \dots, A_n\}$ as follows: If $A_i \in att(F_j)$, $V'_j[i] = 1$, where the function $att(F_j)$ returns the attribute set associated with document F_j ; otherwise,

$V'_j[i] = 0$. The attribute vector of a data user V'_Q can be constructed based on the user's attributes in the similar manner. At last, we adopt the widely used "TF-IDF" measurement to calculate the relevance score between a document F_j and a query \mathcal{Q} as follows:

$$RScore(F_j, \mathcal{Q}) = RScore(V_j, V_Q) = V_j \cdot V_Q \quad (2)$$

It can be observed that the attribute vectors are not employed when calculating the relevance scores between a document and a query. This is reasonable considering that we need to return the legal documents of the query rather than the documents that have similar attributes with the query.

C. Attribute-based Retrieval Feature and ARF Tree

To improve the search efficiency of multi-keywords search process, a height-balanced index tree named ARF tree is built based on the document vectors. Specifically, the document vectors are organized as clusters according to their similarities. Each node in the tree represents a cluster composed of a set of document vectors or sub-clusters. An ARF vector is a quintuple summarization about a cluster. Given P documents $\{F_j\}$ where $j = 1, 2, \dots, P$, we assume that a cluster C comprises the document vectors of $\{F_j\}$, i.e., $\{V_j, V'_j\}$ where $j = 1, 2, \dots, P$. Then, the ARF vector of the cluster is defined as follows: $ARF = (P, LS, SS, V_{max}, A_{min})$, where P is the number of document content vectors in the cluster, LS is the linear sum of the P content vectors, i.e., $LS = \sum_{j=1}^P V_j$, SS is the square sum of the P content vectors, i.e., $SS = \sum_{j=1}^P V_j^2$, V_{max} denotes a vector consisting of m values which are calculated as follows:

$$V_{max}[i] = \max(V_1[i], V_2[i], \dots, V_P[i]), i = 1, 2, \dots, m \quad (3)$$

where $V_j[i]$ is the i -th dimensional value of V_j , A_{min} is the common attribute set vector of the documents in the cluster and it can be calculated as follows:

$$A_{min}[i] = V'_1[i] \wedge V'_2[i] \wedge \dots \wedge V'_P[i], i = 1, 2, \dots, n \quad (4)$$

where $V'_j[i]$ is the i -th dimensional value of V'_j . For each pair of bits in V'_i and V'_j , logic operation " \wedge " returns 1 if both the two bits are 1; otherwise, " \wedge " returns 0. For each pair of bits in V'_i and V'_j , logic operation " \vee " returns 1 if either of the two bits is 1; otherwise, " \vee " returns 0. As an example, $(1, 0, 0, 1) \wedge (1, 1, 0, 0) = (1, 0, 0, 0)$; $(1, 0, 0, 1) \vee (1, 1, 0, 0) = (1, 1, 0, 1)$.

In this paper, a search request of a data user contains both a set of keywords \mathcal{W}_Q and a set of attributes S_U associated with the data user. Only the documents, whose attributes are matched with S_U and contents are relevant with \mathcal{W}_Q , are returned to the data user. As a consequence, both the content vectors and the attribute vectors of the documents should be taken into consideration in document search process. The similarity between a pair of documents F_i, F_j with content vectors V_i, V_j and attribute vectors V'_i, V'_j is defined as follows:

$$Sim(F_i, F_j) = \gamma \cdot RScore(V_i, V_j) + (1 - \gamma) \cdot \frac{Length(V'_i \wedge V'_j)}{Length(V'_i \vee V'_j)} \quad (5)$$

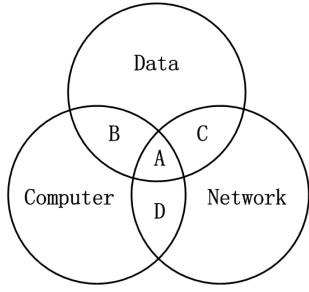


Fig. 2. Assumption of access control mechanism

where $0 \leq \gamma \leq 1$ and $RScore(V_i, V_j)$ is the relevance score between the content vectors of the two documents and it is calculated as:

$$RScore(V_i, V_j) = V_i \cdot V_j \quad (6)$$

γ is a preset parameter to adjust the importance degrees of document vectors and attribute vectors, $Length(V')$ returns the number of non-zero elements in vector V' . Based on an ARF vector, the centroid of a cluster C can be easily calculated as:

$$c = LS/P \quad (7)$$

and the similarity between cluster C and a document F_j is defined as:

$$Sim(C, F_j) = \gamma \cdot RScore(c, V_j) + (1 - \gamma) \cdot \frac{Length(A_{min} \wedge V_j')}{Length(A_{min} \vee V_j')} \quad (8)$$

where $0 \leq \gamma \leq 1$ and $RScore(c, V_j)$ is calculated as:

$$RScore(c, V_j) = c \cdot V_j \quad (9)$$

Further, the radius of cluster C is defined as follows:

$$R = \sqrt{\sum_{j=1}^P (V_j - c)^2 / P} \quad (10)$$

and it also can be calculated by the ARF vector as follows:

$$R = \sqrt{(SS - LS^2/P)/P} \quad (11)$$

Theorem 1 (ARF Additivity Theorem) : If we merge two disjoint clusters with ARF vectors: $ARF_1 = (P_1, LS_1, SS_1, V_{max1}, A_{min1})$, $ARF_2 = (P_2, LS_2, SS_2, V_{max2}, A_{min2})$, the ARF vector of the combined cluster is:

$$\begin{aligned} ARF &= ARF_1 + ARF_2 \\ &= (P_1 + P_2, LS_1 + LS_2, SS_1 + SS_2, V_{max}, A_{min}) \end{aligned}$$

where $V_{max}[i] = \max(V_{max1}[i], V_{max2}[i])$, $A_{min} = A_{min1} \wedge A_{min2}$.

Proof : The proof consists of straightforward algebra.

D. DBDH Assumption

Let $\mathbb{G}_0, \mathbb{G}_1$ be two groups of prime order p and g is a generator of \mathbb{G}_0 . The operator e is a bilinear map between \mathbb{G}_0 and \mathbb{G}_1 as specified in Section IV.B. The challenger chooses $a, b, c, t \in \mathbb{Z}_p$ at random. Then the challenger flips a fair binary coin v and if $v = 1$, it generates a BDH 5-tuple $(g, A = g^a, B = g^b, C = g^c, T = e(g, g)^{abc})$; otherwise, if $v = 0$, the

challenger outputs a random 5-tuple $(g, A = g^a, B = g^b, C = g^c, T = e(g, g)^t)$. The adversary must then output a guess v' of v .

An adversary, \mathcal{B} , has at least an ε advantage in solving the DBDH problem if

$$|Pr[\mathcal{B}(g, g^a, g^b, g^c, e(g, g)^{abc}) = 1] - Pr[\mathcal{B}(g, g^a, g^b, g^c, e(g, g)^t) = 1]| \geq 2\varepsilon$$

where the probability is over the randomly chosen a, b, c, t and the random bits consumed by \mathcal{B} . For the convenience of expression, we denote that $\mathcal{P}_{BDH} = \{(g, g^a, g^b, g^c, e(g, g)^{abc})\}$ and $\mathcal{R}_{BDH} = \{(g, g^a, g^b, g^c, e(g, g)^t)\}$.

Definition 1: The DBDH assumption holds if no probabilistic polynomial-time adversary has at least ε advantage in solving the above game.

E. Selective-Set Security Game

In this paper, we employ the Selective-Set Security Game [21], [28], [41] to prove our scheme's security. The game is composed of six phases and they are presented as follows.

Init. The adversary declares an access tree with a set of attributes S that he wants to be challenged upon.

Setup. The challenger runs the Setup algorithm presented in Section IV to generate the public parameters which are provided to the adversary.

Query Phase 1. The adversary is allowed to issue queries to obtain the secret keys of any access structure \mathbb{A}^* with attribute set S' , where $S \not\subseteq S'$. The secret keys are generated by the challenger through the $KeyGen(MSK, S')$ algorithm.

Challenge. The adversary provides two different messages M_0 and M_1 with equal length to the challenger. The challenger randomly flips a coin $\mu \in \{0, 1\}$ and encrypts M_μ with attribute set S . At last the encrypted message is sent to the adversary.

Query phase 2. The query phase 1 is repeated.

Guess. Based on the obtained information, the adversary output a guess μ' of μ .

We say that our scheme is secure if all the polynomial time adversaries have at most a negligible advantage in the game, where the advantage of the adversary is defined as $|Pr(\mu' = \mu) - \frac{1}{2}|$. Otherwise, we say that the adversary wins the game.

IV. HIERARCHICAL ATTRIBUTE-BASED DOCUMENT ENCRYPTION

A. Monotone Hierarchical Access Tree

Let $\mathcal{A} = \{A_1, A_2, \dots, A_n\}$ be a set of attributes. A collection $\mathbb{A} \subseteq 2^{\mathcal{A}}$ is monotone: Given $\forall B, C$, if $B \in \mathbb{A}$ and $B \subseteq C$, then $C \in \mathbb{A}$. A monotone access structure of a document is a monotone collection \mathbb{A} comprised of non-empty subsets of \mathcal{A} , i.e., $\mathbb{A} \subseteq 2^{\mathcal{A}} \setminus \{\emptyset\}$. The sets in \mathbb{A} are called authorized sets and the sets not in \mathbb{A} are called unauthorized sets. In this paper, we restrict our attention to monotone access structure which is practical considering the characteristics of the problem stated previously.

In this paper, we assume that a file associated with several attributes can be only accessed by the data users who possess all the basic attributes of the file. As an example shown in Fig. 2, the whole document set is divided into three categories

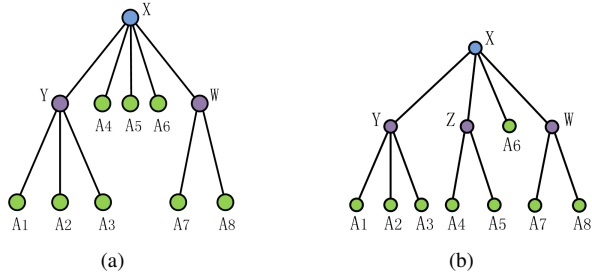


Fig. 3. Examples of access trees.

including computer, network and data related documents. Some documents may own two or three attributes such as the documents in region A, B, C and D . Under our assumption, the crossing region A can be accessed by the data users who own all the three roles of computer researcher, network researcher and data researcher; region B can be accessed by the data users who own the roles of data and computer researcher; region C can be accessed by the data users who own the roles of data and network researcher; region D can be accessed by the data users who own the roles of network and computer researcher. Apparently, under our assumption, the access structure of a document is monotone. Take region B as an example, a data user who owns the attribute of data and computer researcher can access B and then any other data users who have at least these two attributes can also access region B .

Let \mathcal{T} be a monotone hierarchical access tree representing an integrated access structure for a set of documents. The collection of all the access trees is called the access structure of the whole document collection. In this paper, each non-leaf node of the tree represents a threshold “AND” gate and associates with a set of attributes which are represented by the leaf nodes. For convenience, some functions are defined as follows. The number of the child nodes of a non-leaf node x is denoted as num_x . The function $att(x)$ denotes the associated attributes with the node x and in addition, $att(F_i)$ also returns the attribute set associated with document F_i . Each node in the tree is assigned with a numerical identifier and the function $index(x)$ returns the identifier of node x . In addition, $index(F_i)$ returns the identifier of F_i . Note that, each non-leaf node has a unique numerical identifier and the leaf nodes that represent the same attribute in different access trees share a same numerical identifier. Each node in an access tree may contain some files identifiers and the corresponding files will be encrypted by this node. The function $file(x)$ return the file identifiers contained in node x .

We say that node Y in the access tree \mathcal{T} matches a set of attributes S if and only if the attribute set of Y equals to S . As shown in Fig. 3(a), Y matches S if and only if $S = \{A_1, A_2, A_3\}$ and we denotes it as $\mathcal{T}_Y(S) = 0$. If there is no node in the tree can match S , we check whether a node in the tree can cover S . We say that node X covers S if X cannot match S and the leaf child nodes of X compose a superset of S . We denote $\mathcal{T}_X(S) = 1$ if node X covers S . As shown in Fig. 3(a), node Y covers S if $S = \{A_1, A_2\}$ and

Algorithm 1 BuildingAccessStructure.

Input: Document collection $\mathcal{F} = \{F_1, F_2, \dots, F_N\}$ with attribute sets $\{att(F_1), att(F_2), \dots, att(F_N)\}$
Output: A set of access trees $S_{\mathcal{T}}$

- 1: Sort the files in \mathcal{F} in descending order based on the number of their attributes and obtain $\mathcal{F}' = \{F'_1, F'_2, \dots, F'_N\}$;
- 2: $S_{\mathcal{T}} = null$;
- 3: **for** $i = 1 : N$ **do**
- 4: $S = att(F'_i)$;
- 5: Scan the access trees in order;
- 6: **for** the scanned access tree \mathcal{T} in $S_{\mathcal{T}}$ **do**
- 7: **if** node Y in \mathcal{T} matches S , i.e., $\mathcal{T}_Y(S) = 0$ **then**
- 8: Insert the identifier of F'_i into node Y ;
- 9: **break**;
- 10: **else if** node X in \mathcal{T} covers S , i.e., $\mathcal{T}_X(S) = 1$ **then**
- 11: Build a new node Z and let the created node Z be the child of X , and further the leaf nodes associated with S are inserted to Z ; meanwhile, the leaf nodes are deleted from X ;
- 12: Insert the identifier of F'_i into the new node Z ;
- 13: **break**;
- 14: **end if**
- 15: **end for**
- 16: **if** the identifier of F'_i has not been inserted into an access tree **then**
- 17: Build a new access tree for F'_i based on its attributes and insert the identifier of F'_i to the root node;
- 18: Insert the tree to $S_{\mathcal{T}}$;
- 19: **end if**
- 20: **end for**

node X covers S if $S = \{A_4, A_5, A_6\}$.

We construct the access structure of a document collection in an incremental way and an access tree is constructed by continuously splitting the tree in a top-down manner. In the initial, we sort the documents in decent order based on the number of their attributes. Apparently, the attribute set of the first document must be a root node of an access tree and the identifier of the document is inserted to the root node. Given a set of access trees, we discuss how to insert a new document F_i 's identifier into them. The attribute set of the new document $att(F_i)$ can be divided into three categories, i.e., being matched by a node in the access trees, being covered by a node in the access trees or neither being matched or covered by a node in the access trees. We first need to scan the access trees until finding a node that matches $att(F_i)$. If the node exists, the identifier of the new document $index(F_i)$ is inserted to the node. Otherwise, we need to rescan the access trees until find a node X that can cover $att(F_i)$. If the node exists, a new node Z is built in the tree to match $att(F_i)$ and insert $index(F_i)$ into Z . Specifically, node Z is inserted to the access tree as a child node of X and the leaf nodes related with $att(F_i)$ is inserted into node Z . Meanwhile, we need to delete the leaf nodes from node X . As an example, if we insert $\{A_4, A_5\}$ into the tree presented in Fig. 3(a), the updated access tree is shown in Fig. 3(b). At last, if $att(F_i)$ neither is matched or covered by a node in the trees, we build a new access tree for F_i and insert $index(F_i)$ into the root node. The above process is iterated until all the document identifiers are inserted into the access trees. All the access trees compose the access structure of the whole document collection.

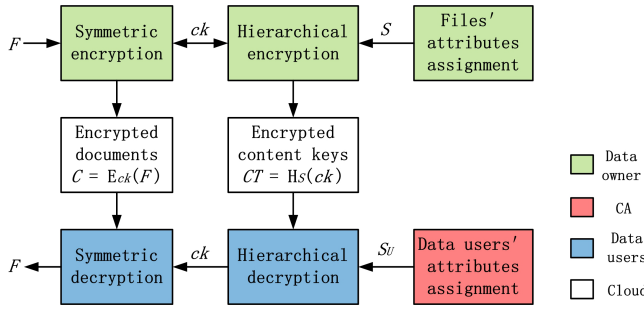


Fig. 4. The flow chart of document encryption and decryption

The pseudo-code of incrementally generating the hierarchical access trees for a document collection \mathcal{F} is presented in Algorithm 1. Based on Algorithm 1, a set of integrated access trees are constructed for the documents. Note that, all the nodes in an access tree compose a monotone access structure and each access tree contains several document identifiers. All the documents in a tree can be encrypted together, which will be discussed in Section IV.B. The identifier of the node x in an access tree is assigned as follows:

1. If x associated with attribute A_i is a leaf node, its numerical identifier is i .
2. If x is a non-leaf node and associated with a set of ordered attributes $\{A_i, A_j, \dots, A_k\} (i < j < \dots < k)$, its numerical identifier is $ij \dots k$.

In this way, each non-leaf node in the access structure has a unique identifier and apparently the leaf nodes associated with a same attribute share a same identifier.

B. Hierarchical Document Encryption

We first describe the system model of hierarchical attribute-based document encryption scheme as shown in Fig. 4. The data owner first selects a set of content keys $ck = \{ck_1, ck_2, \dots, ck_N\}$ which are used to encrypt the documents in \mathcal{F} symmetrically. Then, the content keys are hierarchically encrypted by the attributes assigned by the data owner. The encrypted documents, access structure and encrypted content keys are outsourced to the cloud server. In addition, the index structure of the document collection is also stored in the cloud server to support document search and it will be discussed in Section V. Once the encrypted search results are sent to the data users, they decrypt the content keys by their secret keys and further decrypt the documents based on the decrypted content keys. In the following, we mainly discuss how to encrypt the content keys in detail.

We first introduce the conceptions of bilinear map and Lagrange interpolation which are involved in our scheme. Let \mathbb{G}_0 and \mathbb{G}_1 be two multiplicative cyclic groups of prime order p . Let g be a generator of \mathbb{G}_0 and e be a bilinear map, $e : \mathbb{G}_0 \times \mathbb{G}_0 \rightarrow \mathbb{G}_1$ with the following properties:

1. Bilinearity: For all $u, v \in \mathbb{G}_0$ and $a, b \in \mathbb{Z}_p$, $e(u^a, v^b) = e(u, v)^{ab}$.
2. Non-degeneracy: $e(g, g) \neq 1$.
3. Distributivity: For $u, v, w \in \mathbb{G}_0$ and $a, b, c \in \mathbb{Z}_p$, $e(u^a, v^b w^c) = e(u^a, v^b) e(u^a, w^c)$.

In addition, \mathbb{G}_0 is a bilinear group if the group operations in \mathbb{G}_0 and the bilinear map $e : \mathbb{G}_0 \times \mathbb{G}_0 \rightarrow \mathbb{G}_1$ are both efficiently computable. The Lagrange Coefficient $\Delta_{i,S}$ for $i \in \mathbb{Z}_p$ and a set, S , of elements in \mathbb{Z}_p is defined as $\Delta_{i,S}(x) = \prod_{j \in S, j \neq i} \frac{x-j}{i-j}$. In addition, a hash function $H : \{0, 1\}^* \rightarrow \mathbb{G}_0$ is employed to map the string attributes to a random group element in \mathbb{G}_0 .

The detailed process of encrypting the documents is presented in the following:

Setup. Each document in \mathcal{F} is assigned with a set of attributes and the access structure of the document collection is constructed based on Algorithm 1. A set of content keys $ck = \{ck_1, ck_2, \dots, ck_N\}$ are randomly selected for the files in \mathcal{F} which are used to encrypt the files symmetrically. Then the setup algorithm chooses a bilinear group \mathbb{G}_0 with g as a generator, a bilinear map $e : \mathbb{G}_0 \times \mathbb{G}_0 \rightarrow \mathbb{G}_1$ and two random numbers $\alpha, \beta \in \mathbb{Z}_p$. The public key is published as:

$$PK = (\mathbb{G}_0, g, h = g^\beta, e(g, g)^\alpha)$$

and the master secret key MSK is (β, g^α) .

Encrypt(PK, ck, S_T). For each attribute A_i in \mathcal{A} , we first randomly select a unique secret number $s_i \in \mathbb{Z}_p$. $s_i \in \mathbb{Z}_p$. Then we choose a secret number sk_x for each node x in the access trees. In each access tree, these secret numbers for the nodes are chosen in a bottom-up manner, starting from the leaf nodes to the root node. The number sk_x of the leaf node x associated with attribute A_i is set as s_i . Then for the non-leaf node x with a set of child nodes S_x , the secret number sk_x is computed as $sk_x = \sum_{z \in S_x} sk_z \Delta_{i, S'_x}(index(x))$ where $i = index(z)$, $S'_x = \{index(z), z \in S_x\}$, $index(x)$ is the numerical identifier of node x . By iterating the above process, each node in the access structure can be assigned with a secret number.

Then, the content keys are encrypted by the secret numbers of the nodes in the access trees. As presented in Algorithm 1, each node x contains a set of file identifiers $\{f_m, \dots, f_n\}$ which can be returned by $file(x)$. We encrypt all the corresponding content keys $\{ck_m, \dots, ck_n\}$ by the same secret number sk_x . Specifically, for each access tree \mathcal{T} in S_T , let Y be the set of leaf nodes in \mathcal{T} . All content keys related with \mathcal{T} are encrypted together and the ciphertext is constructed as follows:

$CT = (\mathcal{T}, \forall x \in \mathcal{T}, f_i \in file(x) : C_x^* = g^{sk_x}, \tilde{C}_i = ck_i \cdot e(g, g)^{\alpha \cdot sk_x}, \forall y \in Y : C_y = h^{sk_y}, C'_y = H(att(y))^{sk_y})$. Note that, several leaf nodes y_1, y_2, \dots, y_d of different access trees $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_d$ may share a same attribute A_i and in this case, $C_{y_1} = C_{y_2} = \dots = C_{y_d} = h^{s_i}$, $C'_{y_1} = C'_{y_2} = \dots = C'_{y_d} = H(A_i)^{s_i}$. Therefore, in the ciphertext of the whole document collection, only $|\mathcal{A}|$ (i.e., the number of attributes) records of C_y and C'_y need to be stored.

KeyGen(MSK, S). The key generation algorithm takes a set of attributes S as input and output a secret key that identifies the set. We first chose a random $r \in \mathbb{Z}_p$, and then random $r_j \in \mathbb{Z}_p$ for each attribute $A_j \in S$. Then the keys are computed as follows:

$$SK = (D = g^\alpha \cdot h^r, \forall A_j \in S : D_j = g^r \cdot H(A_j)^{r_j}, D'_j = h^{r_j})$$

Decrypt(CT, SK). We employ a recursive algorithm $DecryptNode(CT, SK, x)$ to decrypt the content keys. This

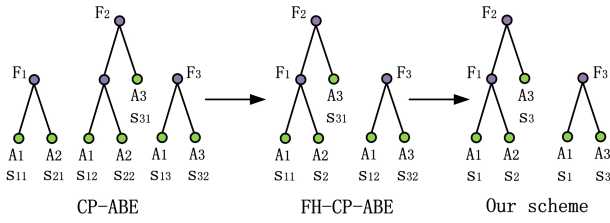


Fig. 5. Comparison of CP-ABE, FH-CP-ABE and our scheme

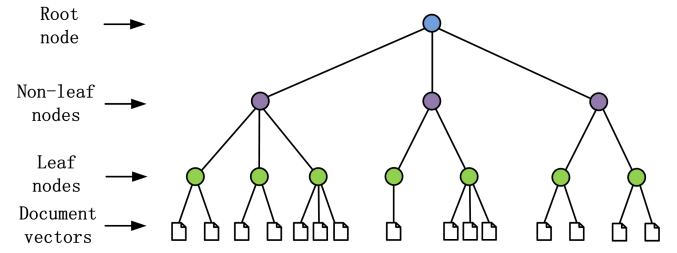


Fig. 6. An ARF tree

algorithm takes a ciphertext CT , a private key SK which is associated with a set of attributes S , and a node x from \mathcal{T} as input.

If the node x is a leaf node, let $A_i = att(x)$, and if $A_i \in S$, the algorithm is defined as follows:

$$\begin{aligned} DecryptNode(CT, SK, x) &= \frac{e(D_i, C_x)}{e(D'_i, C'_x)} \\ &= \frac{e(g^r \cdot H(A_i)^{r_i}, h^{sk_x})}{e(h^{r_i}, H(A_i)^{sk_x})} \\ &= \frac{e(g^r, h^{sk_x})e(H(A_i)^{r_i}, h^{sk_x})}{e(h^{r_i}, H(A_i)^{sk_x})} \\ &= e(g, g)^{r \cdot sk_x} \end{aligned}$$

If $A_i \notin S$, we define $DecryptNode(CT, SK, x) = \perp$.

When x is a non-leaf node, the algorithm is operated recursively. Specifically, it processes as follows: we denote the set of x 's children nodes by S_x . For each node $z \in S_x$, it calls $DecryptNode(CT, SK, z)$ and stores the output as F_z . If at least one $F_z = \perp$, the function $DecryptNode(CT, SK, x)$ returns \perp . Otherwise, we denote $i = index(z)$, $S'_x = \{index(z), z \in S_x\}$ and compute

$$\begin{aligned} F_x &= \prod_{z \in S_x} F_z^{\Delta_{i, S'_x}(index(x))} \\ &= \prod_{z \in S_x} (e(g, g)^{r \cdot sk_z})^{\Delta_{i, S'_x}(index(x))} \\ &= (e(g, g)^{r \cdot \sum_{z \in S_x} sk_z \cdot \Delta_{i, S'_x}(index(x))}) \\ &= e(g, g)^{r \cdot sk_x} \end{aligned}$$

If a data user with a set of attributes S that satisfies the subtree \mathcal{T}_x , the data user can calculate $A = F_x = e(g, g)^{r \cdot sk_x}$ and then each content key ck_i encrypted by node x with sk_x can be decrypted by computing $\tilde{C}_i / (e(C_x^*, D) / A) = \tilde{C}_i / (e(g^{sk_x}, g^\alpha \cdot h^r) / e(g, g)^{r \cdot sk_x}) = ck_i$. At last, all the documents encrypted by ck_i can be decrypted. Otherwise, the data user cannot decrypt the documents.

Note that, in the encryption phase, the secret numbers of the nodes are chosen in a bottom-up manner which is totally different from existing schemes such as CP-ABE and FH-CP-ABE. An advantage of this manner is that all the same attributes in different access trees share a same secret number and this can significantly improve the flexibility of encryption, decryption and secret keys distribution. As an example shown in Fig. 5, three files F_1, F_2, F_3 are associated with attribute sets $\{A_1, A_2\}$, $\{A_1, A_2, A_3\}$ and $\{A_1, A_3\}$, respectively. In CP-ABE, the three files are encrypted individually and attribute A_1 is related with three random secret number s_{11}, s_{12}, s_{13} , A_2 is related with s_{21}, s_{22} , A_3 is related with random secret

number s_{31}, s_{32} . In FH-CP-ABE, file F_1, F_2 share an access structure and they are encrypted together. File F_3 is encrypted individually. In this case, attribute A_1 is related with two secret number s_{11}, s_{12} . Attribute A_2 is related with s_2 and attribute A_3 is related with s_{31}, s_{32} . In our scheme, each attribute is related with only one secret number.

V. EFFICIENT RETRIEVAL OVER ENCRYPTED DOCUMENT COLLECTION

In this section, an efficient retrieval scheme over encrypted document collection is designed and we first describe the process of constructing the ARF tree. Then a depth-first searching algorithm of the ARF tree is designed and in addition, it can be operated in a parallel manner flexibly. Given a collection of documents $\mathcal{F} = \{F_1, F_2, \dots, F_N\}$, each document needs to be scanned for one time and the number of each keyword is recorded. Then a normalized vector for the document is generated based on the keyword dictionary \mathcal{W} as discussed in Section III.B. The attribute vector of a document can be built based on attribute dictionary \mathcal{A} and the associated attributes assigned by the data owner. Organizing the document vectors properly can significantly improve the search efficiency. In some encrypted document retrieval schemes [17], [18], the document content vectors are organized randomly and the search complexity is $O(N)$, where N is the number of documents. To improve search efficiency, in some other schemes [15], [16], the vectors are organized based on their relative similarities and they can obtain sub-linear search efficiency. However, the search accuracy cannot be guaranteed. In our scheme, the similarity between a pair of documents is calculated based on both the content vectors and attribute vectors. The proposed scheme can always obtain the accurate search results with at least a sub-linear search efficiency.

For convenience sake, we first describe the structure of an ARF tree briefly. An ARF tree is presented in Fig. 6 and it can be observed that the ARF tree is a height-balanced multi-way tree. An ARF tree has three main parameters including branching factors K_1, K_2 and threshold T which are preset by the data owner. A leaf node L_i contains at most K_1 document vectors and it is defined as follows:

$$L_i = (ARF, child_1, \dots, child_j), 1 \leq j \leq K_1$$

where ARF is the ARF vector of the cluster, $child_j$ is a pointer to the j -th document vector in the cluster. Each leaf node represents a micro cluster composed of a set of document

vectors. Each non-leaf node NL_i contains at most K_2 child nodes and it is defined as follows:

$$NL_i = (ARF, ARF_1, child_1, \dots, ARF_j, child_j), 1 \leq j \leq K_2$$

where ARF is the ARF vector of the whole cluster represented by NL_i , ARF_j is the ARF vector of the j -th sub-cluster and $child_j$ is a pointer to the child node representing the sub-cluster. Therefore, a non-leaf node represents a cluster made up of all the sub-clusters represented by its child nodes. Further, the cluster of a leaf node must satisfy a threshold requirement: the radius of the cluster which can be calculated by (11) has to be less than T .

We construct the ARF tree in an incremental manner which is similar to the construction process of the CF tree [42]. The process of inserting a document F_j with vector $\{V_j, V'_j\}$ into the ARF tree is presented as follows:

- *Identifying the appropriate leaf node:* Starting from the root, F_j recursively descends the ARF tree by choosing the most similar child node according to the similarity scores between F_j and the sub-clusters as defined in (8) until it reaches a leaf node.

- *Modifying the leaf node:* When F_j reaches a leaf node L_i , it tests whether L_i can “absorb” $\{V_j, V'_j\}$ without violating the constraints of K_1 and T . If so, $\{V_j, V'_j\}$ is inserted into L_i and the ARF vector of L_i is updated based on Theorem 1 as discussed in Section III.C. If not, we must split L_i to two leaf nodes. Node splitting is done by choosing the farthest pair of document vectors based on (5) as seeds, and then redistributing the remaining document vectors based on the closest criteria. The ARF vectors of the two new leaf nodes need to be recalculated.

- *Modifying the path from the root node to the leaf node:* After inserting $\{V_j, V'_j\}$ into a leaf node, we need to update the ARF vector for all the nodes on the path to the leaf node L_i . In the absence of a split, this simply involves updating ARF vectors based on Theorem 1. A leaf node split requires us to insert a new leaf node to the parent node. If the parent node has space for the new leaf node, we just need to insert the new leaf node into it and then update ARF vector for the parent node. In general, however, we may have to split the parent node as well, and so up to the root. If the root is split, the tree height increases by one.

Except for K_1 , K_2 , and T , the parameter γ can also affect the structure of the ARF tree. If γ is set to 1, the documents will be organized based on their content only and the associated attributes are ignored. On the contrary, if we set γ as 0, the attributes of the documents decide the ARF tree and the content of the documents are not employed. In general, we can set γ as a number between 0 and 1 to balance the important degrees of documents’ contents and attributes.

Another challenge is searching the top- k relevant documents whose attributes are covered by the data users. We design a depth-first search algorithm for the ARF tree and the pseudo-code is presented in Algorithm 2. For convenience, some symbols and functions are first defined as follows:

- *kthScore* - The smallest relevance score in current result list $RList$ which stores the most k relevant legal accessed

Algorithm 2 DepthFirstSearch.

Input: an ARF tree with root r , a query vector V_Q , an attribute vector V'_Q of the data user

Output: The most k relevant legal document vectors

```

1:  $u \leftarrow r$ ;
2: while  $u$  is not a leaf node do
3:   for all the child nodes  $v$  of node  $u$  do
4:     Calculate the relevance scores between  $v$  with  $V_Q$  by  $RScore(v, V_Q)$ ;
5:     Check whether the attribute set  $A_{v,min}$  is covered by  $V'_Q$  by comparing  $Length(A_{v,min} \wedge V'_Q)$  and  $Length(A_{v,min})$ ;
6:      $u \leftarrow$  the most relevant child node whose attributes are covered by  $V'_Q$ ;
7:   end for
8: end while
9: Select the most relevant  $k$  document vectors in the leaf node  $u$  whose attributes are covered by  $V'_Q$  and construct  $RList$ ;
10:  $Stack.push(r)$ ;
11: while  $Stack$  is not empty do
12:    $u \leftarrow Stack.pop()$ ;
13:   if the node  $u$  is not a leaf node then
14:     if  $RScore(V_{u,max}, V_Q) > kthScore$  and  $Length(A_{u,min} \wedge V'_Q) = Length(A_{u,min})$  then
15:       Sort the child nodes of  $u$  in ascent order based on the relevant scores with  $V_Q$  whose attribute sets are covered by  $V'_Q$ ;
16:       Push the children of  $u$  into  $Stack$  in order, i.e., the most relevant child is latest inserted into  $Stack$ ;
17:     end if
18:   else
19:     Calculate the relevance scores between the document vectors in the leaf node with  $V_Q$  and compare their attributes with  $V'_Q$ ;
20:     Update  $RList$ ;
21:   end if
22: end while
23: return  $RList$ 

```

document vectors with V_Q and the corresponding relevance scores in order.

- *RScore(u, V_Q)* - The relevance score between the cluster represented by node u and a query vector V_Q is defined as $RScore(u, V_Q) = c \cdot V_Q$ where c is the center of the cluster.

- *Stack* - We employ the variable $Stack$ to store the nodes which need to be searched in the future. In addition, $Stack.push(u)$ inserts node u into $Stack$ and $Stack.pop()$ returns the latest inserted node.

- *Length(V')* - This function returns the number of non-zero elements in attribute vector V' . For two attribute vectors V'_1 and V'_2 , we can test whether V'_1 is covered by V'_2 by checking whether $Length(V'_1 \wedge V'_2) = Length(V'_1)$. If $Length(V'_1 \wedge V'_2) = Length(V'_1)$, V'_1 is covered by V'_2 ; otherwise, V'_1 is not covered by V'_2 .

As shown in line 1 to line 9 in Algorithm 2, we first need to initialize $RList$ by finding the most similar leaf node. Then, as shown in line 10 to line 22, the paths in the tree needed to be searched are selected by criteria $RScore(V_{u,max}, V_Q) > kthScore$ and $Length(A_{u,min} \wedge V'_Q) = Length(A_{u,min})$. This is reasonable considering that if $RScore(V_{u,max}, V_Q) \leq kthScore$ for a cluster, it is impossible that any document vector in the cluster can be a candidate of the search result because the elements in V_Q and V_{max} are naturally nonnega-

tive. In addition, if $Length(A_{u,min} \wedge V'_Q) \neq Length(A_{u,min})$ for a cluster, all the attributes of the documents in the cluster cannot match that of the data user. As a consequence, this cluster is also unnecessary to be searched. However, if a leaf node is searched, the result list $RList$ needs to be updated. In this way, quite many paths are pruned and hence the search efficiency greatly improves. Once the top- k relevant documents are located in the ARF tree, the corresponding encrypted documents are sent to the data user. Apparently, these legal documents can be decrypted by the data user and then the document query process is completed.

We can further improve the search efficiency by operating the searching process in parallel. In the search process, all the processors need to share the same result list $RList$. Assume that there is a set of processors $\mathcal{P} = \{p_1, p_2, \dots, p_l\}$ and given a search request, an idle processor p_i is used to find most relevant leaf node on the tree and initialize $RList$. Then, all the necessary search paths are selected based on criteria $V_Q \cdot V_{max} > kthScore$ and $Length(A_{min} \wedge V'_Q) = Length(A_{min})$. If the search process can be continued on q search paths and there are more than q idle processors, any q processors are selected and each processor is responsible for searching a child path. If there is $q'(q' < q)$ idle processors, they search the latest inserted q' children paths in $Stack$. Once an idle processor appears, it continue to search the node generated by $Stack.pop()$. At last, the most relevant k encrypted files (i.e., the search result SR) are sent to the data user and they are decrypted by the secret key of the data user.

Though the document retrieval efficiency is greatly improved based on the ARF tree, a accompanying challenge is how to protect the privacy of the document vectors in the index structure and query vectors. Fortunately, this problem has been widely discussed and researched [15]–[17], [37]. In this paper, we strictly employ the method in [18] to protect the security of the document vectors while maintaining the searchability.

VI. SECURITY ANALYSIS

In the document retrieval system, the cloud server and CA center are assumed to be trustable. In this section, we focus on the security of the proposed hierarchical document encryption scheme and its security mainly involves two aspects including document confidentiality and content keys confidentiality. The documents are encrypted based on symmetric encryption schemes (e.g., AES) with content keys and their security is out of the scope in this paper. In this section, we analyze the security of the content keys which are encrypted by the proposed hierarchical encryption scheme. We provide the Decisional Bilinear Diffie-Hellman [28], [41], [43] assumption (DBDH) in Section III.D and Selective-Set Security Game is given Section III.E. In this section, we reduce the security of the content keys to the hardness of the DBDH and prove the security of the proposed scheme under the Selective-Set Security Game.

Theorem 2: Under the DBDH assumption, no polynomial time adversary can win the Selective-Set Security Game.

Proof: Suppose there exists an polynomial adversary Adv that can break our scheme with an advantage ε . We can design

a simulator \mathcal{B} that can play the DBDH game with an advantage $\frac{\varepsilon}{2}$. The game is executed as follows:

First, the challenger chooses $\mathbb{G}_0, \mathbb{G}_1, g, a, b, c, t$ and a bilinear map e as specified in Section IV.B. Then he randomly flips a fair binary coin v and if $v = 1, T = e(g, g)^{abc}$, i.e., $(g, A = g^a, B = g^b, C = g^c, T = e(g, g)^{abc}) \in \mathcal{P}_{BDH}$; otherwise, if $v = 0, T = e(g, g)^t$, i.e., $(g, A = g^a, B = g^b, C = g^c, T = e(g, g)^t) \in \mathcal{R}_{BDH}$. The challenger sends $(g, A, B, C, T) = (g, g^a, g^b, g^c, T)$ to the simulator \mathcal{B} . The simulator \mathcal{B} now plays the role of challenger in the security game. Then, the security game are executed as follows:

Init. The adversary Adv submits the simulator \mathcal{B} a set of attributes S that it wants to be challenged upon.

Setup. The simulator \mathcal{B} sets $\alpha = \alpha' + ab$ where α' is randomly selected from \mathbb{Z}_p and it computes $e(g, g)^\alpha = e(g, g)^{\alpha'} \cdot e(g, g)^{ab}$. It further sets $h = g^\beta = g^b = B$. At last, the public key PK is sent to the adversary Adv .

Query Phase 1. The adversary Adv queries \mathcal{B} the secret keys SK of any access structure \mathbb{A}^* with a set of attributes S' and $S \not\subseteq S'$. The simulator \mathcal{B} randomly selects a number $r' \in \mathbb{Z}_p$ and set $r = r' - a$. Then it calculates $D = g^{\alpha'} \cdot h^r = g^{\alpha' + ab} \cdot g^{b(r' - a)} = g^{\alpha'} \cdot h^{r'}$. For each attribute $A_j \in S'$, \mathcal{B} randomly chooses $r_j \in \mathbb{Z}_p$ and calculates $D_j = g^{(r' - a)}$. $H(A_j)^{r_j} = \frac{g^{r'}}{A} \cdot H(A_j)^{r_j}$ and $D'_j = B^{r_j}$. At last, the secret key SK is sent to the adversary Adv .

Challenge. For convenience sake, we assume that only one file is encrypted and consequently the ciphertext can be simplified as $CT = (\mathcal{T}, C_x^*, \tilde{C}_i, \forall y \in S' : C_y = B^{sk_y}, C'_y = H(att(y))^{sk_y})$. The adversary Adv submits two messages M_0 and M_1 with equal lengths to \mathcal{B} . The simulator \mathcal{B} randomly flips a coin $\mu \in \{0, 1\}$ and encrypts M_μ with attribute set S . Let $C_x^* = g^{sk_x} = g^c = C$. Suppose that the simulator is given a BDH tuple, that is $T = e(g, g)^{abc}$. Then we have $\tilde{C}_i = M_\mu \cdot e(g, g)^{\alpha c} = M_\mu \cdot e(g, g)^{abc} \cdot e(g, g)^{\alpha' c} = M_\mu \cdot T \cdot e(g, g)^{\alpha' c}$. We see that the ciphertext is a valid encryption of M_μ . Otherwise, we have that $T = e(g, g)^t$ is a random element of \mathbb{G}_1 . In that case the ciphertext will give no information about the simulator's choice of μ . At last, the CT is sent to Adv .

Query phase 2. The query phase 1 is repeated.

Guess. The adversary Adv makes a guess μ' of μ based on the obtained information. At the same time, the simulator \mathcal{B} also makes the corresponding guess of v in playing the DBDH game based on the different results the adversary Adv guessed. If $\mu' = \mu$, \mathcal{B} outputs guess $v' = 1$ in playing the DBDH game and points out that the challenger given 5-tuple to it which is selected from \mathcal{P}_{BDH} . If $\mu' \neq \mu$, \mathcal{B} outputs guess $v' = 0$ in playing the DBDH game and points out that the challenger given 5-tuple to it which is selected from \mathcal{R}_{BDH} .

The probability that the simulator \mathcal{B} successes in playing the DBDH game between simulator and challenger is calculated as follows.

If $v = 1$, the challenger generates a BDH tuple $(g, g^a, g^b, g^c, e(g, g)^{abc})$, i.e. $(g, A, B, C, T) \in \mathcal{P}_{BDH}$. Then we see that CT is a valid encryption of M_μ and by definition, in this case the adversary Adv has a non-negligible advantage ε to guess the correct μ' , whose probability of success can be

TABLE I. Comparison of CP-ABE, FH-CP-ABE and our scheme

Component	CP-ABE	FH-CP-ABE	Our Scheme
Encryption Time	$[2(\mathbb{A}_{C_1} + \dots + \mathbb{A}_{C_N}) + N]C_{G_0} + NC_{G_1} + NC_e$	<i>Null</i>	$(2 \mathcal{A} + \rho N)C_{G_0} + NC_{G_1} + NC_e$
Decryption Time	$2[(t_1 + \dots + t_N) + N]C_{G_1} + [2(\mathbb{A}_{C_1} + \dots + \mathbb{A}_{C_N}) + N]C_e$	<i>Null</i>	$2N(\rho + 1)C_{G_1} + (2 \mathcal{A} + N)C_e$
The Size of <i>PK</i>	$3L_{G_0} + L_{G_1}$	$3L_{G_0} + L_{G_1}$	$2L_{G_0} + L_{G_1}$
The Size of <i>MSK</i>	$L_{\mathbb{Z}_p} + L_{G_0}$	$L_{\mathbb{Z}_p} + L_{G_0}$	$L_{\mathbb{Z}_p} + L_{G_0}$
The Size of <i>SK</i>	$[2(\mathbb{A}_{C_1} + \dots + \mathbb{A}_{C_N}) + N]L_{G_0}$	$[2(\mathbb{A}_{C_1} + \dots + \mathbb{A}_{C_N}) + N]L_{G_0}$	$(2 \mathcal{A} + 1)L_{G_0}$
The Size of <i>CT</i>	$[2(\mathbb{A}_{C_1} + \dots + \mathbb{A}_{C_N}) + N]L_{G_0} + NL_{G_1}$	<i>Null</i>	$(2 \mathcal{A} + \rho N)L_{G_0} + NL_{G_1}$

calculated as $Pr[\mu' = \mu | (g, A, B, C, T) \in \mathcal{P}_{BDH}] = \frac{1}{2} + \varepsilon$.

If $v = 0$, the challenger builds a random 5-tuple $(g, g^a, g^b, g^c, e(g, g)^t)$, i.e. $(g, A, B, C, T) \in \mathcal{R}_{BDH}$. Then we have that T is a random element of \mathbb{G}_1 . The adversary Adv did not get any information about the message M_μ , so there is no advantage to guess the correct μ' . As a consequence, the adversary can make a correct choice with a probability $\frac{1}{2}$. Therefore, the probability of success for the simulator is $Pr[\mu' = \mu | (g, A, B, C, T) \in \mathcal{R}_{BDH}] = \frac{1}{2}$.

At last, the overall advantage of \mathcal{B} in playing the DBDH game can be calculated as $\frac{1}{2}Pr[\mu' = \mu | (g, A, B, C, T) \in \mathcal{P}_{BDH}] + \frac{1}{2}Pr[\mu' = \mu | (g, A, B, C, T) \in \mathcal{R}_{BDH}] - \frac{1}{2} = \frac{\varepsilon}{2}$. Based on the definition of DBDH assumption, we can infer that our scheme is secure. The theorem is proved.

VII. PERFORMANCE EVALUATION

In this section, we evaluate the performance of the hierarchical document encryption scheme and in addition the search efficiency of the ARF tree. We first analyze the efficiency theoretically and then conduct experiments to verify the analysis result.

A. Theoretical Analysis

We compare our hierarchical encryption scheme with CP-ABE scheme in [21] and FH-CP-ABE scheme [28], and some definitions are defined first. We assume that C_{G_i} ($i = 0, 1$) is the time cost of an operation on the group such as exponentiation or multiplication. Let \mathbb{Z}_p be the group $\{0, 1, \dots, p-1\}$ and C_e be the cost of a bilinear map operation e . Let N be the number of documents in the collection, ρ be a parameter that associated with Algorithm 1 and ρN is the number of the nodes in all the access trees. Considering that a set of file identifiers share a node in the access trees, ρ is naturally smaller than 1. Let $\mathcal{A}, \mathbb{A}_u, \mathbb{A}_{C_i}$ be the attribute dictionary, the attributes associated with the data user and document F_i respectively. Let t_i be the number of interior nodes in the access tree of file F_i . In addition, we define $|*|$ as the number of elements in $*$, L_* as the length of an element in $*$.

In the analysis, we assume that the data owner has N document files and their content keys are encrypted by CP-ABE, FH-CP-ABE and our scheme. Note that, we focus on the encryption process of the content keys rather than that of the documents which are encrypted by the content keys symmetrically. We further assume that a data user is responsible for decrypting all the documents and the analysis result is presented in Table 1. For a large document set, we have $|\mathcal{A}| \ll (|\mathbb{A}_{C_1}| + \dots + |\mathbb{A}_{C_N}|)$ and $\rho N < N \ll |t_1| + \dots + |t_N|$.

As a consequence, our scheme performs better than CP-ABE in time costs of encryption and decryption, and the sizes of *PK*, *SK* and *CT*. The two schemes have same performance in the size of *MSK*. In conclusion, our scheme can improve time and storage efficiency compared with CP-ABE. For a constant attribute set \mathcal{A} and parameter ρ , the encryption time, decryption time and size of *CT* all increase linearly with the number of documents in our scheme. The sizes of the keys are independent of the document collection. In addition, our scheme outperforms FH-CP-ABE in terms of the size of *PK* and *SK* and they have similar performance in terms of *MSK*. Considering that FH-CP-ABE is designed to encrypt a set of documents with incremental attribute sets, i.e., $\mathbb{A}_{C_1} \supseteq \mathbb{A}_{C_2} \supseteq \dots \supseteq \mathbb{A}_{C_N}$, it is impossible to accurately predict the time cost of encryption and decryption and the size of *CT* for a document collection with randomly assigned attribute sets. As a consequence, we will further compare our scheme with FH-CP-ABE by simulation in Section VII.B.

The organization structure of the document collection affects the search efficiency significantly. The keyword balanced binary (KBB) tree [18] can provide accurate search result. However, the document vectors are randomly inserted into the tree and they are organized chaotically. Some similar vectors may locate very far in the tree and some totally different vectors may be neighbors with each other. Consequently, the interior nodes in the tree can provide very limited information to lead a query vector to the area with a set of strongly relevant document vectors. On the contrary, the vectors in the ARF tree are organized strictly according to their similarities and similar vector can always compose a cluster in spite of the vectors' input order. The query vector can easily locate a cluster that contains relevant document vectors. The search proportion is defined as the proportion that the document vectors being searched in a search process and it is calculated by the number of the searched nodes to the number of all the nodes in the tree. A basic comparison between the two trees is presented in Fig. 7. All the document vectors are randomly generated in 2D and 3D space. To be fair, we ignore the attributes of data user and documents considering that the KBB tree does not support attribute constrained search. It can be observed that the ARF tree outperforms KBB tree significantly in both 2D and 3D spaces. Specifically, the search proportion of ARF tree is about 5% to 10% to that of KBB tree.

B. Experimental Simulation

We conduct a thorough experimental evaluation for the proposed document retrieval scheme on a real world data set: the

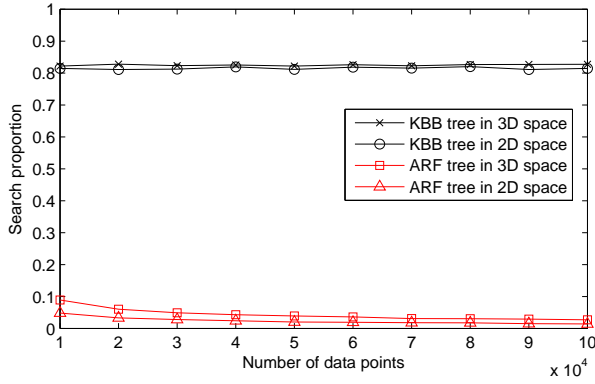


Fig. 7. Search proportion of KBB tree and ARF tree

Algorithm 3 AttributeGeneration.**Input:** $\mathcal{A} = \{C_1, C_2, C_3, C_4\}, \mathcal{F}, p_r (0.25 \leq p_r \leq 1)$ **Output:** The attributes of each document

```

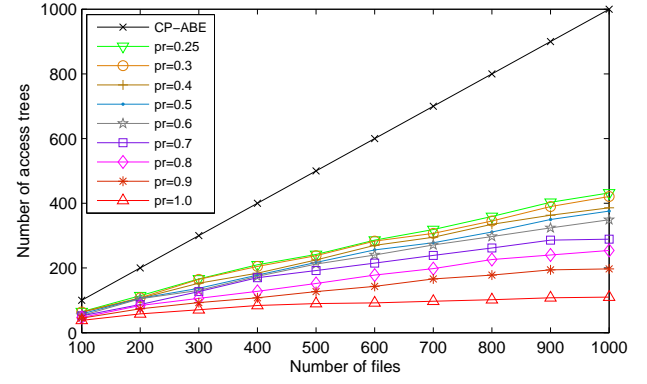
1: for each document  $F_i \in \mathcal{F}$  do
2:    $Att = \emptyset$ ;
3:   Randomly select a number  $m$  from  $\{1, 2, 3, 4, 5\}$ ;
4:   Randomly select an attribute  $A_n$  from  $\mathcal{A}$  and we assume that
      $A_n \in C_k, k = 1, 2, 3, 4$ ;
5:   Insert  $A_n$  to  $Att$ ;
6:   for  $i = 2 : m$  do
7:     Randomly generate a number  $p'_r (0 \leq p'_r \leq 1)$  and if  $p'_r \leq p_r$ ,
       randomly select an attribute  $A_q$  from  $C_k$ ; otherwise,
       uniformly randomly select an attribute  $A_q$  from  $\mathcal{A}$ ;
8:     Insert  $A_q$  to  $Att$ ;
9:   end for
10:  The attributes in  $Att$  is defined as the attributes of document  $F_i$ ;
11: end for

```

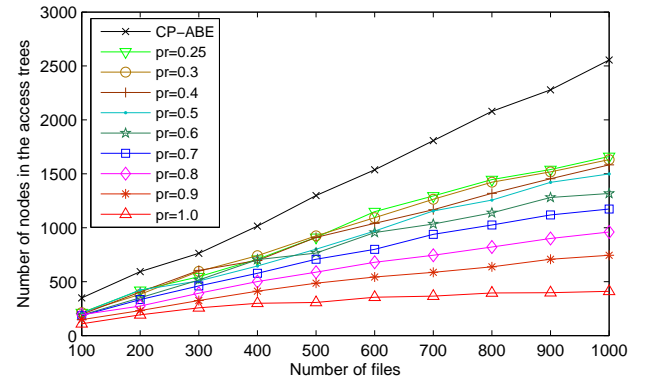
Enron Email Data Set [44]. The data set is first processed and 1,000 records are randomly chosen as our experiment corpus. We implement the hierarchical encryption scheme based on CP-ABE toolkit and Java Pairing-Based Cryptography library [45]. The document search process is implemented based on Java language. All the following experiments are conducted on a 2.6 GHZ Intel Core i5 processor, Windows 7 operating system with a RAM of 4G.

1) *Effectiveness of the Integrated Access Trees:* The attribute set is defined as $\mathcal{A} = \{A, B, \dots, Z\}$ which is composed of 26 letters. Then, all the attributes are divided into 4 categories, i.e., $C_1 = \{A, B, \dots, G\}, C_2 = \{H, I, \dots, N\}, C_3 = \{O, P, \dots, T\}, C_4 = \{U, V, \dots, Z\}$. The associated attributes of a document is randomly generated through Algorithm 3. We assume that each document has at least 1 attribute and at most 5 attributes. As shown in line 5 of Algorithm 3, the attributes of a document tend to belong to one attribute category with a large probability p_r . This is reasonable considering that the attributes are associated with each other and if a set of attributes are strongly related, they are likely to belong to a document together. For example, if a document is related with “computer”, it is more likely to be also related with “network” rather than other attributes such as “economic” and “finance”.

Parameter p_r affects the access trees greatly as presented



(a)



(b)

Fig. 8. Number of access trees and that of nodes in the trees with different p_r and number of files

in Fig.8. For a constant p_r , the number of the access trees monotonously increases with the number of files as shown in Fig.8(a). When p_r is set to 1.0, all the attributes of a file fall in a sub-category of \mathcal{A} and in this case the number of access trees is the smallest. Note that, a small number of access trees can lead a high encryption and decryption efficiency, because many documents share an access tree and they can be encrypted together in this case. When we decrease p_r from 1.0 to 0.3, the attributes of a file are more and more likely to be selected from the whole attribute set \mathcal{A} randomly and the diversity of the documents' attributes increases. Consequently, the number of the access trees increases. In the worst case, i.e., p_r is set to 0.25 and the attributes of a file are totally randomly selected from \mathcal{A} , the number of access trees is the largest with a constant number of files. In CP-ABE, each document has an access tree and the number of all the access trees equals to the number of files which is much larger than that of the proposed scheme. As shown in Fig.8(b), the number of nodes in the access trees has similar pattern with the number of access trees and the proposed scheme always performs better than CP-ABE.

We further analyze the distribution of files in the access trees and simulation result with $N = 1,000$ is provided in Fig. 9. The access trees are first descendingly sorted according to their sizes, i.e., the number of nodes in the trees, and then the numbers of files in the trees are calculated. It can be observed

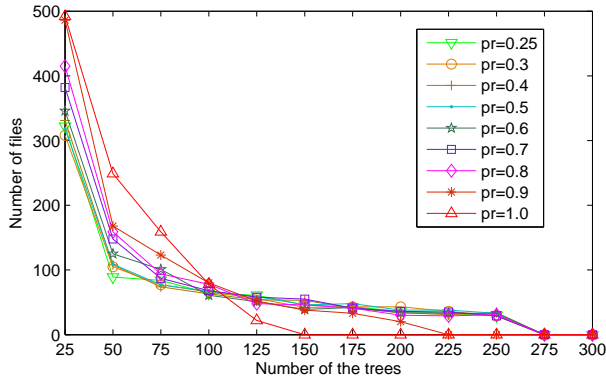


Fig. 9. Distribution of files in the access trees

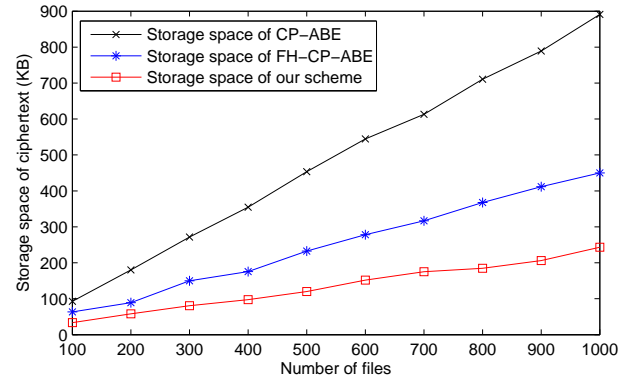


Fig. 11. Storage space of the ciphertext CT

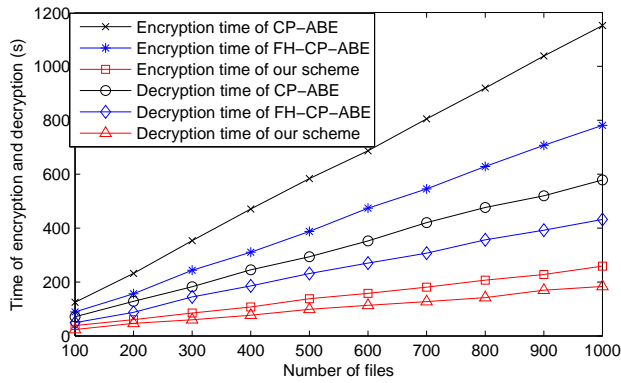


Fig. 10. Efficiency of encryption and decryption

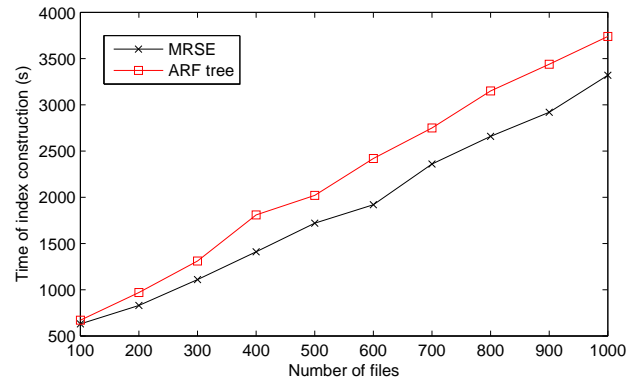


Fig. 12. Construction time of an ARF tree

that about 30% to 50% files are covered by the 25 largest trees and about 40% to 80% files are covered by the 50 largest trees. In addition, the files trend to aggregate with each other to some larger trees with the increasing of p_r . When we set p_r to 1, more than 90% files are covered by the largest 100 access trees and, most of the other trees contain a small number of nodes and they may cover 1 or 2 files. Without loss of generality, in the following, we assume that p_r equals to 0.9.

2) *Efficiency of Hierarchical Document Encryption*: The time consumptions of encrypting and decrypting the whole document collection are presented in Fig.10. In CP-ABE, each document is encrypted and decrypted individually. Consequently, the time of both encryption and decryption increases almost linearly with the number of files. On the contrary, a set of files in our scheme share an access tree and they are encrypted and decrypted together. The encryption and decryption time increases logarithmically with the number of files. Apparently, the proposed scheme is much more time efficient than CP-ABE. Though the FH-CP-ABE performs slightly better than CP-ABE, it cannot efficiently encrypt and decrypt a document collection considering that the number of integrated access trees are much larger than that of our scheme.

The storage space of the ciphertext is presented in Fig.11. Note that, only the encrypted content keys are considered in this experiment and the symmetrically encrypted documents are not considered. The storage space of CP-ABE linearly increases with the number of files and it can be explained

by the fact that each file has a content secret key which is encrypted individually. In our scheme, if a set of file have similar attribute sets, they may share an access structure and their content keys are related with each other. In addition, a set of files can share a same content key if they have the same attribute sets. Consequently, the proposed scheme is more space-efficient than CP-ABE. Similar to the efficiency of encryption and decryption, FH-CP-ABE performs better than CP-ABE and worse than our scheme.

3) *Efficiency of Document Retrieval*: Except for providing an efficient document encryption scheme, we also improve the search efficiency compared with MRSE. Note that, in our simulation, the index structures of both MRSE and ARF are plaintext. The construction time of an ARF tree is strongly related with the number of files and it is presented in Fig.12. The index construction times of both the two schemes linearly increase with the number of files. This can be explained by the fact that most time is consumed in the process of generating document vectors (about 3.2 seconds/file). The ARF tree consumes slightly more time than MRSE, because the document vectors need to be inserted into the tree.

Another measurement of our scheme is the search efficiency. In the Enron Email Data Set, the documents have no attribute which should be assigned by the data owner. In general, the attributes of the documents are related with their contents. However, in Algorithm 3, the attributes of a document are randomly selected and they may mislead the ARF tree

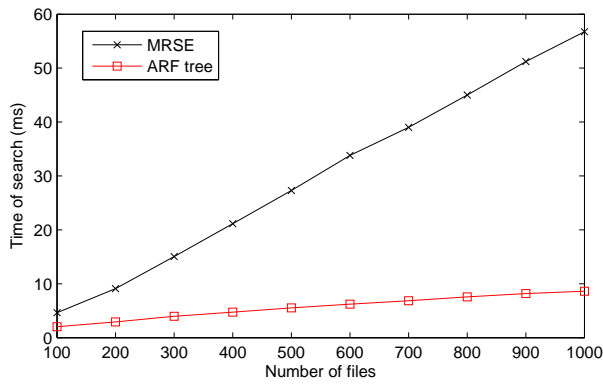


Fig. 13. Search time of a query

construction process. Consequently, for convenience, in the following we set γ equals to 1 when constructing the ARF tree. In addition, k is set as 10 (i.e., 10 encrypted documents are returned for a query). However, the attributes are employed in the document search process and the simulation result is provided in Fig.13. Apparently, the search time in MRSE linearly increases with the number of files considering that the document vectors are organized randomly and all the document vectors need to be scanned for one time. However, the ARF tree organizes the files based their similarities which greatly improve the search efficiency. Specifically, quite a number of the search paths are pruned in the search process and ARF tree has logarithmic time consumption with the number of files.

VIII. CONCLUSION

In this paper, we consider a new encrypted document retrieval scenario in which the data owner wants to control the documents in fine-grained level. To support this service, we first design a novel hierarchical attribute-based document encryption scheme to encrypt a set of documents together that share an integrated access structure. Further, the ARF tree is proposed to organize the document vectors based on their similarities. At last, a depth-first search algorithm is designed to improve the search efficiency for the data users which is extremely important for large document collections. The performance of the approach is thoroughly evaluated by both theoretical analysis and experiments.

The proposed scheme can be further improved in several aspects: First, in this paper, we assume that each node in the access trees represent an “AND” gate and this limits the flexibility of assigning the attributes to the documents. In the future, we will attempt to introduce “OR” gates into the access trees. Second, the access structure of the document collection is generated in a greedy manner and we will check whether it can be further optimized to decrease the number of access trees. In addition, the revocation method of the data users’ attributes needs to be designed. Third, the update strategy of the ARF tree should be proposed. Though the ARF tree naturally supports inserting new nodes to the tree, the method of deleting a node from the tree didn’t provided. Fourth, a new document collection, in which each file is associated with a

set of proper attributes, should be developed and a thorough experiment should be conducted on the collection to test the affection of parameter γ on the approach.

REFERENCES

- [1] K. Ren, C. Wang, and Q. Wang, “Security challenges for the public cloud,” *IEEE Internet Computing*, vol. 16, pp. 69–73, Jan. 2012.
- [2] D. X. Song, D. Wagner, and A. Perrig, “Practical techniques for searches on encrypted data,” in *Security and Privacy, 2000. SandP 2000. Proceedings. 2000 IEEE Symposium on*, pp. 0–44, 2002.
- [3] E. J. Goh, “Secure indexes,” *Cryptology ePrint Archive*, <http://eprint.iacr.org/2003/216>, 2003.
- [4] R. Curtmola, J. Garay, S. Kamara, and R. Ostrovsky, “Searchable symmetric encryption: improved definitions and efficient constructions,” in *ACM Conference on Computer and Communications Security*, pp. 79–88, 2006.
- [5] J. Li, Y. Shi, and Y. Zhang, “Searchable ciphertext-policy attribute-based encryption with revocation in cloud storage,” *International Journal of Communication Systems*, vol. 30, no. 1, 2017.
- [6] Y. Miao, J. Ma, X. Liu, X. Li, Q. Jiang, and J. Zhang, “Attribute-based keyword search over hierarchical data in cloud computing,” *IEEE Transactions on Services Computing*, vol. PP, no. 99, pp. 1–1, 2017.
- [7] 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 *ACM Workshop on Storage Security and Survivability, Storagess 2007, Alexandria, Va, Usa, October*, pp. 7–12, 2007.
- [8] 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, pp. 1467–1479, Aug. 2012.
- [9] S. Zerr, D. Olmedilla, W. Nejdl, and W. Siberski, “Zerber +r : top-k retrieval from a confidential index,” in *International Conference on Extending Database Technology: Advances in Database Technology*, pp. 439–449, 2009.
- [10] P. Golle, J. Staddon, and B. Waters, “Secure conjunctive keyword search over encrypted data,” *Lecture Notes in Computer Science*, vol. 3089, pp. 31–45, 2004.
- [11] B. Dan and B. Waters, “Conjunctive, subset, and range queries on encrypted data,” in *Theory of Cryptography Conference*, pp. 535–554, 2007.
- [12] A. Lewko, T. Okamoto, A. Sahai, K. Takashima, and B. Waters, “Fully secure functional encryption: attribute-based encryption and (hierarchical) inner product encryption,” in *International Conference on Theory and Applications of Cryptographic Techniques*, pp. 62–91, 2010.
- [13] Y. Miao, J. Ma, X. Liu, X. Li, Z. Liu, and H. Li, “Practical attribute-based multi-keyword search scheme in mobile crowdsourcing,” *IEEE Internet of Things Journal*, vol. PP, no. 99, pp. 1–1, 2017.
- [14] Y. Miao, J. Ma, X. Liu, Q. Jiang, J. Zhang, L. Shen, and Z. Liu, “Vcksm: Verifiable conjunctive keyword search over mobile e-health cloud in shared multi-owner settings,” *Pervasive and Mobile Computing*, vol. 40, pp. 205–219, 2017.
- [15] C. Chen, X. Zhu, P. Shen, J. Hu, S. Guo, Z. Tari, and A. Zomaya, “An efficient privacy-preserving ranked keyword search method,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 27, pp. 951–963, Apr. 2016.
- [16] Z. Fu, K. Ren, J. Shu, X. Sun, and F. Huang, “Enabling personalized search over encrypted outsourced data with efficiency improvement,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 27, pp. 2546–2559, Sep. 2016.
- [17] 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, pp. 222–233, Jan. 2014.
- [18] Z. Xia, X. Wang, X. Sun, and Q. Wang, “A secure and dynamic multi-keyword ranked search scheme over encrypted cloud data,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 27, pp. 340–352, Jan. 2016.
- [19] A. Sahai and B. Waters, “Fuzzy identity-based encryption,” in *International Conference on Theory and Applications of Cryptographic Techniques*, pp. 457–473, 2005.
- [20] J. Hur and K. N. Dong, “Attribute-based access control with efficient revocation in data outsourcing systems,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 22, no. 7, pp. 1214–1221, 2010.
- [21] J. Bethencourt, A. Sahai, and B. Waters, “Ciphertext-policy attribute-based encryption,” in *IEEE Symposium on Security and Privacy*, pp. 321–334, 2007.

- [22] V. Goyal, A. Jain, O. Pandey, and A. Sahai, *Bounded Ciphertext Policy Attribute Based Encryption*. DBLP, 2008.
- [23] F. Guo, Y. Mu, W. Susilo, and D. S. Wong, "Cp-abe with constant-size keys for lightweight devices," *Information Forensics and Security IEEE Transactions on*, vol. 9, pp. 763–771, May. 2014.
- [24] Y. Yang, J. K. Liu, K. Liang, K. K. R. Choo, and J. Zhou, *Extended Proxy-Assisted Approach: Achieving Revocable Fine-Grained Encryption of Cloud Data*, vol. 9327. Springer International Publishing, Sep. 2015.
- [25] K. Liang, H. A. Man, J. K. Liu, W. Susilo, D. S. Wong, G. Yang, Y. Yu, and A. Yang, "A secure and efficient ciphertext-policy attribute-based proxy re-encryption for cloud data sharing," *Future Generation Computer Systems*, vol. 52, pp. 95–108, Nov. 2015.
- [26] Y. S. Rao, "A secure and efficient ciphertext-policy attribute-based signcryption for personal health records sharing in cloud computing," *Future Generation Computer Systems*, vol. 67, pp. 133–151, Feb. 2017.
- [27] J. Li, W. Yao, Y. Zhang, H. Qian, and J. Han, "Flexible and fine-grained attribute-based data storage in cloud computing," *IEEE Transactions on Services Computing*, vol. 10, no. 5, pp. 785–796, 2017.
- [28] S. Wang, J. Zhou, J. K. Liu, J. Yu, J. Chen, and W. Xie, "An efficient file hierarchy attribute-based encryption scheme in cloud computing," *IEEE Transactions on Information Forensics and Security*, vol. 11, no. 6, pp. 1265–1277, 2016.
- [29] G. Wang, Q. Liu, and J. Wu, "Hierarchical attribute-based encryption for fine-grained access control in cloud storage services," in *ACM Conference on Computer and Communications Security*, pp. 735–737, 2010.
- [30] Z. Wan, J. Liu, and R. H. Deng, "Hasbe: A hierarchical attribute-based solution for flexible and scalable access control in cloud computing," *IEEE Transactions on Information Forensics and Security*, vol. 7, pp. 743–754, Apr. 2012.
- [31] H. Deng, Q. Wu, B. Qin, J. Domingo-Ferrer, L. Zhang, J. Liu, and W. Shi, "Ciphertext-policy hierarchical attribute-based encryption with short ciphertexts," *Information Sciences*, vol. 275, pp. 370–384, Aug. 2014.
- [32] E. Luo, Q. Liu, and G. Wang, "Hierarchical multi-authority and attribute-based encryption friend discovery scheme in mobile social networks," *IEEE Communications Letters*, vol. 20, pp. 1772–1775, Sep. 2016.
- [33] J. Li, Y. Wang, Y. Zhang, and J. Han, "Full verifiability for outsourced decryption in attribute based encryption," *IEEE Transactions on Services Computing*, 2017.
- [34] H. Qian, J. Li, Y. Zhang, and J. Han, "Privacy-preserving personal health record using multi-authority attribute-based encryption with revocation," *International Journal of Information Security*, vol. 14, no. 6, pp. 487–497, 2015.
- [35] J. Li, W. Yao, J. Han, Y. Zhang, and J. Shen, "User collusion avoidance cp-abe with efficient attribute revocation for cloud storage," *IEEE Systems Journal*, 2017.
- [36] Y. Guo, J. Li, Y. Zhang, and J. Shen, "Hierarchical attribute-based encryption with continuous auxiliary inputs leakage," *Security and Communication Networks*, vol. 9, no. 18, 2016.
- [37] W. K. Wong, D. W. Cheung, B. Kao, and N. Mamoulis, "Secure knn computation on encrypted databases," in *ACM SIGMOD International Conference on Management of Data*, pp. 139–152, 2009.
- [38] G. A. Miller, "Wordnet: a lexical database for english," *Communications of the Acm*, vol. 38, pp. 39–41, Nov. 1995.
- [39] J. Li, X. Lin, Y. Zhang, and J. Han, "Ksf-oabe: Outsourced attribute-based encryption with keyword search function for cloud storage," *IEEE Transactions on Services Computing*, vol. 10, no. 5, pp. 715–725, 2017.
- [40] C. D. Manning and P. Raghavan, *Introduction to Information Retrieval*, vol. 1. Cambridge University Press, 2010.
- [41] V. Goyal, O. Pandey, A. Sahai, and B. Waters, "Attribute-based encryption for fine-grained access control of encrypted data," in *ACM Conference on Computer and Communications Security*, pp. 89–98, 2006.
- [42] T. Zhang, R. Ramakrishnan, and M. Livny, "birch: an efficient data clustering method for very large databases," pp. 103–114, in *Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data*, 1996.
- [43] B. Waters, "Ciphertext-policy attribute-based encryption: An expressive, efficient, and provably secure realization," *Lecture Notes in Computer Science*, vol. 2008, pp. 321–334, 2015.
- [44] C. Project, "Enron email dataset," <https://www.cs.cmu.edu/enron/>, 2005.
- [45] A. D. Caro and V. Iovino, "jpbcc: Java pairing based cryptography," in *Computers and Communications*, pp. 850–855, Jun. 2011.



Na Wang received the B.S. and M.S. degrees in Mathematics from Xinjiang Normal University, China, in 2012 and 2015 respectively. She is currently working toward the Ph.D. degree in the School of Mathematical Sciences, Xiamen University, China. Her research interests include cryptography, message sharing and information security issues in distributed and cloud systems.



Junsong Fu received the B.E. degree from Beijing Jiaotong University, Beijing, China, in 2012. He is currently working toward the Ph.D. degree in the Key Laboratory of Communication and Information Systems, Beijing Jiaotong University. His research interests include in-network data processing, secret sharing and information privacy issues in distributed systems and Internet of Things.



IEEE. He is a fellow of IEEE.

Bharat Bhargava is a Professor of Computer Science at Purdue University. He is the editor-in-chief of four journals and serves on over ten editorial boards of international journals. Prof. Bhargava is the founder of the IEEE Symposium on Reliable and Distributed Systems, IEEE conference on Digital Library, and the ACM Conference on Information and Knowledge Management. Prof. Bhargava has published hundreds of research papers and has won five best paper awards in addition to the technical achievement award and golden core award from



security and cloud computing.

Jiwen Zeng is a Professor of School of Mathematic Science, Xiamen University. He obtained Ph.D degree from Beijing University in 1995. As academic visiting professor, he once visited Birmingham University of Unite Kingdom, York University of Canada, Jena University of Germany and Pingdong University in Taiwang, China. His research work is supported by natural science foundation of China. He is evaluation expert of state natural science foundation in Mathematics and its application. His research interests include mathematics, information